

# Formal Security Analysis of the OpenID FAPI 2.0: Accompanying a Standardization Process

Pedram Hosseyni, Ralf Küsters, and Tim Würtele  
Institute of Information Security  
University of Stuttgart, Germany  
{pedram.hosseyni, ralf.kuesters, tim.wuertele}@sec.uni-stuttgart.de

## Abstract

In recent years, the number of third-party services that can access highly-sensitive data has increased steadily, e.g., in the financial sector, in eGovernment applications, or in high-assurance identity services. Protocols that enable this access must provide strong security guarantees.

A prominent and widely employed protocol for this purpose is the OpenID Foundation’s FAPI protocol. The FAPI protocol is already in widespread use, e.g., as part of the UK’s Open Banking standards and Brazil’s Open Banking Initiative as well as outside of the financial sector, for instance, as part of the Australian government’s Consumer Data Rights standards.

Based on lessons learned from FAPI 1.0, the OpenID Foundation has developed a completely new protocol, called FAPI 2.0. The specifications of FAPI 2.0 include a concrete set of security goals and attacker models under which the protocol aims to be secure.

Following an invitation from the OpenID Foundation’s FAPI Working Group (FAPI WG), we have accompanied the standardization process of the FAPI 2.0 protocol by an in-depth formal security analysis. In this paper, we report on our analysis and findings.

Our analysis incorporates the first formal model of the FAPI 2.0 protocol and is based on a detailed model of the web infrastructure, the Web Infrastructure Model, originally proposed by Fett, Küsters, and Schmitz. Our analysis has uncovered several types of attacks on the protocol, violating the aforementioned security goals set by the FAPI WG. We subsequently have worked with the FAPI WG to fix the protocol, resulting in several changes to the specifications. After adapting our model to the changed specifications, we have proved the security properties to hold under the strong attacker model defined by the FAPI WG.

This technical report embodies the full version of our CSF 2024 paper with the same title.

## I. INTRODUCTION

Web-based authentication and authorization is ubiquitous. Many websites and applications can be used by logging in with a so-called Identity Provider, for instance, “Login with Google” or “Login with Facebook”, generally dubbed “social login”, or more generally, *Single Sign-On (SSO)*. It is also possible to authorize applications, including websites, but also IoT devices, such as routers and smart TVs, to access resources managed by the Identity Provider. Such resources include email addresses [24, 74], documents/calendars/pictures/movies stored in the cloud [59, 61], YouTube accounts [81], or development repositories [74]. A widely used protocol family for these authorization and authentication use cases are the OAuth 2.0 and OpenID Connect protocols [79, 82].

While plain OAuth 2.0 and OpenID Connect are suitable for typical low-risk use cases (e.g., social login), many use cases have emerged in high-risk settings for both authorization and authentication scenarios: Third party services can be authorized to, e.g., get access to bank transaction histories for monitoring and feedback [4, 10], trigger financial transactions [62, 63, 80], access cars [12], perform health-related actions like managing electronic prescriptions [23], or access medical records [19, 76]. In such high-risk use cases, attacks that enable malicious actors to access resources or impersonate end-users not only have more severe consequences than in classical low-risk settings, but such use cases also require more robust protocols and overall stronger security guarantees. For example, when using SSO to manage access to health records, it is important to not only recognize the same user again at a later point, which is typical for social network SSO, but to provide the full legal identity of the user. Likewise, protocols for high-risk use cases should be robust, i.e., provide security even if some messages or relevant values leak to an attacker, e.g., through leaked server log files [1].

To provide a higher level of security for such use cases, the OpenID Foundation developed FAPI 1.0, which is based on OAuth 2.0 and OpenID Connect, but uses many additional mechanisms to increase their security, e.g., to guarantee authorization and authentication even if the attacker can misconfigure certain endpoints or certain TLS-protected messages leak to the attacker, e.g., via log files. Meeting many ecosystems’ needs, FAPI 1.0 is now in widespread use, e.g., as part of the UK’s Open Banking standards [62], in the Australian Government’s (mandatory to support) *Consumer Data Standards* [16] which govern customer interaction with banks, energy, and telecommunications companies, with more industry sectors to follow. Further FAPI 1.0 uses include Brazil’s Open Finance and Open Insurance programs [63, 64], companies like yes.com [80] with an ecosystem of more than 1,000 actively participating banks, acting as trusted identity providers that can be used by sites and apps to login users, identify natural persons, obtain account information, and initiate payments, as well as US-based Financial Data Exchange FDX with more than 42 million users [27], and New Zealand’s core payment clearing house payments.nz [40].

The high security goals FAPI 1.0 aims to achieve have been formally analyzed before [38], uncovering several attacks and proposing fixes for those. Based on the experiences with FAPI 1.0, including interoperability and implementation aspects, the OpenID Foundation is currently standardizing a successor named FAPI 2.0, which comprises a completely new protocol (see also Section IV-C).

FAPI 2.0 is a framework of specifications, with the core protocol specified in the *FAPI 2.0 Security Profile* [32]. Another important specification for our purposes in this framework is the *FAPI 2.0 Attacker Model* [29], which captures the security goals that the protocol aims to fulfill, along with assumptions on the attacker capabilities, resulting in a strong attacker, far exceeding standard attacker models for protocol analysis (see Section II-C). In the following, we will often refer to the FAPI 2.0 Security Profile and FAPI 2.0 Attacker Model as just FAPI 2.0.<sup>1</sup>

Despite being a very recent standard, FAPI 2.0 is expected to be adopted soon in many important ecosystems, with several of the aforementioned FAPI 1.0 users already having committed to switch to FAPI 2.0.

Given the importance of FAPI 2.0 and its current and future use in high-risk environments, the FAPI WG has asked us to accompany the standardization process with a formal analysis, providing feedback early on and throughout the development of FAPI 2.0. We hence have performed a detailed formal security analysis of the FAPI 2.0 Security Profile with the security goals specified in the FAPI 2.0 Attacker Model. Our analysis is based on the Web Infrastructure Model (WIM) [33], a symbolic Dolev-Yao style model of the Web infrastructure. The WIM is the most comprehensive and detailed model of the Web infrastructure to date. Other models of the Web and its infrastructure, such as those by Pai et al. [65], Kumar [49, 50], or Bansal et al. [6] are significantly more abstract and limited by the tools upon which they are based.

For our work, we used the WIM as-is in its most recent published version [51], except for an additional HTTP header for DPoP proofs (see below for DPoP and Appendix D-I).

For the analysis of FAPI 2.0, we have created a formal model of the FAPI 2.0 Security Profile and formalized the security properties stated in the FAPI 2.0 Attacker Model and incorporated the assumptions on the attacker laid out therein (based on the state of the specifications at that time). We coordinated these steps with the FAPI WG to ensure a faithful modeling. The process of proving the properties within the model has revealed several attacks that break the goals of the protocol. We have proposed fixes and improvements to the specifications, which the FAPI WG appreciated and amalgamated into the official specifications, resulting in substantial changes thereof.

We adapted our formal model and formal security properties according to the changes discussed with the FAPI WG and were finally able to prove that the desired properties hold true within the model. Hence, our analysis reflects the latest official version of the FAPI 2.0 specifications [29, 32]. We note that while our analysis uncovered new attacks, we also found known attack patterns which the FAPI WG is familiar with and tried to avoid. This highlights the importance of a systematic formal analysis, which makes it possible to detect subtle flaws even in very complex protocols, where it is easy to overlook such flaws.

**Contributions.** In summary, our contributions are:

- We provide the first formal model of FAPI 2.0, along with a formalization of the security goals set forth in the FAPI 2.0 Attacker Model.
- Our analysis has uncovered several attacks, i.e., violations of the security goals under the attacker model defined by the FAPI WG.
- We propose fixes and improvements and worked with the FAPI WG to incorporate them, resulting in significantly modified and improved FAPI 2.0 specifications.
- We adapted our formal model to reflect the improved specifications and were then able to prove the formalized security goals.
- We have accompanied the development of the specifications from an early stage on and were able to support the standardization process with security recommendations before widespread deployment of the protocol.

**Structure of This Paper.** We first give a detailed description of the FAPI 2.0 protocol in Section II and then present the attacks that we have discovered in the process of our formal analysis in Section III. We describe our formal model and formal security theorem along with a proof sketch in Section IV. The full formal model and proofs are given in the appendix. Related work is discussed in Section V. We conclude in Section VI.

## II. FAPI 2.0 PROTOCOL AND SECURITY GOALS

In the following, we describe the FAPI 2.0 protocol and the accompanying FAPI 2.0 Attacker Model as of the beginning of our work. Recall that the FAPI 2.0 Attacker Model is part of the FAPI 2.0 specification framework. After consultation with the FAPI WG, we used FAPI 2.0 specifications as of June 1st, 2022 as a basis for our analysis ([30, 31], the FAPI 2.0 Security Profile used to be called *Baseline Profile*). We discuss changes to the specifications made due to our findings since then in Section III.

<sup>1</sup>While in the original FAPI 1.0 protocol, “FAPI” stands for “Financial-grade API”, the scope and expected uses of FAPI 2.0 reach far beyond the financial sector, thus, FAPI 2.0 is not an acronym anymore.

## A. Overview of FAPI 2.0

In a nutshell, FAPI 2.0 allows a user (also called resource owner) to grant a *client* application access to their data stored at a *resource server* (RS), by means of an *authorization server* (AS) which is responsible to manage access to the user’s data. In addition, the AS may provide the client with information on the user’s identity at the AS. For example, FAPI 2.0 may be used to grant an account aggregation service (client) read access to a user’s account balance at various banks (RSs), with services of these banks (ASs) managing such access (such services are in use today, e.g., [4, 10, 62, 63]).

On a high level, a FAPI 2.0 protocol run, also called *flow* or *grant*, advances as follows: A user visits a website or uses an application of the client *c*, which wants to access data of the user stored at the RS. Since the user’s data at the RS is managed by an AS *AS*, *c* contacts AS with some initial information, e.g., what kind of data *c* want to access. AS replies with an internal reference to the current flow, which *c* then forwards to the user’s browser while also instructing the browser to visit a website of AS to proceed. Once the user, or more precisely, their browser, visits that AS website, the user is asked to authenticate, e.g., with username and password, and to authorize the client’s request. If the user consents, AS instructs the user’s browser to return to the client website or application, passing on a value called the *authorization code*. Once the client receives that authorization code, it can contact AS and exchange the authorization code for so-called tokens. There are two types of tokens in FAPI 2.0: *ID Tokens* and *Access Tokens*. An id token contains information to identify the user, e.g., an email address or username with which the user is registered at the AS. This id data can be used by the client to authenticate users in the context of the client application. An access token, on the other hand, can be used by the client to request users’ resources from an RS, e.g., account balances. Upon receiving such a request, an RS verifies the access token’s validity. Depending on the access token format, this may include checking a signature on the access token or using so-called *token introspection*, which means that the RS queries the AS for validity information on a given access token.

## B. The FAPI 2.0 Protocol in Detail

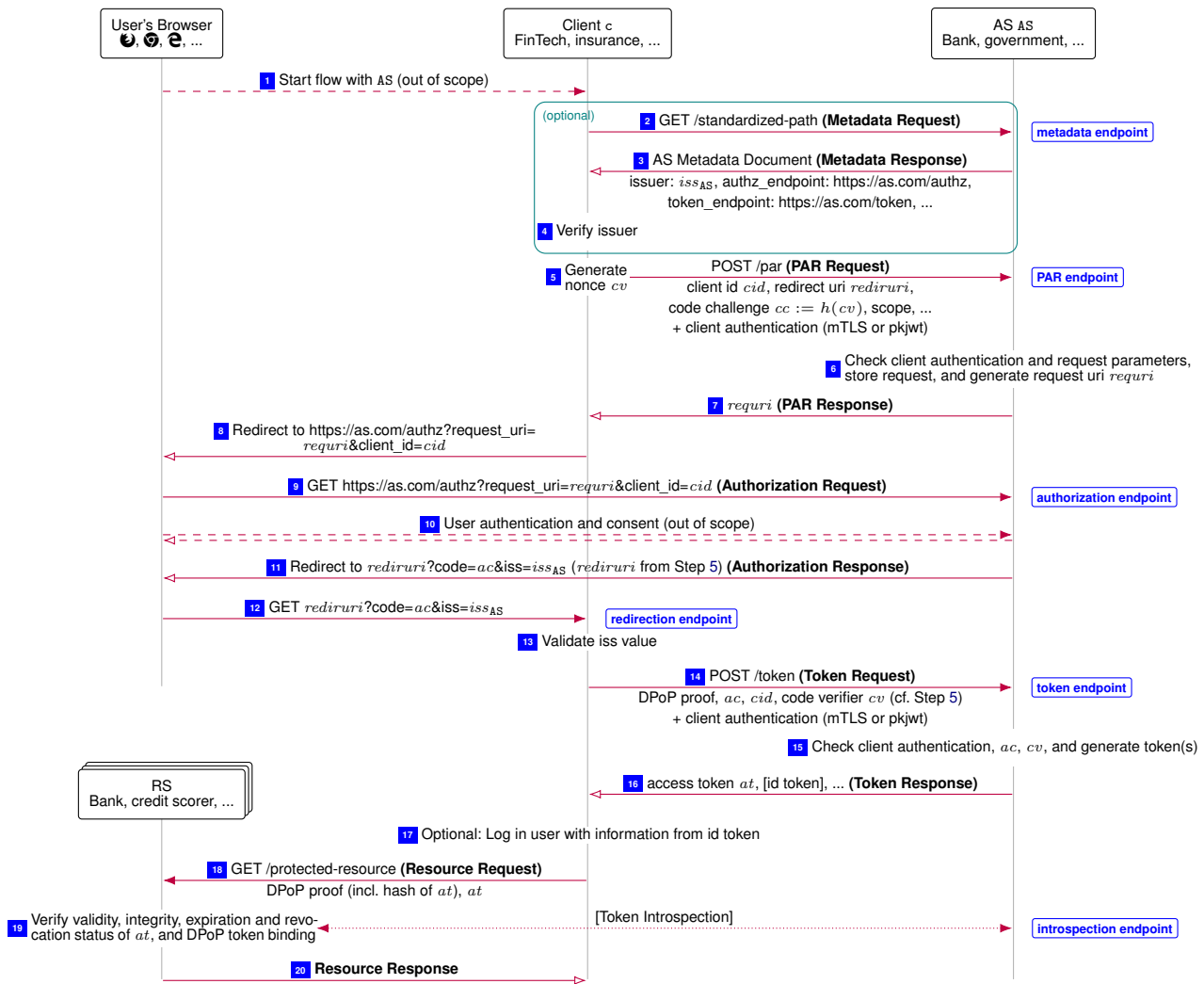
In the following, we describe a FAPI 2.0 protocol flow in detail (depicted in Figure 1). The flow is initiated by a user visiting the website or using an application of a client *c*, typically expressing the wish to authorize *c* using a certain AS *AS*, e.g., by clicking a “Login with AS” button (Step 1).

FAPI 2.0 assumes that *c* received the so-called *issuer identifier*  $iss_{AS}$  of AS (e.g., via configuration). That issuer identifier is used in FAPI 2.0 and other protocols to uniquely identify AS [75]. However, to complete a FAPI 2.0 flow, *c* needs additional knowledge on AS, e.g., endpoint URLs. If *c* does not yet know all necessary values (e.g., via configuration), it can proceed by fetching so-called *Authorization Server Metadata* [48, 71] from AS (Step 2). However, this step is optional. Like all other communication in FAPI 2.0, this exchange is done via HTTPS, i.e., is protected by TLS. The metadata returned by AS includes URIs of the relevant endpoints, supported cryptographic algorithms, and similar information, along with the issuer identifier of AS (Step 3). Once the client acquired the metadata, it verifies the aforementioned issuer identifier to prevent mix-ups, e.g., due to injection attacks.

Once the required values are available, *c* assembles a *Pushed Authorization Request* (PAR) [57] and sends it to AS (Step 5). This PAR request contains everything needed by AS to provide the user with sufficient information in Step 10 such that the user can make an informed decision on whether to grant *c* access to their data. This information includes: 1) a *client id* *cid*, uniquely identifying *c* at AS. 2) A *scope* value, describing what data *c* wants to access, e.g., “read transactions”, and whether *c* requests an id token to be issued. 3) A *redirect uri* *rediruri*, which is used by AS in Step 11 to redirect the user’s browser back to *c*. 4) A *code challenge*, i.e., a hash  $h(cv)$  of a client chosen nonce *cv*, which is used in Step 14 to verify that the client requesting a token is the same client that sent the PAR request (even if the PAR request leaks). This mechanism is called *Proof Key for Code Exchange* (PKCE) [72]. 5) Client authentication information (see below for a description).

Upon receiving the PAR request, AS verifies the client authentication, the presence of the parameters explained above, and checks whether the requested scope can be granted to the client (under the policies of AS). If all these checks pass, AS creates a random *request uri* *requri* and stores the requested scope, *cid*,  $cc := h(cv)$ , *rediruri*, and *requri* (Step 6); *requri* will be used as a reference to the PAR data in Step 9 and is therefore sent to *c* in the PAR response (Step 7). Client *c* then redirects the user’s browser to AS, adding *requri* and *cid* as request parameters (Step 8). Following that redirect, the user’s browser visits AS and in doing so, forwards *requri* and *cid*, hence providing information on the user’s context (i.e., the current flow) to AS (Step 9). The user now authenticates at AS and reviews the access requested by *c* (Step 10), the exact details of this step are up to the AS and out of scope of FAPI 2.0. If the user consents, AS generates a random *authorization code* *ac* and stores it with the PAR data from Step 5. AS then redirects the user’s browser back to the *rediruri* of *c* (stored in Step 6), and includes *ac* as well as an *iss* value [75] (i.e., the issuer identifier  $iss_{AS}$ ) as parameters (Steps 11 and 12).

Once *c* has received the browser’s (redirected) request, it validates the *iss* value by comparing it to the issuer identifier of the AS to which *c* sent the PAR request in Step 5 to prevent mix-up attacks [36, 55, 58, 75]. If this check passes, *c* sends a *token request* to AS (Step 14). This token request contains the authorization code *ac* from Step 12, client id *cid*, a *code verifier* *cv*, i.e., the nonce from Step 5, and client authentication similar to Step 5. Furthermore, *c* must include information for access token sender constraining, which we describe below.



**Figure 1.** FAPI 2.0 Security Profile protocol flow (with DPoP sender constraining)

When AS receives that token request, it verifies the client authentication, presence of a sender constraining method, and validity of the authorization code and code verifier (Step [15]). The latter is verified by checking whether  $h(cv) = cc$ , with  $cc$  being the code challenge stored in Step [6] and  $cv$  being the code verifier from the token request. The code  $ac$  is then invalidated and AS generates an access token  $at$  (and id token if requested) and sends them back to  $c$  in Step [16].

Given an id token,  $c$  may now log in the user with whatever identity the user has at AS, e.g., a user name (Step [17]). This allows clients to offer SSO to their users.

Using the access token  $at$ ,  $c$  can request user's resources at an RS as follows: in the resource request (Step [18]),  $c$  must include  $at$  as well as corresponding information for access token sender constraining (see below). The RS then has to verify  $at$ 's validity, integrity, expiration, and revocation status, as well as the sender constraining information (Step [19]). Except for the sender constraining, FAPI 2.0 does not specify how RSs should perform those (nonetheless mandatory) checks. Currently, there are two widely-adopted methods to do so [60]: token introspection [69], and structured access tokens, which contain the necessary information and are typically signed by the AS [8, 46]. With token introspection, the RS sends the access token to the introspection endpoint of the AS which issued the token, to which the AS answers with information on the validity of the token and on the public key to which the access token is bound.

**Client Authentication.** FAPI 2.0 requires ASs to authenticate clients at the PAR and token endpoints (Steps [5] and [14]) using *Mutual-TLS (mTLS)* or *private\_key\_jwt*. In both cases, clients need to be registered with the AS beforehand. With mTLS [11], the client presents a TLS certificate containing the client's identity, e.g., its client id, during TLS connection establishment. With *private\_key\_jwt* [70], the client adds a signed *JSON Web Token (JWT)* [44, 46, 47] to its messages. This JWT contains the client's id at the AS, the issuer identifier of the AS, and a nonce, and is signed with a private key of the client.

**Access Token Sender Constraining.** When issuing an access token (Steps [14]–[16]), a FAPI 2.0 AS is required to bind the

token to a key of the client who requested it. Likewise, the RS must verify this binding when it receives a resource request (Step [19]). FAPI 2.0 defines two methods to establish and verify such a binding: *OAuth 2.0 Demonstrating Proof-of-Possession at the Application Layer (DPoP)* [39], which is shown in Figure 1, and mTLS [11]. In both cases, the access token is bound to a client key pair, e.g., by including a hash of the public key in the token, and the client has to include a proof of possession of the private key when using the access token.

With DPoP, the token request (Step [14]) must include a *DPoP proof*: a signed JWT *dpopJWT*, containing the URL to which it is sent, a nonce, and a public verification key  $\text{pub}(k)$  (of the client’s choice). *dpopJWT* is signed using the corresponding private key  $k$ . The AS then binds the access token to  $\text{pub}(k)$ . When requesting resources (Step [18]), the client has to include another DPoP proof—signed with  $k$ —which must contain a hash of the access token in addition to the aforementioned items.

With mTLS, the AS binds the access token to the public key included in the client’s TLS certificate, which the client presents during connection establishment in Step [14]. When using the access token (Step [18]), the client presents the same certificate during the TLS connection establishment (which includes a proof of possession of the corresponding private key).

We emphasize that client authentication and access token sender constraining are independent of each other, including the key material. E.g., a client which uses mTLS to authenticate may use DPoP for sender constraining, and a client can authenticate with `private_key_jwt` and at the same time use mTLS for sender constraining. I.e., there are four possible combinations.

### C. Security Goals and Attacker Model

Along with the actual protocol specification, the FAPI WG developed the FAPI 2.0 Attacker Model [30] which outlines security goals and assumptions on attackers under which these goals are expected to hold. As before, we describe the state as of June 1st, 2022 here and discuss changes made since then in Section III. The formalized security properties and modeling of attacker assumptions are presented in Section IV.

**Authorization Goal.** The authorization goal states that no attacker should be able to access resources belonging to an honest user. In addition, the FAPI 2.0 Attacker Model states that this goal is “fulfilled if no attacker can successfully obtain and use an access token” issued for an honest user.

**Authentication Goal.** The authentication goal is fulfilled when no attacker is able to log in at a client under the identity of an honest user.

**Session Integrity for Authorization Goal.** Session integrity goals aim to prevent attackers from tricking users into using attacker’s resources or identities. Hence, the session integrity for authorization goal ensures users cannot be forced to use resources of the attacker.

**Session Integrity for Authentication Goal.** Similar to the session integrity for authorization goal, the session integrity for authentication goal is fulfilled if no attacker can force an honest user to be logged in under an identity of the attacker.

**Attacker Assumptions.** In the following, we summarize the aforementioned attacker assumptions laid out in the FAPI 2.0 Attacker Model. We stress that these (strong) assumptions are part of and justified by the specification.

A1. The attacker controls the network, i.e., can intercept, block, and tamper with all messages sent over the network. In particular, the attacker can also reroute, reorder, and create (from its knowledge) new messages. However, the attacker cannot break cryptography unless it learned the respective keys. Nevertheless, the attacker can pose as any party (and any network participant) in the protocol. In addition, the attacker can also send links to (honest) users which are then visited by these users. See [30, Sec. A1, A1a, A2].

A2. The attacker can read authorization requests in plain (cf. Step [9] in Figure 1). See [30, Sec. A3a].

A3. The attacker can read authorization responses in plain (Step [12]). See [30, Sec. A3b].

A4. The attacker can trick the client into using an attacker-controlled token endpoint URL, i.e., one for which the attacker can obtain a valid TLS certificate (other endpoints, e.g., PAR, are not affected). Hence, the attacker can read token requests (Step [14]) in plain and construct arbitrary token responses from its knowledge (Step [16]). However, this assumption only applies to clients which do not use the AS metadata mechanism. See [30, Sec. A5].

A5. Resource requests (Step [18]) leak to the attacker in plain. See [30, Sec. A7].

A6. Resource responses (Step [20]) leak to the attacker in plain. See [30, Sec. A7].

A7. The attacker can modify resource responses (Step [20]). I.e., the attacker can replace an honest RS’ resource response with its own message without the client noticing, even though that response is protected by TLS. Note that this does *not* give the attacker the ability to replace arbitrary messages in TLS connections, but is limited to resource responses. See [30, Sec. A8].

## III. ATTACKS

We formally modeled the FAPI 2.0 specifications (as of June 1st, 2022) and then formalized and tried to prove the security goals according to the FAPI 2.0 Attacker Model [30] under the attacker assumptions outlined in the same document. In the course of this analysis, we have uncovered a number of attacks, i.e., violations of the security goals laid out in the FAPI 2.0 Attacker Model. We have discussed these findings with the FAPI WG and worked with them to resolve the issues, resulting in a number of changes to the specifications which we explain here. The formal model presented in Section IV, for which we prove security, incorporates

these changes. While, to the best of our knowledge, Attacker Token Injection, the Client Impersonation attacks, and DPOP Proof Replay are completely new attacks, interestingly, for the other attacks (Browser Swapping, Cuckoo’s Token, Authorization Request Leak), similar attack patterns have been reported for related protocols [38, 54]. This emphasizes the importance of systematic, formal analysis, as even seasoned experts overlook known attack patterns for complex protocols.

Due to space constraints, we present additional attacks and their fixes as well as inconsistencies in the specifications discovered during our analysis in [Appendix B](#).

### A. Attacker Token Injection

This two-phased attack violates both session integrity goals and requires attacker assumptions [A4](#) (token endpoint misconfiguration) and [A5](#) (resource requests leak, see [Section II-C](#)). In the first phase, the attacker, posing as a user, completes a flow with an honest client  $c_{\text{hon}}$  and honest AS  $AS_{\text{hon}}$ . During this flow, the attacker uses [A5](#) to obtain  $at_{\text{att}}$ , i.e., an access token bound to keys of  $c_{\text{hon}}$ , issued by  $AS_{\text{hon}}$  for resources of the attacker. For the second phase, the attacker uses [A4](#), such that  $c_{\text{hon}}$  uses an attacker-controlled token endpoint (instead of  $AS_{\text{hon}}$ ’s token endpoint). Hence, when an honest user  $u$  starts a flow with  $c_{\text{hon}}$  and  $AS_{\text{hon}}$ , the attacker receives  $c_{\text{hon}}$ ’s token request, to which the attacker answers with  $at_{\text{att}}$  and an id token constructed by the attacker for an attacker identity. Client  $c_{\text{hon}}$  then logs  $u$  in under the attacker’s identity and uses attacker resources in a session with  $u$ .

At its core, this attack is possible because the attacker can trick an honest client into using an attacker-controlled token endpoint (due to [A4](#)). We describe a fix, our communication with the FAPI WG, and resulting changes to the specifications at the end of the next subsection.

### B. Client Impersonation Attacks

This attack violates the authorization goal and requires attacker assumption [A4](#) (token endpoint misconfiguration, see [Section II-C](#)). The attack targets an honest client  $c_{\text{hon}}$ , which authenticates itself to an honest AS  $AS_{\text{hon}}$  with the `private_key_jwt` method (see [Section II-B](#)). Using [A4](#), the attacker modifies  $c_{\text{hon}}$ ’s configuration such that  $c_{\text{hon}}$  uses an attacker-controlled token endpoint (instead of  $AS_{\text{hon}}$ ’s token endpoint).

In a nutshell, the attacker first obtains two valid client authentication JWTs  $pkjwt_1$  and  $pkjwt_2$  for  $c_{\text{hon}}$  at  $AS_{\text{hon}}$ . Afterwards, the attacker uses those JWTs to impersonate  $c_{\text{hon}}$  at  $AS_{\text{hon}}$  and obtains an access token issued by  $AS_{\text{hon}}$  for resources of an honest user  $u$ , with the token being bound to a key of the attacker, i.e., the attacker can use this token at an RS to access  $u$ ’s resources.

To obtain said JWTs, the attacker starts a flow with  $c_{\text{hon}}$ , selects  $AS_{\text{hon}}$ , and authenticates at  $AS_{\text{hon}}$  (with an attacker id). However, once the attacker receives  $c_{\text{hon}}$ ’s token request (due to the misconfigured token endpoint), no further actions are taken. The token request sent by  $c_{\text{hon}}$  contains  $pkjwt_1$  to authenticate  $c_{\text{hon}}$  at  $AS_{\text{hon}}$ , which has now leaked to the attacker and has never been sent to  $AS_{\text{hon}}$ , i.e.,  $pkjwt_1$  is still valid. The attacker can repeat this to obtain  $pkjwt_2$ .

With such valid client authentication JWTs for  $c_{\text{hon}}$  at  $AS_{\text{hon}}$ , the attacker can proceed as depicted in [Figure 2](#): An honest user  $u$  starts a flow with  $c_{\text{hon}}$  and expresses their wish to use an AS  $AS_{\text{att}}$ , identified by issuer identifier  $iss_{\text{att}}$ , which happens to be controlled by the attacker ([Step 2](#)). Hence,  $c_{\text{hon}}$  sends a PAR request to  $AS_{\text{att}}$ , containing its client id  $cid$  at  $AS_{\text{att}}$ , a redirect uri, a code challenge  $h(cv_{\text{hon}})$  for a code verifier  $cv_{\text{hon}}$  chosen by  $c_{\text{hon}}$ , etc. as described in [Section II-B](#) ([Step 3](#)). Instead of replying to this request immediately, the attacker now poses as  $c_{\text{hon}}$  towards  $AS_{\text{hon}}$ : the attacker sends its own PAR request to  $AS_{\text{hon}}$ , assembled from an attacker-chosen code challenge  $h(cv_{\text{att}})$ , a redirect uri  $rediruri_{\text{att}}$  pointing to a URL of  $AS_{\text{att}}$ , the client id  $cid'$  of  $c_{\text{hon}}$  at  $AS_{\text{hon}}$ , and  $pkjwt_1$  ([Step 4](#)). Upon receiving the attacker’s PAR request,  $AS_{\text{hon}}$  validates the request and replies with  $requiri$  ([Step 5](#)), which the attacker forwards to  $c_{\text{hon}}$  (now again in the role of  $AS_{\text{att}}$  towards  $c_{\text{hon}}$ ) in [Step 6](#) in response to  $c_{\text{hon}}$ ’s original PAR request from [Step 3](#).

Client  $c_{\text{hon}}$  now instructs  $u$  to visit  $AS_{\text{att}}$  for authentication ([Step 7](#)). However, instead of the usual login page,  $AS_{\text{att}}$  responds with a page luring the user into clicking a link, e.g., by explaining that it is now cooperating with  $AS_{\text{hon}}$  ([Step 9](#)). This link points to  $AS_{\text{hon}}$ ’s authorization endpoint and contains the parameters  $requiri$  and  $cid'$ , i.e., the client id of  $c_{\text{hon}}$  at  $AS_{\text{hon}}$ ; instead of a link,  $AS_{\text{att}}$  could also just redirect  $u$  to  $AS_{\text{hon}}$  directly. Hence,  $u$  ends up authenticating at  $AS_{\text{hon}}$  and authorizes  $c_{\text{hon}}$  – recall that  $u$  expects to authorize  $c_{\text{hon}}$  here ([Step 11](#)). Following  $u$ ’s consent,  $AS_{\text{hon}}$  redirects  $u$  with authorization code  $ac$  to  $rediruri_{\text{att}}$  received in [Step 4](#), i.e., to an attacker-controlled location ([Step 12](#)), which  $u$  follows ([Step 13](#)).

Having received  $ac$ , the attacker can now construct a valid token request (using  $pkjwt_2$ ) as shown in [Step 14](#), and subsequently receives an access token  $at$  for  $u$ ’s resources. Note that from  $AS_{\text{hon}}$ ’s point of view,  $u$  authorized the attacker, posing as  $c_{\text{hon}}$  towards  $AS_{\text{hon}}$ , to receive  $at$ . In addition, recall that the DPOP key to which  $AS_{\text{hon}}$  binds  $at$  is chosen by the sender of the token request ([Step 14](#)), i.e., the attacker. Hence,  $at$  is bound to a key of the attacker, and can be used at an RS to access  $u$ ’s resources, thus breaking the authorization goal.

We stress that  $u$  authenticates at an honest, trusted AS and authorizes not only an honest and trusted client, but also exactly the client  $u$  expected to authorize. There are some variants of this attack with slightly different preconditions, but similar outcome, which we describe in [Appendix B-B](#).

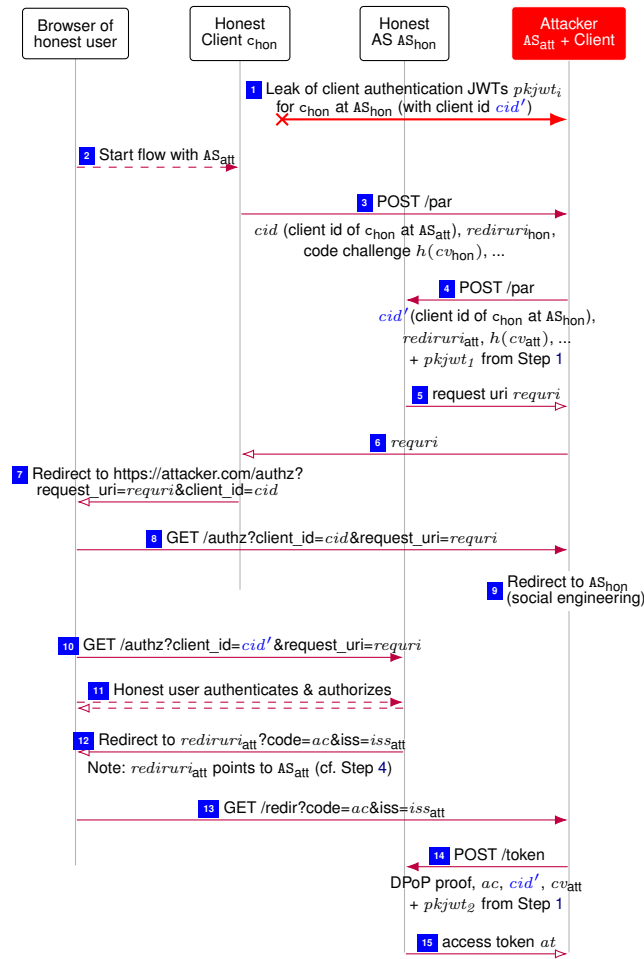


Figure 2. Client impersonation attack

These attacks emerged when we tried to prove that a client authentication JWT for authentication of an honest client at an honest AS  $AS_{hon}$  (i.e., the JWT contains the issuer identifier of  $AS_{hon}$ ) cannot leak to an attacker (Lemma 9).

**Fix.** Since the possible misconfiguration of an honest client’s token endpoint is the root cause for both the attacker token injection and the client impersonation attacks and FAPI 2.0 ASs are already required to serve a metadata document, our proposed fix of mandating clients to request and use this metadata was adopted by the FAPI WG [20, 25, 42]. Recall  $\mathcal{A}4$ : such token endpoint misconfiguration is only considered for clients which do not use the AS metadata mechanism.

### C. DPoP Proof Replay

With attacker assumption  $\mathcal{A}5$  (resource requests leak) from Section II-C, the attacker can read resource requests in plain and hence can try to replay them at the RS (cf. Step [18] in Figure 1), thus violating the authorization goal (see Section II-C) if the RS accepts the replayed request.

When DPoP sender constraining is used, the attacker can indeed replay the client’s DPoP proof (using the attacker’s TLS keys for the underlying connection): neither FAPI 2.0, nor DPoP [39] itself, nor the specifications on which DPoP is built [45, 46, 70] mandate for DPoP proofs to be strictly one-time use, hence, the RS does not reject the replayed proof.

Our initial attempts to prove the authorization property (see Definition 2) revealed a second variant, which also violates the authorization goal: with  $\mathcal{A}1$  (network attacker), the attacker can additionally block the honest client’s request, i.e., from the point of view of the RS, the attacker’s resource request is not even a replay, and hence, replay protection does not completely prevent this attack.

Note that either variant of this attack is not possible with mTLS sender constraining since the attacker cannot even establish an mTLS connection with the client’s public mTLS key (to which the access token is bound when mTLS sender constraining is used).

**Fix.** In our discussions with the FAPI WG, it became clear that the FAPI WG formulated attacker assumption  $\mathcal{A}5$  with leaks of RS server log files in mind. Hence, the FAPI 2.0 Attacker Model was changed to clarify that resource requests leak *after* processing by the RS [13]. However, this only resolves the problem resulting from the attacker blocking the honest client’s request, but does

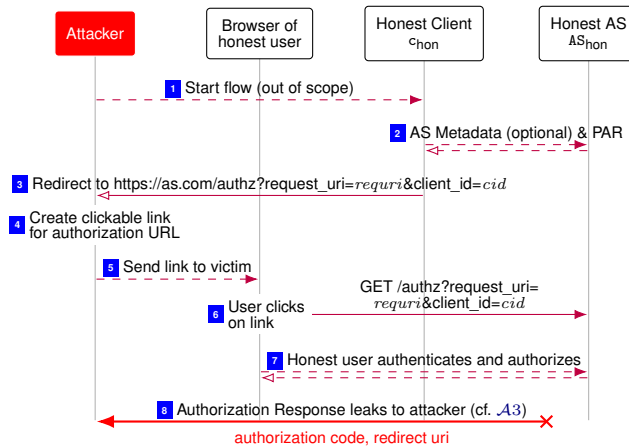


Figure 3. Browser swapping attack

not prevent the replay attack. Hence, we proposed to fix this attack by mandating the use of *resource server-provided nonces* with strict one-time use enforcement by the RS [39, Sec. 9]. The server-provided nonce mechanism is an optional part of DPoP in which—in a challenge-response manner—the client first requests a nonce from the RS which the client then has to include in its DPoP proof. We validated effectiveness of this fix in our formal model.

The FAPI WG acknowledged the attack [22], and added a description of the attack, as well as several options to fix it, to the specification.

#### D. Browser Swapping Attack

In this attack, the attacker violates the authorization and authentication goals by combining attacker assumptions  $\mathcal{A}1$  (network attacker) and  $\mathcal{A}3$  (authorization responses leak, see Section II-C).

On a high level, the attacker poses as a user and starts a flow with an honest client  $c_{hon}$  and honest AS  $AS_{hon}$ , but tricks an honest user  $u$  into logging in at  $AS_{hon}$  and to authorize access for  $c_{hon}$ . After  $u$  authorized  $c_{hon}$ , the attacker continues its session with  $c_{hon}$  (as if the attacker authenticated and authorized at  $AS_{hon}$ ). Hence, the attacker gets logged in at  $c_{hon}$  under the identity of  $u$ , and  $c_{hon}$  provides access to  $u$ 's resources to the attacker.

The detailed attack flow is depicted in Figure 3: in Step 1, the attacker initiates a flow at  $c_{hon}$  with  $AS_{hon}$ . Then,  $c_{hon}$  and  $AS_{hon}$  exchange the PAR as described in Section II-B (Step 2). Following this,  $c_{hon}$  instructs its user, i.e., the attacker, to visit the authorization endpoint of  $AS_{hon}$  with the request uri  $requri$  from Step 2 and  $c_{hon}$ 's client id  $cid$ . However, instead of following this redirect, the attacker creates a clickable link, pointing to  $AS_{hon}$ 's authorization endpoint, with  $requri$  and  $cid$  as parameters (Step 4). The attacker then sends this link to  $u$  in Step 5, e.g., in an email or as part of an attacker website.

Once  $u$  follows that link (e.g., by means of social engineering, see  $\mathcal{A}1$ ) in Step 6,  $u$  will be asked to authenticate (if not already logged in at  $AS_{hon}$ ) and to authorize  $c_{hon}$  to access  $u$ 's resources (Step 7). We emphasize that, similar to  $AS_{hon}$ ,  $u$  might trust the (honest) client  $c_{hon}$ . Once  $u$  consents, the attacker blocks all further communication for  $u$  (see  $\mathcal{A}1$ ).

Now recall  $\mathcal{A}3$ : through a leaking authorization response, the attacker can obtain the authorization code  $ac$  (Step 8). With  $ac$ , the attacker visits  $c_{hon}$ 's redirect uri, so  $c_{hon}$  can subsequently exchange  $ac$  for an access (and id) token at  $AS_{hon}$ . From  $c_{hon}$ 's point of view, these tokens are associated with  $c_{hon}$ 's session with the attacker. However, due to Step 7, these tokens are issued for the (honest user)  $u$ 's resources (and identity).

At the heart of this attack is the lack of a strong connection between the sessions user–client and user–AS. We discovered it when trying to prove the authorization property (see Definition 2): while proving that in a flow between an honest client, honest AS, and honest user, the client does not leak the user's resources, we have to prove, among others, that the authorization code associated with the flow cannot be sent to the client's redirection endpoint by the attacker.

**Fix.** The FAPI WG acknowledged this attack and after several discussions, there was consensus that this attack cannot be fixed with currently deployed methods [9, 66]. This decision is documented and explained along with the attack in the specifications [29, Sec. 6.5.7]. Consistent with this decision, the FAPI WG also removed attacker assumption  $\mathcal{A}3$ , i.e., FAPI 2.0 no longer claims to fulfill its security goals when authorization responses leak.

#### E. Cuckoo's Token Attack

In this attack, the attacker leverages attacker assumption  $\mathcal{A}5$  (resource requests leak) to violate the authorization goal (see Section II-C).

Said attacker assumption allows the attacker to obtain an access token  $at_{hon}$  from an honest flow, i.e.,  $at_{hon}$  was issued by an honest AS  $AS_{hon}$  for an honest client  $c_{hon}$  on behalf of an honest user  $u$  (and thus  $at_{hon}$  is bound to keys of  $c_{hon}$ ). By then injecting



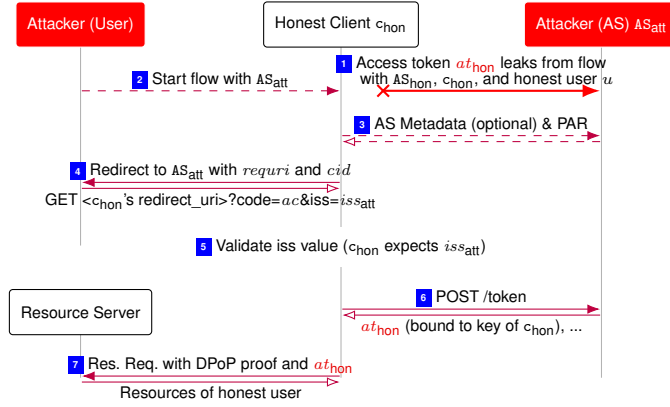


Figure 4. Cuckoo’s token attack

$at_{hon}$  into a flow between the attacker as user,  $c_{hon}$ , and an attacker-controlled AS, the attacker (as user of  $c_{hon}$ ) gains access to  $u$ ’s resources.

Figure 4 shows the attack in detail: the attacker first acquires an access token  $at_{hon}$  from an honest flow (Step 1). Such a token may, for example, be obtained by the attacker through observing a resource request (see A5). Due to access token sender constraining, the attacker cannot use  $at_{hon}$  directly at an RS (recall:  $at_{hon}$  is bound to keys of  $c_{hon}$ ). Instead, the attacker, posing as a user towards  $c_{hon}$ , now starts a flow with  $c_{hon}$ , selecting an AS  $AS_{att}$  controlled by the attacker (Step 2).  $c_{hon}$  initiates the flow as usual with PAR (Step 3), followed by instructing the user, i.e., attacker, to visit  $AS_{att}$  (Step 4). Since  $AS_{att}$  is controlled by the attacker, the authentication and authorization steps can be skipped and the attacker (posing as a user) immediately “redirects” itself to  $c_{hon}$  with issuer identifier  $iss_{att}$  (which identifies  $AS_{att}$ ) and an arbitrary authorization code  $ac$  (Step 4). Upon receiving that message,  $c_{hon}$  validates the issuer identifier, which succeeds:  $c_{hon}$  is and wants to be in a flow with  $AS_{att}$ . As usual,  $c_{hon}$  now sends a token request to  $AS_{att}$ , which responds with the previously acquired  $at_{hon}$  (Step 6). Recall that  $at_{hon}$  was issued for and is bound to keys of  $c_{hon}$  to access resources of  $u$ . Hence,  $c_{hon}$  can use  $at_{hon}$  at an honest RS. But since  $c_{hon}$  associates  $at_{hon}$  with the session between  $c_{hon}$  and the attacker posing as a user, this gives the attacker access to  $u$ ’s resources through  $c_{hon}$ .

At its core, this attack exploits the lack of binding between access token and AS from the client’s point of view. Note: the client is mandated to handle the access token as an opaque value, i.e., cannot perform any checks on the token.

**Fix.** We proposed to fix this attack by mandating the client to include an AS issuer identifier in each resource request (which would be the attacker AS’ identifier in the example above). The RS can then compare this issuer identifier sent by the client with the actual issuer of the access token (which, in our example, would be some different, honest AS). Note that in order to verify the token’s validity, the RS already needs a way to know which AS originally issued the token. The FAPI WG acknowledged the attack [17] and added a description of the attack, as well as our proposed fix, to the specification [26]. However, the FAPI WG decided that the “preconditions for this attack do not apply to many ecosystems and require a powerful attacker” [32, Sec. 5.6.5] and hence made implementation of a fix optional. Note that while this fix seems to be a small change, there is no standardized way to send the issuer identifier of the AS, as well as no standardized way for the RS to get a value to compare against [32, Sec. 5.6.5]. Hence, mandating such a change would require substantial standardization efforts.

#### IV. FORMAL ANALYSIS

In this section, we describe our formal analysis. We start with a primer on the WIM, continue with a description of our formal model of FAPI 2.0, including its limitations, and our formalization of a FAPI 2.0 Web System. We then discuss the most important differences and technical challenges of our work compared to prior work on authorization protocols using the WIM. This is followed by a description of the formalized security properties and a proof sketch for the authorization property; see the appendix for the full formal model and full proofs of all security properties.

##### A. The Web Infrastructure Model

FAPI 2.0 is a Web-based protocol, and the interaction between browsers and Web servers introduces potential attack surfaces, e.g., by cross-site requests, in-browser communication, malicious scripts, insecure headers, or redirections. To account for attacks originating from the browser and complex interactions inside browsers, as well as between parties, we analyze FAPI 2.0 based on the Web Infrastructure Model (WIM) [33], which is the most detailed formal model of the Web infrastructure to date. The WIM is a Dolev-Yao (DY) style pen-and-paper web model and requires manual analysis. It has successfully been applied to several web standards, to uncover previously unknown attacks and to prove security properties [21, 33, 34, 36–38] (see Section V). So far

no mechanized analysis framework has such a comprehensive model of the Web. Mechanizing such a very detailed model, from scratch or on top of existing tools, is a big challenge by itself and out of the scope of this work.

In the following, we give a high-level overview of the WIM closely following the summary in [36], with the full model given in [Appendix G](#): the WIM is designed independently of a specific Web application and closely mimics published (de-facto) standards and specifications for the Web, for example, the HTTP/1.1 and HTML5 standards and associated (proposed) standards. The WIM defines a general communication model, and, based on it, Web systems consisting of Web browsers, DNS servers, and Web servers as well as Web and network attackers.

**Communication Model.** The main entities in the model are (*atomic*) *processes*, which are used to model browsers, servers, and attackers. Each process listens to one or more (IP) addresses. Processes communicate via *events*, which consist of a message as well as a receiver and a sender address. In every step of a run, one event is chosen non-deterministically from a “pool” of waiting events and is delivered to one of the processes that listens to the event’s receiver address. The process can then handle the event and output new events, which are added to the pool of events, and so on.

As usual in DY models (see, e.g., [2]), messages are expressed as formal terms over a signature  $\Sigma$ . The signature contains constants (for (IP) addresses, strings, nonces) as well as sequence, projection, and function symbols (e.g., for encryption/decryption and signatures). For example, in the Web model, an HTTP request is represented as a term  $r$  containing a nonce, an HTTP method, a domain name, a path, URI parameters, headers, and a message body. For example, a request for the URI <http://example.com/s?p=1> is represented as

$$r := \langle \text{HTTPReq}, n_1, \text{GET}, \text{example.com}, /s, \langle \langle p, 1 \rangle \rangle, \langle \rangle, \langle \rangle \rangle$$

where the body and the headers are empty. An HTTPS request for  $r$  is of the form  $\text{enc}_a(\langle r, k' \rangle, \text{pub}(k_{\text{example.com}}))$  where  $k'$  is a fresh symmetric key (a nonce) generated by the sender of the request (typically a browser); the responder is supposed to use this key to encrypt the response.

The *equational theory* associated with  $\Sigma$  is defined as usual in DY models. The theory induces a congruence relation  $\equiv$  on terms, capturing the meaning of the function symbols in  $\Sigma$ . For instance, the equation in the equational theory which captures asymmetric decryption is  $\text{dec}_a(\text{enc}_a(x, \text{pub}(y)), y) = x$ . With this, we have that, for example,

$$\text{dec}_a(\text{enc}_a(\langle r, k' \rangle, \text{pub}(k_{\text{example.com}})), k_{\text{example.com}}) \equiv \langle r, k' \rangle$$

i.e., these two terms are equivalent w.r.t. the equational theory.

A (*DY*) *process* consists of a set of addresses the process listens to, a set of states (terms), an initial state, and a relation that takes an event and a state as input and (non-deterministically) returns a new state and a sequence of events. The relation models a computation step of the process. It is required that the output can be computed (more formally, derived in the usual DY style) from the input event and the state.

The so-called *attacker process* is a DY process which records all messages it receives and outputs all events it can possibly derive from its recorded messages. Hence, an attacker process carries out all attacks any DY process could possibly perform. Attackers can corrupt other parties at any time; corrupted parties behave like the attacker process.

A *script* models JavaScript running in a browser. Scripts are defined similarly to DY processes. When triggered by a browser, a script is provided with state information, corresponding to the (browser) data available to JavaScript in real browsers. The script then outputs a term representing a new internal state and a command to be interpreted by the browser (see also the specification of browsers below). Similarly to an attacker process, the so-called *attacker script* may output everything that is derivable from its input.

A *system* is a set of processes. A *configuration*  $(S, E, N)$  of this system consists of the states  $S$  of all processes in the system, the pool of waiting events  $E$ , and an infinite sequence of unused nonces  $N$ . Systems induce *runs*, i.e., sequences of configurations, where each configuration is obtained by delivering one of the waiting events of the preceding configuration to a process, which then performs a computation step. Such a transition is called *processing step* and denoted by

$$(S, E, N) \xrightarrow[p \rightarrow E_{\text{out}}]{e_{\text{in}} \rightarrow p} (S', E', N').$$

Here, the process  $p$  processes the event  $e_{\text{in}}$  and creates the output events  $E_{\text{out}}$  which are added to the pool of waiting events of the next configuration.

A *Web system* formalizes the Web infrastructure and Web applications. It contains a system consisting of honest and attacker processes. Honest processes can be Web browsers, Web servers, or DNS servers. Attackers can be either *Web attackers* (who can listen to and send messages from their own addresses only) or *network attackers* (who may listen to and spoof all addresses and therefore are the most powerful attackers). A Web system further contains a set of scripts (comprising honest scripts and the attacker script) and a mapping of these scripts to strings. A Web system also defines the pool of initial events, which typically only contains so-called trigger events, which trigger pre-defined actions (see below for an example for pre-defined browsers actions).

**Web Browsers.** An honest browser is thought to be used by one honest user, who is modeled as part of the browser. User actions, such as following a link, are modeled as non-deterministic actions of the Web browser. User credentials are stored in the initial

state of the browser and are given to the respective Web pages, i.e., scripts. Besides user credentials, the state of a Web browser contains (among others) a tree of windows and documents, cookies, and Web storage data (localStorage and sessionStorage).

A *window* inside a browser contains a set of *documents* (one being active at any time), modeling the history of documents presented in this window. Each represents one loaded Web page and contains (among others) a script and a list of subwindows (modeling iframes). The script, when triggered by the browser, is provided with all data it has access to, such as a (limited) view on other documents and windows, certain cookies, and Web storage data. Scripts then output a command and a new state. This way, scripts can navigate or create windows, send XHRs and postMessages, submit forms, set/change cookies and Web storage data, and create iframes. Navigation and security rules ensure that scripts can manipulate only specific aspects of the browser’s state, according to the Web standards.

A browser will typically send DNS and HTTP(S) requests as well as XHRs, and it processes the responses. Several HTTP(S) headers are modeled, including, for example, cookie, location, strict transport security (STS), and origin headers. A browser, at any time, can also receive a trigger message upon which the browser non-deterministically chooses an action, for instance, to trigger a script in some document.

**Generic HTTPS Server.** The WIM defines a generic HTTPS server model which can be instantiated by application models. The generic server provides some generic functionality, e.g., a function for sending HTTPS requests, which internally handles DNS resolution and key management for symmetric transportation keys. The generic server also provides placeholder functions, e.g., for processing HTTPS requests and responses, which need to be instantiated by the application model.

## B. Overview of FAPI 2.0 Model

We created the application-specific model in the WIM based on the FAPI 2.0 specifications. For this, we instantiated the generic HTTPS server model provided by the WIM to create detailed formal models for ASs, clients, and RSs, including models of scripts. Our formal model covers all essential mechanisms used by FAPI 2.0, e.g., AS Metadata, PAR, PKCE, the mTLS and `private_key_jwt` client authentication methods, DPoP and mTLS token sender constraining, and ID tokens. Our model also covers both structured and opaque access tokens—which can be different for each flow (the token type is chosen non-deterministically for each token request)—and token introspection.

Using client authentication and token sender constraining as examples, we exemplify how our model closely follows the specifications and show (parts of) some concrete messages within our model. We then continue with details on how we incorporated the attacker model and fixes described in [Section III](#) into our formal model, followed by a discussion of the limitations of our model. More details specific to the AS, client, and RS models are given in [Appendix A](#) and we provide the full formal models in [Appendix D](#).

**Client Authentication.** Within the model, the ASs only accept PAR and token requests if they are client-authenticated, as mandated by the specifications [32, Sec. 5.3.1.2. No. 4], and clients always add client authentication when sending such requests [57, Sec. 2], [32, Sec. 5.3.2.1. No. 2]. Our model supports both mTLS and `private_key_jwt` authentication [32, Sec. 5.3.1.1. No. 6, Sec. 5.3.2.1. No. 2]. To illustrate how we model the `private_key_jwt` method, we show a client authentication JWT. Let *cid* be the client identifier of the client *c* at the AS *AS*, and let *iss<sub>AS</sub>* be the issuer identifier of AS. Within the model, `signkey(c)` denotes a private signing key of *c*. With these values, a client authentication JWT is—closely following the specification [70, Sec. 9]—represented by the term  $\text{sig}([\text{iss}: \text{cid}, \text{sub}: \text{cid}, \text{aud}: \text{iss}_{\text{AS}}], \text{signkey}(c))$ , where *iss* is the issuer of the JWT, *sub* is the subject that is being authenticated, and *aud* is the audience value (i.e., the AS for which the JWT is being created).

**Access Token Sender Constraining.** Our AS model returns access tokens that are sender constrained by mTLS or DPoP [32, Sec. 5.3.1.1. No. 4, 5], and the client model sends the token request to the AS and the resource request to the RS with the corresponding proof-of-possession [32, Sec. 5.3.2.1. No. 1]. As mandated by FAPI 2.0, our model of RSs requires access tokens to be sender constrained using one of these methods [32, Sec. 5.3.3. No. 5]. We illustrate how we model a DPoP proof in a token request, where, according to the DPoP specification [39, Sec. 4.1], a client *c* adds a DPoP proof to the HTTP headers of its request:  $\text{headers}[\text{DPoP}] := \text{dpopProof}$ , where *dpopProof* is a signed JWT as defined in [39, Sec. 4.2] (not to be confused with a client authentication JWT):  $\text{dpopProof} := \text{sig}(\text{dpopJwt}, \text{signkey}(c))$  with  $\text{dpopJwt} := [\text{headers}: [\text{jwk}: \text{pub}(\text{signkey}(c))], \text{payload}: [\text{htm}: \text{POST}, \text{htu}: \text{tokenEndpoint}]]$ . The *jwk* value is the public key that the receiver can use for verifying the signature, the *htm* value is the HTTP method of the request in which the DPoP proof is included, e.g., GET or POST, and the *htu* value is the URL of that request (but with empty parameters and fragment).

1) *Optional Fixes:* As explained in [Section III](#), the FAPI WG decided to make some fixes resulting from our attacks optional. In order to prove FAPI 2.0 secure, these fixes are required, hence our results only apply to implementations also employing these fixes. Here we describe how we deal with these fixes in the model.

**Cuckoo’s Token Attack.** As described in [Section III-E](#), the working group added (optional) countermeasures for preventing the attack. However, none of these fixes are mandatory, and each of them might mask other attacks. Thus, we decided to model a minimal fix that specifically targets this attack: right before sending the resource request to an RS, the client model checks whether the requested resource is managed by the AS from which it got the access token.

**DPoP Proof Replay.** For preventing this attack, we modeled replay protection using server-provided DPoP nonces (see also Section III-C). We require these nonces to be one-time use only, i.e., the RS model invalidates them after one use.

2) *FAPI 2.0 Attacker Model:* We initially modeled all attacker assumptions (see Section II-C), but as described in Section III, the FAPI WG decided to remove some of them in response to our analysis. We describe how we modeled the final assumptions (under which our proofs hold) and refer to Appendix C for the removed ones:

*A1* is part of the WIM, see Section IV-A.

*A2* (authorization requests leak) is modeled by leakage at the client: After creating the authorization request URI, the client can non-deterministically decide whether to leak it. In the leak case, it non-deterministically chooses an IP address and sends out the request to this IP (in plain). Thus, the client model leaks the client identifier and the request URI value of the authorization request.

*A5* (resource requests leak) was initially modeled by leaking the resource request at the RS after receiving the request (and in the case of opaque tokens, before receiving the introspection response). As described in Section III-C, this enables the attacker to replay DPoP proofs. In addition to using server-provided DPoP nonces, we adapted the model (in line with changes to the FAPI 2.0 Attacker Model, see Section III-C) such that the resource request leaks *after* the RS responds with the resource response (i.e., after invalidating the DPoP nonce).

3) *Limitations of Our Model:* While our model covers most of the FAPI 2.0 specifications, there are a few things which we handle on a very abstract level or do not model at all (besides the inherent abstractions of the WIM, such as details of TLS).

**Error Handling.** FAPI 2.0 and several of the underlying specifications define a set of error messages, e.g., when client authentication fails. These are not represented in our model: if a process encounters an error condition, it just aborts the current processing step without output and without changes to the process state. We note that none of the specifications mandate a certain behavior upon *receiving* an error message.

**Rich Authorization Requests.** *Rich Authorization Requests (RAR)* [56] is a mechanism which allows clients to be more precise about what kind of access they request. However, FAPI 2.0 does not mandate the use of RAR, and RAR semantics are subject to individual AS policies. Hence, RAR does not allow for a general treatment, and hence, is not covered by our model.

**Modeled Grant Types.** The OAuth 2.0 framework [41], as well as OpenID Connect [70], define various grant types, such as the implicit grant, the hybrid grant, and the authorization code grant. FAPI 2.0 explicitly excludes all of them, except for the client credentials grant and the authorization code grant, where only the latter is required to be supported by FAPI 2.0 ASs and clients. In the client credentials grant, the client takes the role of the user, hence, removing the additional interaction between user and client. So, this grant is subsumed by the authorization code grant, which is why we only model and analyze the latter.

**Refresh Tokens.** Even though FAPI 2.0 allows for refresh tokens (see [41]) to be used, we do not model them. Instead, we model regular access (and id) tokens as having an indefinite lifetime, i.e., they never expire.

**Native Clients with Loopback Redirect.** FAPI 2.0 ASs must reject PAR requests with a redirect uri using the http scheme, i.e., where the client's redirection endpoint is not protected by TLS (cf. Step 5 in Figure 1). There is, however, one exception to this rule: if the client is a native client running on the same device as the browser and is using loopback interface (i.e., device-local) redirection, it may use an http redirect uri (see [18] and [31, Sec. 4.3.1.2 No. 8]). In our AS model, we do not allow http redirect uris.

### C. Comparison to WIM Analyses of Related Protocols

We here describe specific differences between our work and previous WIM analyses of FAPI 1.0 [38], OAuth 2.0 [36], and OpenID Connect [37].

Since FAPI 2.0 is a completely redesigned protocol, obviously, previous analyses do not apply to FAPI 2.0. The differences are also apparent by the new vulnerabilities and attacks we have found on FAPI 2.0.

While FAPI 2.0 is based on OAuth 2.0 and OpenID Connect, it not only contains many additional mechanisms, such as PAR, DPoP, mTLS, PKCE, Authorization Server Issuer Identification, AS metadata and so on, but also aims to be secure under a much stronger attacker, see also below.

Similarly, while FAPI 1.0, the predecessor of FAPI 2.0, was formally analyzed [38], FAPI 2.0 is a very different protocol: in comparison to FAPI 1.0, several security mechanisms have been removed, most notably *OAuth Token Binding* [43] and *JWT Secured Authorization Response Mode (JARM)* [53]. On the other hand, FAPI 2.0 adds new mechanisms, such as DPoP and PAR, which have not undergone any formal treatment so far. Also, FAPI 2.0 mandates client authentication, removes support for public clients as well as the hybrid flow, and introduces a stronger attacker model.

In the following, we briefly discuss further details and differences.

**Token Introspection.** Previous work on OAuth 2.0 and OpenID Connect does not consider the resource server, and while prior work on FAPI 1.0 does, token introspection is modeled only on a very abstract level: instead of querying the AS, the RS uses an idealized introspection oracle. In contrast, our AS and RS models contain a detailed representation of token introspection as defined in RFC 7662 [69]. Hence, we have to reason about additional messages, their contents, as well as integrity, authentication, and secrecy thereof.

**Structured Access Tokens.** Prior work on OAuth 2.0, OpenID Connect, and FAPI 1.0 does not contain a model of structured access tokens and their use, in particular at the RS. Note that besides the additional options for flows introduced by structured access tokens, modeling them also requires additional keys to be distributed.

**Number of ASs Supported by an RS.** As mentioned above, prior work on OAuth 2.0 and OpenID Connect did not model the RS. In the FAPI 1.0 analysis in [38], the RS supports one fixed AS, whereas our RS model supports an arbitrary set of ASs, some of which may be corrupted. Hence, our analysis considers a broader set of cases and possible mix-ups.

**RS Authentication.** RFC 7662 [69] mandates for RSs to authenticate to the AS during token introspection. Since previous work on OAuth 2.0, OpenID Connect, and FAPI 1.0 did not model token introspection, our model is the first (not just among WIM analyses) incorporating RS authentication. Similar to structured access tokens above, this requires distribution of, and reasoning about, additional secrets.

**DPoP Access Token Sender Constraining.** The DPoP [39] mechanism has only recently been standardized and hence, has not been part of any formal analysis to the best of our knowledge. With its additional signatures and subtle details regarding the signed data, in particular, the exact contents of the `htu` value (see [39, Sec. 4.2]), DPoP requires careful modeling of the AS, the RS, and the client. In addition, DPoP requires key material to be distributed between client, AS, and RS (as well as additions to the token introspection and structured access token models).

**PAR.** Similar to DPoP above, PAR [57] has only recently been standardized and has thus not been part of any formal analysis. While PAR does not require additional key material, it adds more options and messages to the FAPI 2.0 flows, hence necessitating reasoning about them.

**JWT Client Authentication.** While the `public_key_jwt` client authentication method has been standardized as part of OpenID Connect [70], it has not received any formal security analysis so far, including WIM analyses of OpenID Connect and FAPI 1.0 (FAPI 1.0 is based on OpenID Connect and explicitly allows for `public_key_jwt` client authentication). Hence, our work is the first to consider this kind of client authentication.

**Attacker Model.** Compared to prior WIM work on OAuth 2.0, OpenID Connect, and FAPI 1.0, our analysis considers a significantly stronger attacker model. Where [36] and [37] consider a fairly standard Dolev-Yao attacker (i.e., similar to  $\mathcal{A}_1$ , see Section II-C), [38] considers a stronger attacker: the attacker in [38] is a network attacker similar to our  $\mathcal{A}_1$ , which also may access authorization requests (cf.  $\mathcal{A}_2$ ) and authorization responses (cf.  $\mathcal{A}_3$ ) in plain, as well as force honest clients to use attacker-controlled token endpoints (cf.  $\mathcal{A}_4$ ). In addition, the authors of [38] consider access tokens leaking to the attacker, which is something that is subsumed by our attacker assumption  $\mathcal{A}_5$  (resource request leak). However, our attacker model is even stronger in that (1)  $\mathcal{A}_5$  not only leaks the access token, but also all other values in the resource request, in particular DPoP proofs, (2)  $\mathcal{A}_6$  leaks resource responses, and (3)  $\mathcal{A}_7$  even allows for the attacker to modify resource responses. To the best of our knowledge, neither (1), nor (2), nor (3) have previously been considered in any formal analysis of authorization protocols.

**Session Integrity Properties.** While prior works on OAuth 2.0 and OpenID Connect do consider session integrity properties similar to ours, these are only proven in a setting with *web attackers*, i.e., the attacker may corrupt any party, but does not control the network and cannot spoof sender addresses. Similar properties have also been proven for FAPI 1.0 *with* a network attacker, however, their property requires the use of *OAuth 2.0 Token Binding* [43], a mechanism to bind authorization codes and access tokens to TLS connections, which is quite different from DPoP and mTLS access token sender constraining.

#### D. FAPI 2.0 Web System

As outlined in Section IV-A, a web system formalizes the overall system covered by our analysis. We call  $(\mathcal{W}, \mathcal{S}, \text{script}, E^0)$  a *FAPI 2.0 web system with network attacker*, or short  $\mathcal{F}API$ , if the components of the web system are defined as follows:

$\mathcal{W} = \text{Hon} \cup \text{Net}$  consists of a network attacker process (in `Net`), a finite set `B` of web browsers, a finite set `C` of web servers for the clients, a finite set `AS` of web servers for the authorization servers, and a finite set `RS` of web servers for the resource servers, with `Hon` := `B`  $\cup$  `C`  $\cup$  `AS`  $\cup$  `RS`. The honest processes behave as outlined in Section IV-A and Section IV-B (formally specified in Appendix D). DNS servers are subsumed by the network attacker and are therefore not modeled explicitly.

$\mathcal{S}$  contains the attacker script  $R^{\text{att}}$  outlined in Section IV-A (formally specified in Definition 55), the script `script_as_form` for logging in the resource owner at the AS, and the script `script_client_index` for selecting an AS and initiating a flow at the client as outlined in Appendix A (both scripts are formally specified in Appendix D). The mapping script maps these scripts to strings: `script_client_index`  $\mapsto$  `script_client_index`, `script_as_form`  $\mapsto$  `script_as_form`,  $R^{\text{att}}$   $\mapsto$  `att_script`.

$E^0$  is the initial set of events and contains only trigger events (see Section IV-A).

Note that we prove security properties about all such FAPI 2.0 web systems with network attacker, and hence, systems with an arbitrary number of browsers, clients, ASs, and RSs, and in which each party can have an arbitrary number of parallel, interleaving protocol sessions. Also recall from Section IV-A that the attacker can corrupt honest parties at any time.

## E. Formal Security Properties

As described in [Section II-C](#), the FAPI WG wants the protocol to meet authorization, authentication, and session integrity goals. In the following, we present our formalized authorization and authentication properties, capturing the corresponding security goals. We present the formalized session integrity properties in [Appendix E](#). We conclude this section with our main security theorem.

1) *Authorization*: Recall that informally, authorization means that an attacker should never be able to access resources of honest users (unless the user authorized such access). We highlight that this statement covers many different scenarios, for example, that the attacker cannot use leaked access tokens at the RS and cannot, by some mix-up, force an honest client to use an access token associated with an honest user in a session with the attacker.

We first formalize that an access token was issued by AS  $AS$ , is bound to a public key  $k$ , and is associated with a user identity  $id$ , and give our authorization property afterwards.

*Definition 1 (Access Token bound to Key, AS and ID)*. Let  $\mathcal{F}API$  be a FAPI 2.0 web system with network attacker,  $k \in \mathcal{T}_{\mathcal{N}}$  be a term,  $AS \in \mathcal{AS}$  an AS, and  $id \in \mathcal{ID}$  a user identity.<sup>2</sup> We say that a term  $t$  is an *access token bound to  $k$ ,  $AS$ , and  $id$*  in configuration  $(S, E, N)$  of a run of  $\mathcal{F}API$ , if there is an entry  $rec \in \mathcal{S}(AS).records$  (i.e., in the state of  $AS$ ) such that  $rec[access\_token] \equiv t$  and  $rec[subject] \equiv id$  and  $(rec[cnf] \equiv [jkt: hash(k)]) \vee (rec[cnf] \equiv [x5t\#S256: hash(k)])$ .

Informally, this means that  $t$  is stored in the state of  $AS$ , together with the identity  $id$  and the hash of  $k$ . If the key is stored under the name  $jkt$ , then the token is bound via DPoP, otherwise, it is bound via mTLS.

The following authorization property captures the following: if an honest RS  $rs$  provides access to a resource  $r$  of an honest resource owner with user identity  $id$  managed by an honest AS  $AS$ , then the following holds true: (I)  $rs$  has received a request for accessing the resource  $r$  with an access token  $at$  in the same (which is possible if the token  $at$  is structured and can be verified by the RS immediately) or in a previous processing step (if the token  $at$  is opaque to the RS and it thus performed token introspection), and  $rs$  created the resource when receiving the resource request (see [Appendix A-C](#) on how our model manages resources). (II) The token  $at$  is bound to some key  $k$ ,  $AS$ , and the user identity  $id$  (see [Definition 1](#)). (III) If  $k$  is the key of an honest client, then the attacker cannot derive the resource.

*Definition 2 (Authorization Property)*. We say that a FAPI 2.0 web system with network attacker  $\mathcal{F}API$  is *secure w.r.t. authorization* iff for every run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of  $\mathcal{F}API$ , every RS  $rs \in \mathcal{RS}$  that is honest in  $S^n$ , every identity  $id \in \mathcal{S}_0^{rs}.ids$  with  $b = \text{ownerOfID}(id)$  being an honest browser in  $S^n$ , every processing step  $Q = (S^Q, E^Q, N^Q) \xrightarrow[rs \rightarrow E_{out}^Q]{e_{in}^Q \rightarrow rs}$   $(S^{Q'}, E^{Q'}, N^{Q'})$  in  $\rho$ , every  $resourceID \in \mathcal{S}$  with  $AS = \text{authorizationServerOfResource}^{rs}(resourceID)$  being honest in  $S^Q$ , it holds true that:<sup>3</sup>

If  $\exists r, x, y, k, m_{resp}. \langle x, y, \text{enc}_s(m_{resp}, k) \rangle \in \mathcal{E}_{out}^Q$  such that  $m_{resp}$  is an HTTP response,  $r := m_{resp}.body[resource]$ , and  $r \in \mathcal{S}^{Q'}(rs).resourceNonce[id][resourceID]$ , then

- (I)  $\exists$  a processing step  $P = s_i \xrightarrow[rs \rightarrow E_{out}^P]{e_{in}^P \rightarrow rs} s_{i+1}$  such that
  - a) either  $P = Q$ , or  $P$  is prior to  $Q$  in  $\rho$ , and
  - b)  $e_{in}^P$  is an event  $\langle x, y, \text{enc}_a(\langle m_{req}, k_1 \rangle, k_2) \rangle$  for some  $x, y, k_1$ , and  $k_2$  where  $m_{req}$  is an HTTP request which contains a term (access token)  $t$  in its Authorization header, i.e.,  $t \equiv m_{req}.headers[Authorization].2$ , and
  - c)  $r$  is a fresh nonce generated in  $P$  at the resource endpoint of  $rs$  in [Line 46](#) of [Algorithm 12](#).
- (II)  $t$  is bound to a key  $k \in \mathcal{T}_{\mathcal{N}}$ ,  $AS$ , and  $id$  in  $S^Q$  (see [Definition 1](#)).
- (III) If there exists a client  $c \in \mathcal{C}$  such that  $k \equiv \text{pub}(\text{signkey}(c))$  or  $k \equiv \text{pub}(\text{tlskey}(d_c))$  and  $d_c \in \text{dom}(c)$ , and if  $c$  is honest in  $S^n$ , then  $r$  is not derivable from the attackers knowledge in  $S^n$ , i.e.,  $r \notin d_0(S^n(\text{attacker}))$ .

2) *Authentication*: Recall that the authentication goal states that an attacker should not be able to log in at an honest client under the identity of an honest user. In our model, the client sets a cookie that we call *service session id* at the browser after a successful login. The client model stores the service session id in its `sessions` state subterm, and associates with it the identity that is logged in to the session (the identity is taken from an ID token, see [Appendix A-B](#) for more details on the login of a user at a client). On a high level, our formalized property states that an attacker should not be able to derive the service session id for a session at an honest client where an honest identity is logged in, as long as the identity is managed by an honest AS. We stress that this not only covers that a cookie set at the browser of the honest user does not leak, but that there is no way in which the attacker can log in at an honest client as an honest user.

<sup>2</sup>ID is a set of terms of the form  $\langle name, domain \rangle$ , where  $name$  is a string, the user name, and  $domain$  is a domain (usually of an AS). We also define a mapping  $\text{ownerOfID}: \mathcal{ID} \rightarrow \mathcal{B}$  which maps an identity to the browser whose user owns the identity (users are modeled as part of their browser). Likewise, we define a mapping  $\text{governor}: \mathcal{ID} \rightarrow \mathcal{AS}$  mapping identities to “their” AS. See [Appendix D-B](#) and [Appendix D-D](#) for details.

<sup>3</sup> $\text{authorizationServerOfResource}^{rs}$  is a mapping from resource ids to the authorization server that manages the respective resource, see [Definition 10](#).

We start with an auxiliary definition, capturing that the client logged in a user with a service session id, before presenting the authentication property itself.

*Definition 3 (Service Sessions).* We say that there is a *service session identified by a nonce  $n$  for a user identity  $id$  at some client  $c$*  in a configuration  $(S, E, N)$  of a run  $\rho$  of a FAPI 2.0 web system with network attacker  $\mathcal{FAPI}$  iff there exists some session id  $x$  and a domain  $d \in \text{dom}(\text{governor}(id))$  such that

$$S(c).\text{sessions}[x][\text{loggedInAs}] \equiv \langle d, id \rangle$$

and

$$S(c).\text{sessions}[x][\text{serviceSessionId}] \equiv n.$$

*Definition 4 (Authentication Property).* We say that a FAPI 2.0 web system with network attacker  $\mathcal{FAPI}$  is *secure w.r.t. authentication* iff for every run  $\rho$  of  $\mathcal{FAPI}$ , every configuration  $(S, E, N)$  in  $\rho$ , every  $c \in C$  that is honest in  $S$ , every identity  $id \in \text{ID}$  with  $\text{AS} = \text{governor}(id)$  being an honest AS (in  $S$ ) and with  $b = \text{ownerOfID}(id)$  being an honest browser in  $S$ , every service session identified by some nonce  $n$  for  $id$  at  $c$ ,  $n$  is not derivable from the attacker's knowledge in  $S$  (i.e.,  $n \notin d_\emptyset(S(\text{attacker}))$ ).

3) *Security Theorem:* As described in [Section II-C](#), the protocol aims to fulfill authorization, authentication, and session integrity properties. Thus, our overall security theorem is the conjunction of all four properties, where as mentioned session integrity for authorization and authentication is defined in [Appendix E](#).

*Theorem 1.* Every FAPI 2.0 web system with network attacker  $\mathcal{FAPI}$  fulfills authorization, authentication, session integrity for authentication, and session integrity for authorization.

We highlight that we prove this theorem for the powerful attacker described in [Section IV-B2](#) within a faithful formal model that includes the fixes described in [Section III](#). We also emphasize that our analysis takes into account many Web features that can be the root of attacks: e.g., the browser model allows for the execution of scripts loaded from different websites/origins at the same time, possibly with malicious scripts. The model also considers fine-grained behavior of HTTP redirects,<sup>4</sup> several security-critical headers, as well as subtleties of various cookie attributes, which, for example, could result in vulnerable session management, and in-browser communication using `postMessages`, just to name a few of the Web features considered in our analysis. Thus, our analysis excludes attacks that arise from features of the Web infrastructure. Our proof of [Theorem 1](#), which, due to space limitations, we give in [Appendix F](#), consists of more than 20 lemmas and of course reasons about the full formal model that we provide in [Appendix D](#).

#### F. Proof Sketch for Authorization

In this section, we summarize our proof of the authorization property. We refer to [Appendix F](#) for the full proof of this as well as the authentication and session integrity properties.

Recall that on a high level, the authorization property requires that the attacker cannot access resources of an honest identity  $id$  stored at an honest RS.

For the **first postcondition**, see [Definition 2 \(I\)](#), we show that whenever the RS gives access to a resource of an identity (the resource owner), it does so in response to a request that contains an access token. We show this for both types of access tokens, i.e., for structured tokens signed by the AS managing the resource, as well as opaque tokens that the RS has to verify with the AS via token introspection. Unlike the other two, this postcondition is quite easy to show.

For the **second postcondition**, see [Definition 2 \(II\)](#), we show that the access token was issued by the AS managing the resource and is bound to some key  $k$  and the identity  $id$  of the resource owner (see [Definition 1](#)). Keep in mind that [Definition 1](#) refers to the state of AS. For structured access tokens, we show this by reasoning on the validity of the token: The RS checks that the token has a valid signature by the AS managing the resource, i.e., the token was indeed created by that AS (we give more details in the full proof, e.g., proving that the signing key of AS is secret). We show that the AS, when creating the token, stores the token in its state, together with  $id$  and  $k$  that it also puts into the structured token (i.e., binds the token to  $k$  and  $id$ ). For opaque tokens, the introspection endpoint of AS returns the identity  $id$  and the key  $k$  (the RS only provides access in case of a successful introspection response). The AS only sends such a response if it previously bound the token to  $k$  and  $id$ .

The **third postcondition**, see [Definition 2 \(III\)](#), states that if the key  $k$  is the private (TLS or DPoP) key of an honest client  $c$ , then the attacker does not get access to the resource. Showing this property is more involved than the previous ones.

We first show that the resource request was indeed sent by  $c$ : If the RS received a DPoP proof (for the key that the token is bound to), then we can apply [Lemma 12](#) which gives us that the DPoP proof can only be derived by  $c$  (and  $rs$ , which just received it). Note that this only holds true with the fix that we propose in [Section III-C](#) (i.e., effective replay protection) and with the updated [A5](#) attacker assumption (see [Section IV-B2](#)). Similarly, we show the effectiveness of mTLS token binding, i.e., that access tokens bound to a TLS key of an honest client can only be used by that client. Overall, we show for both structured and opaque access

<sup>4</sup>For example, FAPI 2.0 excludes code 307 redirects, as they would cause attacks similar to [\[36\]](#)

tokens that the resource request must have been sent by  $c$ . Thus, the RS provides  $c$  access to the resource of  $id$  (recall that all communication between parties happens via secure channels).

We then prove that  $c$  does not provide the attacker access to the resource, i.e., the attacker cannot use  $c$  (e.g., as a user) to access resources of  $id$ . We show this by contradiction, i.e., we assume that  $c$  provides the attacker access, and lead this to a contradiction. In the full proof, we show that this can only happen if the redirection request (Step 12 in Figure 1) was sent by the attacker (recall that  $c$  expects a request to its redirection endpoint with an authorization code  $code$  that  $c$  then uses at the token endpoint of the AS). We then show that  $c$  sends the token request to the correct AS, i.e., to the one that is managing the resource. This is true thanks to the fix for preventing the Cuckoo’s Token attack (see Section IV-B1). As the access token that  $c$  receives from the token endpoint provides access to resources of  $id$ , it follows that  $code$  is also associated with  $id$  (i.e., the AS stored in its state that  $code$  was issued when  $id$  logged in at its authorization endpoint). Lemma 5 now gives us that  $c$  received this exact  $code$  at its redirection endpoint. As the request to the redirection endpoint is assumed to originate from the attacker, we have that the attacker can derive  $code$ . However, this is a contradiction to Lemma 17, which states that the authorization code of a flow between an honest AS and honest client is not derivable by any other process, including the attacker. This concludes the proof.

## V. RELATED WORK

While there has not been much work on the relatively recent and partly still under development FAPI standards and specifications so far, standards like OAuth 2.0 and OpenID Connect, on which FAPI is based, have received quite some attention by security researchers. Besides many studies on implementations and deployments [77–79, 82], OAuth 2.0 has been formally analyzed: Pai et al. [65] built a limited model of OAuth 2.0, lacking many generic web features, for the Alloy finite-state model checker and showed that with their approach, known weaknesses can be found. Chari et al. [14] analyzed the authorization code flow in the UC model and found no attacks, but their model omits many web features. Two more comprehensive formal analyses have been conducted by Bansal et al. [5, 6], using their WebSpi library and ProVerif, in which they modeled various settings of OAuth 2.0, e.g., with CSRF vulnerabilities in ASs and clients. Bansal et al. uncovered several previously unknown attacks on various popular OAuth 2.0 implementations. However, their works main focus was on finding attacks, rather than proving security. The latter was the focus of a formal analysis by Fett et al. [36], in which, despite prior formal analysis efforts, they discovered several new attacks, and were only able to prove security after developing fixes for them. As in our work, the analysis by Fett et al. is based on the WIM, i.e., includes a comprehensive formal model of the web infrastructure. Likewise, OpenID Connect has been analyzed before [37, 52, 58].

However, as described in more detail in Section IV-C, FAPI 2.0 is a completely redesigned protocol, and hence, previous analyses on related protocols do not apply. Finally, as already mentioned, FAPI 2.0 itself has not undergone any formal security analysis so far.

## VI. CONCLUSION

Asked by the OpenID Foundation’s FAPI working group, we accompanied the development of the FAPI 2.0 specifications with a formal security analysis, by creating formal models and formalizing security properties closely following the FAPI 2.0 specifications, including their strong attacker model. During this formal analysis, we discovered several attacks that violate the security goals set by the FAPI 2.0 specifications. These have been reported to the FAPI WG, who acknowledged our attacks. We then worked with them to incorporate many of our proposed fixes and improvements into the official specifications. This finally allowed us to provide formal security proofs of all security properties.

Performing such an analysis in a meaningful model like the WIM ensures that attacks based on many threats and features of the web and their complex interactions can be found and ruled out. We note that our analysis is a pen-and-paper analysis, as current tools do not directly support models that are as detailed as the WIM. Hence, our proofs are not mechanized and, given the complexity of the protocol, are inherently lengthy, hence, tedious to verify. Nevertheless, our formal and systematic analysis has helped improve the standard and gain more insights and confidence in its security. Conversely, we consider mechanized analyses of web protocols as interesting future work.

Surprisingly, besides new vulnerabilities and attacks, we also uncovered some attacks which were similar to known attacks on related protocols, and in particular known within the FAPI WG. These findings further underline the importance of a systematic, formal analysis, as for complex protocols, like FAPI 2.0, also experts easily overlook even known attack patterns.

Overall, our contributions to the standardization process were very welcomed by the FAPI WG and led to significant improvements of an important protocol just in time: FAPI 2.0 is about to be adopted in highly sensitive environments, with millions of users managing bank account data, financial transactions, eGovernment applications, and even health data.

## ACKNOWLEDGMENT

We thank our anonymous reviewers for their invaluable feedback. This research was funded in part by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – 443324941, and by the OpenID Foundation.



## REFERENCES

- [1] *A05:2021 – Security Misconfiguration*. URL: [owasp.org/Top10/A05\\_2021-Security\\_Misconfiguration/](https://owasp.org/Top10/A05_2021-Security_Misconfiguration/).
- [2] M. Abadi and C. Fournet. “Mobile Values, New Names, and Secure Communication”. In: *ACM POPL*. 2001, pp. 104–115.
- [3] *Authorization Request Leaks lead to CSRF*. URL: [bitbucket.org/openid/fapi/issues/534](https://bitbucket.org/openid/fapi/issues/534).
- [4] *Bank Account Aggregation Services Australia*. URL: [www.creditsense.com.au/bank-account-aggregation/](http://www.creditsense.com.au/bank-account-aggregation/).
- [5] C. Bansal et al. “Discovering Concrete Attacks on Website Authorization by Formal Analysis”. In: *IEEE CSF*. 2012.
- [6] C. Bansal et al. “Discovering concrete attacks on website authorization by formal analysis”. In: *Journal of Computer Security* 22.4 (2014), pp. 601–657.
- [7] R. Berjon et al., eds. *HTML5, W3C Recommendation*. 2014. URL: [www.w3.org/TR/html5/](http://www.w3.org/TR/html5/).
- [8] V. Bertocci. *JSON Web Token (JWT) Profile for OAuth 2.0 Access Tokens*. RFC 9068.
- [9] *Browser swap attack explained on 2022-09-28*. URL: [bitbucket.org/openid/fapi/issues/543](https://bitbucket.org/openid/fapi/issues/543).
- [10] *Budget Tracker & Planner*. URL: [mint.intuit.com/](https://mint.intuit.com/).
- [11] B. Campbell et al. *OAuth 2.0 Mutual-TLS Client Authentication and Certificate-Bound Access Tokens*. RFC 8705.
- [12] *Car API platform for connected vehicle data*. URL: [smartcar.com/](https://smartcar.com/).
- [13] *Change attacker model to reflect formal model*. URL: [bitbucket.org/openid/fapi/pull-requests/381](https://bitbucket.org/openid/fapi/pull-requests/381).
- [14] S. Chari et al. *Universally Composable Security Analysis of OAuth v2.0*. Cryptology ePrint Archive, Paper 2011/526. 2011. URL: [eprint.iacr.org/2011/526](https://eprint.iacr.org/2011/526).
- [15] L. Chen et al. *Cookies: HTTP State Management Mechanism*. IETF Draft draft-ietf-httpbis-rfc6265bis-09. 2021.
- [16] *Consumer Data Standards*. URL: [consumerdatastandards.gov.au/](https://consumerdatastandards.gov.au/).
- [17] *Decide on what to do for A. Cuckoo’s Token Attack*. URL: [bitbucket.org/openid/fapi/issues/525](https://bitbucket.org/openid/fapi/issues/525).
- [18] W. Denniss and J. Bradley. *OAuth 2.0 for Native Apps*. RFC 8252.
- [19] *Die elektronische Patientenakte*. URL: [www.bundesgesundheitsministerium.de/elektronische-patientenakte.html](https://www.bundesgesundheitsministerium.de/elektronische-patientenakte.html).
- [20] *Discovery should be mandated for clients*. URL: [bitbucket.org/openid/fapi/issues/536](https://bitbucket.org/openid/fapi/issues/536).
- [21] Q. H. Do et al. “A Formal Security Analysis of the W3C Web Payment APIs: Attacks and Verification”. In: *IEEE S&P*. 2022, pp. 215–234.
- [22] *DPoP & resource leaks*. URL: [bitbucket.org/openid/fapi/issues/533](https://bitbucket.org/openid/fapi/issues/533).
- [23] *Electronic Prescription Service - FHIR API*. URL: [digital.nhs.uk/developer/api-catalogue/electronic-prescription-service-fhir](https://digital.nhs.uk/developer/api-catalogue/electronic-prescription-service-fhir).
- [24] *email – Graph API*. URL: [developers.facebook.com/docs/permissions/reference/email](https://developers.facebook.com/docs/permissions/reference/email).
- [25] *FAPI2SP: Add requirement for RP to use discovery*. URL: [bitbucket.org/openid/fapi/pull-requests/363](https://bitbucket.org/openid/fapi/pull-requests/363).
- [26] *FAPI2SP: Add security consideration for cuckoo’s token attack*. URL: [bitbucket.org/openid/fapi/pull-requests/364](https://bitbucket.org/openid/fapi/pull-requests/364).
- [27] *FDX*. URL: [financialdataexchange.org/](https://financialdataexchange.org/).
- [28] D. Fett. “An Expressive Formal Model of the Web Infrastructure”. PhD thesis. 2018.
- [29] D. Fett. *FAPI 2.0 Attacker Model. 2nd Implementer’s Draft*. 2022. URL: [openid.net/specs/fapi-2\\_0-attacker-model-02.html](https://openid.net/specs/fapi-2_0-attacker-model-02.html).
- [30] D. Fett. *FAPI 2.0 Attacker Model, Commit 209f58a*. OpenID Foundation. 2022. URL: [bitbucket.org/openid/fapi/src/209f58a/FAPI\\_2\\_0\\_Attacker\\_Model.md](https://bitbucket.org/openid/fapi/src/209f58a/FAPI_2_0_Attacker_Model.md).
- [31] D. Fett. *FAPI 2.0 Baseline Profile, Commit 209f58a*. OpenID Foundation. 2022. URL: [bitbucket.org/openid/fapi/src/209f58a/FAPI\\_2\\_0\\_Baseline\\_Profile.md](https://bitbucket.org/openid/fapi/src/209f58a/FAPI_2_0_Baseline_Profile.md).
- [32] D. Fett. *FAPI 2.0 Security Profile. 2nd Implementer’s Draft*. 2022. URL: [openid.net/specs/fapi-2\\_0-security-02.html](https://openid.net/specs/fapi-2_0-security-02.html).
- [33] D. Fett et al. “An Expressive Model for the Web Infrastructure: Definition and Application to the BrowserID SSO System”. In: *IEEE S&P*. 2014, pp. 673–688.
- [34] D. Fett et al. “Analyzing the BrowserID SSO System with Primary Identity Providers Using an Expressive Model of the Web”. In: *ESORICS*. 2015, pp. 43–65.
- [35] D. Fett et al. “SPRESSO: A Secure, Privacy-Respecting Single Sign-On System for the Web”. In: *ACM CCS*. 2015.
- [36] D. Fett et al. “A Comprehensive Formal Security Analysis of OAuth 2.0”. In: *ACM CCS*. 2016.
- [37] D. Fett et al. “The Web SSO Standard OpenID Connect: In-Depth Formal Security Analysis and Security Guidelines”. In: *IEEE CSF*. 2017.
- [38] D. Fett et al. “An Extensive Formal Security Analysis of the OpenID Financial-grade API”. In: *IEEE S&P*. 2019, pp. 1054–1072.
- [39] D. Fett et al. *OAuth 2.0 Demonstrating Proof-of-Possession at the Application Layer (DPoP)*. IETF Draft draft-ietf-oauth-dpop-08. 2022.
- [40] *Governance of NZ payment systems*. URL: [www.paymentsnz.co.nz/](https://www.paymentsnz.co.nz/).
- [41] D. Hardt. *The OAuth 2.0 Authorization Framework*. RFC 6749.
- [42] *Improve attacker model description after introduction of metadata*. URL: [bitbucket.org/openid/fapi/pull-requests/371](https://bitbucket.org/openid/fapi/pull-requests/371).
- [43] M. Jones et al. *OAuth 2.0 Token Binding*. IETF Draft draft-ietf-oauth-token-binding-08. 2018.

- [44] M. Jones et al. *Assertion Framework for OAuth 2.0 Client Authentication and Authorization Grants*. RFC 7521.
- [45] M. Jones et al. *JSON Web Signature (JWS)*. RFC 7515.
- [46] M. Jones et al. *JSON Web Token (JWT)*. RFC 7519.
- [47] M. Jones et al. *JSON Web Token (JWT) Profile for OAuth 2.0 Client Authentication and Authorization Grants*. RFC 7523.
- [48] M. Jones et al. *OAuth 2.0 Authorization Server Metadata*. RFC 8414.
- [49] A. Kumar. “Using automated model analysis for reasoning about security of web protocols”. In: *ACM ACSAC*. 2012.
- [50] A. Kumar. “A Lightweight Formal Approach for Analyzing Security of Web Protocols”. In: *Research in Attacks, Intrusions and Defenses*. Springer, 2014, pp. 192–211.
- [51] R. Küsters et al. *The Web Infrastructure Model (WIM)*. Technical Report. Version 1.0. 2022. URL: [www.sec.uni-stuttgart.de/research/wim/WIM\\_V1.0.pdf](http://www.sec.uni-stuttgart.de/research/wim/WIM_V1.0.pdf).
- [52] W. Li and C. J. Mitchell. “Analysing the Security of Google’s Implementation of OpenID Connect”. In: *Detection of Intrusions and Malware, and Vulnerability Assessment*. Springer, 2016, pp. 357–376.
- [53] T. Lodderstedt and B. Campbell. *Financial-grade API: JWT Secured Authorization Response Mode for OAuth 2.0 (JARM)*. OpenID Foundation. 2018. URL: [openid.net/specs/openid-financial-api-jarm.html](http://openid.net/specs/openid-financial-api-jarm.html).
- [54] T. Lodderstedt et al. *IETF 105 Materials: OAuth 2.0 Security Best Current Practice*. 2019. URL: [datatracker.ietf.org/meeting/105/materials/slides-105-oauth-sessa-oauth-security-topics-00.pdf](http://datatracker.ietf.org/meeting/105/materials/slides-105-oauth-sessa-oauth-security-topics-00.pdf).
- [55] T. Lodderstedt et al. *OAuth 2.0 Security Best Current Practice*. IETF Draft ietf-oauth-security-topics-19. 2021.
- [56] T. Lodderstedt et al. *OAuth 2.0 Rich Authorization Requests*. IETF Draft draft-ietf-oauth-rar-15. 2022.
- [57] T. Lodderstedt et al. *OAuth 2.0 Pushed Authorization Requests*. RFC 9126.
- [58] V. Mladenov et al. “On the security of modern Single Sign-On Protocols: Second-Order Vulnerabilities in OpenID Connect”. In: *CoRR* abs/1508.04324v2 (2015).
- [59] *OAuth 2.0 Scopes for Google APIs*. URL: [developers.google.com/identity/protocols/oauth2/scopes#drive](https://developers.google.com/identity/protocols/oauth2/scopes#drive).
- [60] *OAuth 2.0 Simplified: Token Introspection Endpoint*. URL: [www.oauth.com/oauth2-servers/token-introspection-endpoint/](http://www.oauth.com/oauth2-servers/token-introspection-endpoint/).
- [61] *OAuth Guide*. URL: [developers.dropbox.com/oauth-guide](https://developers.dropbox.com/oauth-guide).
- [62] *Open Banking UK*. URL: [www.openbanking.org.uk/](http://www.openbanking.org.uk/).
- [63] *Open Finance*. URL: [www.bcb.gov.br/en/financialstability/open\\_finance](http://www.bcb.gov.br/en/financialstability/open_finance).
- [64] *Open Finance*. URL: [www.gov.br/susep/pt-br/assuntos/open-insurance](http://www.gov.br/susep/pt-br/assuntos/open-insurance).
- [65] S. Pai et al. “Formal Verification of OAuth 2.0 Using Alloy Framework”. In: *IEEE CSNT*. 2011.
- [66] *Reduced attacker model*. URL: [bitbucket.org/openid/fapi/pull-requests/377](https://bitbucket.org/openid/fapi/pull-requests/377).
- [67] J. Reschke. *The ‘Basic’ HTTP Authentication Scheme*. RFC 7617.
- [68] E. Rescorla and T. Dierks. *The Transport Layer Security (TLS) Protocol Version 1.2*. RFC 5246.
- [69] J. Richer. *OAuth 2.0 Token Introspection*. RFC 7662.
- [70] N. Sakimura et al. *OpenID Connect Core 1.0 incorporating errata set 1*. OpenID Foundation. 2014. URL: [openid.net/specs/openid-connect-core-1\\_0.html](http://openid.net/specs/openid-connect-core-1_0.html).
- [71] N. Sakimura et al. *OpenID Connect Discovery 1.0 incorporating errata set 1*. OpenID Foundation. 2014. URL: [openid.net/specs/openid-connect-discovery-1\\_0.html](http://openid.net/specs/openid-connect-discovery-1_0.html).
- [72] N. Sakimura et al. *Proof Key for Code Exchange by OAuth Public Clients*. RFC 7636.
- [73] N. Sakimura et al. *The OAuth 2.0 Authorization Framework: JWT-Secured Authorization Request (JAR)*. RFC 9101.
- [74] *Scopes for OAuth*. URL: [docs.github.com/en/developers/apps/building-oauth-apps/scopes-for-oauth-apps](https://docs.github.com/en/developers/apps/building-oauth-apps/scopes-for-oauth-apps).
- [75] K. M. zu Selhausen and D. Fett. *OAuth 2.0 Authorization Server Issuer Identification*. RFC 9207.
- [76] *Services provided at Helsenorge*. URL: [www.helsenorge.no/en/about-services-on-helsenorge/how-to-use-helsenorge-services#services-provided-at-helsenorge](http://www.helsenorge.no/en/about-services-on-helsenorge/how-to-use-helsenorge-services#services-provided-at-helsenorge).
- [77] E. Shernan et al. “More Guidelines Than Rules: CSRF Vulnerabilities from Noncompliant OAuth 2.0 Implementations”. In: *Detection of Intrusions and Malware, and Vulnerability Assessment*. Springer, 2015, pp. 239–260.
- [78] S.-T. Sun and K. Beznosov. “The devil is in the (implementation) details: an empirical analysis of OAuth SSO systems”. In: *ACM CCS*. 2012.
- [79] R. Yang et al. “Model-based Security Testing: An Empirical Study on OAuth 2.0 Implementations”. In: *ASIA CCS*. 2016.
- [80] *yes.com*. URL: [yes.com/](http://yes.com/).
- [81] *YouTube Data API Overview*. URL: [developers.google.com/youtube/v3/getting-started](https://developers.google.com/youtube/v3/getting-started).
- [82] Y. Zhou and D. Evans. “SSOScan: Automated Testing of Web Applications for Single Sign-On Vulnerabilities”. In: *USENIX*. 2014, pp. 495–510.

## APPENDIX A DETAILS ON APPLICATION-SPECIFIC MODEL

In this appendix, we continue the overview of our formal model from [Section IV-B](#) with additional details and references to the respective specifications, demonstrating how our model closely follows the specifications.

### A. Authorization Server Model

**Pushed Authorization Request.** We model the PAR endpoint as part of the AS (mandated by FAPI 2.0 [32, Section 5.3.1.2. No. 2]). The endpoint model stores all values and returns a (freshly chosen) request URI value. In particular, as mandated by FAPI 2.0, the model requires a redirect URI and a PKCE code challenge value in the PAR request [32, Section 5.3.1.2. No. 5, 6].

**Token Introspection.** The introspection endpoint part of the AS model expects an opaque access token and returns whether the token is active (and thus, if the token was issued by this AS at all), information on the key to which the token is bound, and a subject identifier, i.e., the identity of the user whose login at the AS lead to the AS issuing the token. As required in [69, Section 2.1], the AS model requires successful authentication of the RS that sent the request. The exact authentication mechanism for RSs at ASs is out of scope for FAPI 2.0. In our model, we use the HTTP Basic Authentication mechanism suggested by [69].

**Login Script.** Upon receiving the authorization request, the AS model responds with a script that models the login page for the user. As a result of executing the script, the browser sends a POST request to the AS with the login credentials (for the current AS) of a user that is using the browser.

**Further Endpoints.** The AS model also comprises an endpoint for server metadata that returns information about the server such as the different endpoint URLs. This endpoint is mandated by FAPI 2.0 [32, Section 5.3.1.1. No. 1]. The authorization endpoint of the AS model requires a request URI value as mandated in [32, Section 5.3.1.2. No. 3]. Furthermore, the model has a JWKS endpoint, where the server responds with the public signature verification key, which is strongly recommended by FAPI 2.0 [32, Section 5.6.3].

### B. Client Model

**Configurations.** As noted above, the client authenticates either via mTLS or the `private_key_jwt` method, and supports sender constraining by either mTLS or DPoP. In the model, the client can have different combinations of client authentication and sender constraining methods for different ASs. However, for a given pair of client and AS, this configuration does not change, reflecting that clients must register with ASs before starting a flow and select a configuration as part of that registration. We note that this registration is out of scope for FAPI 2.0 and we model it as part of the initial states of clients and ASs.

**Starting a Flow.** For starting a flow, the client model provides a script for browsers which triggers a POST request to the client. In our model, this request must contain the domain of an AS, modeling a user selecting an AS. When starting a flow for the first time with an AS, the client fetches the server metadata from the AS (as required by FAPI 2.0 [32, Section 5.3.2.1. No. 9]).

**End-User Authentication.** For each flow, the client decides non-deterministically whether it wants to authenticate the user. If that is the case, the client adds the value `openid` to the `scope` value contained in the pushed authorization request, hence requesting an id token from the AS. If the client model requests an id token, and once it receives the token response, it non-deterministically decides whether to log the user in based on the id token or to redeem the access token. The client logs in the user by creating a fresh cookie (called *service session id* in the client model) associated with the subject identifier contained in the id token.

**Further Client Aspects.** Furthermore, as FAPI 2.0 mandates, the client model always uses pushed authorization requests and PKCE, and always checks the *iss* issuer identifier in the redirection request (at the redirection endpoint) [32, Section 5.3.2.2. No. 2, 3, 4].

Another important detail concerns the handling of id tokens: FAPI 2.0 uses OpenID Connect id tokens, which means that id tokens are signed by their issuer. This of course raises the question of why the honest client in the attacker token injection attack accepts an attacker-constructed id token (see [Section III-A](#)). The reason is that OpenID Connect [70, Sec. 3.1.3.7 No. 6] (and thus FAPI 2.0) allows clients to skip signature verification on id tokens if the token is received directly from the token endpoint over a TLS-protected connection – these conditions are always fulfilled with FAPI 2.0. In the attacker token injection attack, the client contacts an attacker-controlled token endpoint (via a TLS-protected connection) and thus does not verify the id token signature.

### C. Resource Server Model

**Verification of Access Token.** A request for a resource must contain an access token in the HTTP headers [32, Section 5.3.3 No. 1]. The RS model verifies that the token is valid, identifies the resource owner for whose resources the token was issued, and the key to which the token is bound as follows. If the token is a structured token, the RS checks the signature of the token using the public verification key of the AS responsible for the requested resource and retrieves the resource owner information, as well as the key to which the token is bound from the token. Otherwise, the RS model sends the token and RS authentication information to the introspection endpoint of the AS that manages the requested resource. The token introspection response then contains the necessary information.

**Modeling Resources.** As the management of resources is not within the scope of the FAPI 2.0 specifications, we assumed a generic resource management which we describe in the following: The RS model manages different resources identified by URLs, in particular, by the path part of URLs. For a given resource, the RS model knows which AS manages the resource. Note that this is common (and necessary) in real ecosystems, at least when using opaque access tokens: how else can the RS determine where it has to send the token introspection request. The RS model identifies the resource owner through the access token (see above). Then, the RS model creates a fresh nonce, which represents the protected resource of the resource owner. That nonce is then returned the client.

## APPENDIX B ADDITIONAL ATTACKS AND INCONSISTENCIES IN THE FAPI 2.0 ATTACKER MODEL

### A. Authorization Request Leak Attack

In this attack, the authorization request (Step 8 in Figure 1) of an otherwise honest flow between an honest user  $u$ , an honest client  $c_{\text{hon}}$ , and honest AS  $AS_{\text{hon}}$  is blocked by (A1), but leaks to the attacker (A2). The attacker now visits  $AS_{\text{hon}}$ 's authorization endpoint, posing as  $u$  with the leaked authorization request, and logs in using its own (attacker) identity.  $AS_{\text{hon}}$  then answers with an authorization response, containing an authorization code  $code$ , which  $AS_{\text{hon}}$  associates with the attacker identity and the PAR of the initial, honest flow. From the authorization response, the attacker assembles a link pointing to  $c_{\text{hon}}$ 's redirection endpoint with  $code$  and the `iss` parameter as received from  $AS_{\text{hon}}$ , and sends this link to  $u$ . Once  $u$  follows this link,  $c_{\text{hon}}$  exchanges the  $code$  for tokens, logs  $u$  in under the attacker's identity and accesses resources belonging to the attacker in a session with  $u$ , hence violating both session integrity goals.

**Fix.** Upon our notification, the FAPI WG acknowledged the attack [3]. After some discussion, there was consensus that there is no currently deployed technology which can prevent such attacks. However, the FAPI WG wanted to keep A2 in their attacker model. Therefore, the FAPI WG decided to document the attack, as well as (optional) mitigations, which do not prevent, but harden against the attack in practice, in the specifications [32, Section 5.6.6]. For our analysis, we formulated the session integrity properties such that they only apply to flows in which the authorization request does not leak.

### B. Variants of Client Impersonation Attack

The attack described in Section III-B has two additional variants which we describe here.

Recall that in the attack variant described in Section III-B, the user authenticated at an honest, trusted AS and authorized not only an honest and trusted client, but also exactly the client which the user expected to authorize. However, if one assumes that the user does not thoroughly check which client they authorize, this attack even works if they initially contact a different client. In the first variant, we assumed that the user starts a flow with  $c_{\text{hon}}$ , i.e., an honest client of the user's choice, for which the attacker needs leaked client authentication JWTs (we describe how the attacker can obtain such JWTs in Section III-B). This variant, at the cost of the additional assumption that the user does not pay close attention to the client they are authorizing, makes the attack much more viable in practice: the attack can now select an arbitrary (but honest) client for which it can obtain client authentication JWTs. We note that this additional assumption about the user's behaviour is not a particularly strong one: even if the user pays attention, the AS also has to gather and show enough (reliable) information about the client, which can be quite difficult in practice, especially in dynamic environments, in which information AS has on clients are provided by the clients themselves in some automatic registration process.

In yet another variant of the attack, the attacker once again acts as  $c_{\text{hon}}$  towards  $AS_{\text{hon}}$ , but sends a PAR without waiting for the user to initiate a flow. The attacker then constructs a link similar to Step 9 and uses social engineering to make the user click on that link. The remainder of the attack is the same. We note that this attack variant assumes that the attacker can convince the user to not only click a link and authenticate at  $AS_{\text{hon}}$ , but also to authorize  $c_{\text{hon}}$  even though the user did not start a flow at all (however,  $AS_{\text{hon}}$  and  $c_{\text{hon}}$  may be entities the user knows and trusts).

### C. Inconsistencies in Attacker Model

In addition to the attacks described above, we discovered inconsistencies in the FAPI 2.0 Attacker Model, i.e., assumptions on the attacker which immediately violate one or more of the security goals.

**A6 Violates Authorization Goal.** As resource responses contain users' resources, this attacker assumption immediately violates the authorization security goal. When we pointed this out to the FAPI WG, it was decided to drop this attacker assumption [13].

**A7 Violates Session Integrity for Authorization Goal.** If the attacker can tamper with resource responses, the attacker can force users to use attacker resources by replacing (honest) resources in resource responses with attacker resources. This obviously violates the session integrity for authorization goal. As above, the FAPI WG decided to drop the attacker assumption once we pointed out this inconsistency [13].

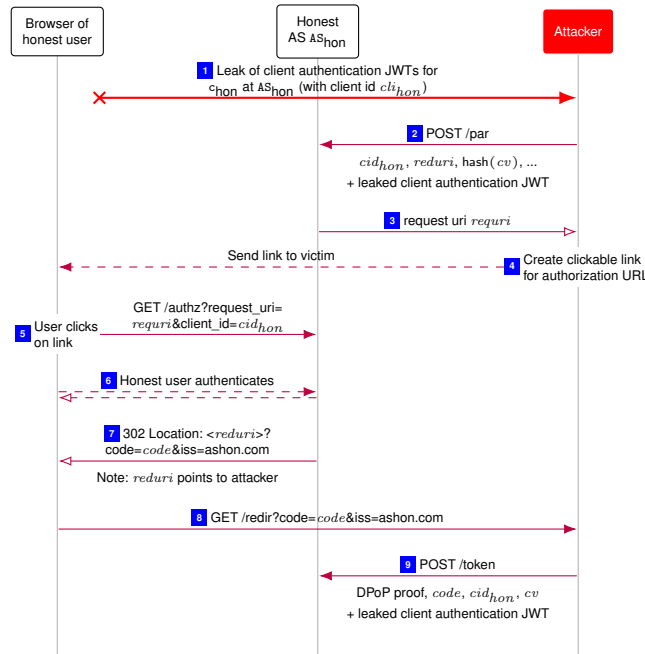


Figure 5. Variant of client impersonation attack

## APPENDIX C MODELING OF REMAINING ATTACKER ASSUMPTIONS

In Section IV-B2, we describe how we model the final assumptions on the attacker (i.e., as of the latest FAPI 2.0 Attacker Model). As mentioned there, our initial model contained all attacker assumptions as of July 1st (listed in Section II-C). We now detail how we incorporated the remaining assumptions (i.e., those which have been removed since July 1st) into our initial model and describe why we removed them from the final model.

$\mathcal{A}_3$  (authorization responses leak) was initially modeled by leakage of the authorization response at the AS, which sent the response in plain to a non-deterministically chosen IP address (in addition to sending the response to the client). As described in Section III-D, we had to remove this leak according to the changed attacker model.

$\mathcal{A}_4$  (attacker can trick client into using an attacker-controlled token endpoint) was initially modeled by the client non-deterministically choosing whether to use the correct token endpoint or a non-deterministically chosen endpoint. After the FAPI WG made it mandatory for clients to use server metadata, we removed this attacker capability (see Section III-B).

$\mathcal{A}_6$  and  $\mathcal{A}_7$ . As mentioned in Appendix B-C, the FAPI WG dropped these attacker assumptions after we reported inconsistencies in the attacker model resulting from these assumptions. Thus, we did not model these two assumptions.

## APPENDIX D FAPI 2.0 MODEL

In this section, we provide the full formal model of the FAPI 2.0 participants. We start with the definition of keys and secrets, as well as protocol participants and identities within the model, followed by how we initialize AS-client relationships and details on how *OAuth 2.0 Mutual TLS for Client Authentication and Certificate Bound Access Tokens* [11] is modeled. We continue with the formal models of the FAPI 2.0 clients (Appendix D-J), the FAPI 2.0 ASs (Appendix D-K), and the FAPI 2.0 RSs (Appendix D-L).

### A. Protocol Participants

We define the following sets of atomic Dolev-Yao processes: AS is the set of processes representing authorization servers. Their relation is described in Appendix D-K. RS is the set of processes representing resource servers, described in Appendix D-L. C is the set of processes representing clients, described in Appendix D-J. Finally, B is the set of processes representing browsers, including their users. They are described in Appendix G-G.

### B. Identities

Identities consist, similar to email addresses, of a user name and a domain part. For our model, this is defined as follows:

*Definition 5.* An identity  $i$  is a term of the form  $\langle name, domain \rangle$  with  $name \in \mathbb{S}$  and  $domain \in \text{Doms}$ . Let  $\text{ID}$  be the finite set of identities. We say that an id is *governed* by the DY process to which the domain of the id belongs. This is formally captured by the mappings  $\text{governor} : \text{ID} \rightarrow \mathcal{W}, \langle name, domain \rangle \mapsto \text{dom}^{-1}(domain)$  and  $\text{ID}^y := \text{governor}^{-1}(y)$ .

### C. Keys and Secrets

The set  $\mathcal{N}$  of nonces is partitioned into disjoint sets, an infinite set  $N$ , and finite sets  $K_{\text{TLS}}, K_{\text{sign}}, \text{Passwords}$ , and  $\text{RScredentials}$ :

$$\mathcal{N} = N \uplus K_{\text{TLS}} \uplus K_{\text{sign}} \uplus \text{Passwords} \uplus \text{RScredentials}$$

These sets are used as follows:

- The set  $N$  contains the nonces that are available for the DY processes
- The set  $K_{\text{TLS}}$  contains the keys that will be used for TLS encryption. Let  $\text{tlskey} : \text{Doms} \rightarrow K_{\text{TLS}}$  be an injective mapping that assigns a (different) private key to every domain. For an atomic DY process  $p$  we define  $\text{tlskeys}^p = \langle \langle d, \text{tlskey}(d) \rangle \mid d \in \text{dom}(p) \rangle$  (i.e., a sequence of pairs).
- The set  $K_{\text{sign}}$  contains the keys that will be used by ASs for signing id and access tokens and by clients to sign JWTs as well as DPOP proofs. Let  $\text{signkey} : \text{AS} \cup \text{C} \rightarrow K_{\text{sign}}$  be an injective mapping that assigns a (different) signing key to every AS and client.
- The set  $\text{Passwords}$  is the set of passwords (secrets) the browsers share with servers. These are the passwords the users use to log in. Let  $\text{secretOfID} : \text{ID} \rightarrow \text{Passwords}$  be a bijective mapping that assigns a password to each identity.
- The set  $\text{RScredentials}$  is a set of secrets shared between authorization and resource servers. RSs use these to authenticate at ASs' token introspection endpoints. Let  $\text{secretOfRS} : \text{Doms} \times \text{Doms} \rightarrow \text{RScredentials}$  be a partial mapping, assigning a secret to some of the RS-AS pairs (with the function arguments in that order).

### D. Passwords

*Definition 6.* Let  $\text{ownerOfSecret} : \text{Passwords} \rightarrow \text{B}$  be a mapping that assigns to each password a browser which *owns* this password. Similarly, we define  $\text{ownerOfID} : \text{ID} \rightarrow \text{B}$  as  $i \mapsto \text{ownerOfSecret}(\text{secretOfID}(i))$ , which assigns to each identity the browser that owns this identity (i.e., this identity belongs to the browser).

### E. Web Browsers

Web browser processes (i.e., processes  $b \in \text{B}$ ) are modeled as described in [Appendix G](#). Before defining the initial states of Web browsers, we introduce the following set (for some process  $p$ ):

$$\text{Secrets}^{b,p} = \{s \mid b = \text{ownerOfSecret}(s) \wedge (\exists i : s = \text{secretOfID}(i) \wedge i \in \text{ID}^p)\}$$

*Definition 7 (Initial Web Browser State for FAPI).* The initial state of a Web browser process  $b \in \text{B}$  follows the description in [Definition 68](#), with the following additional constraints:

- $s_0^b.\text{ids} \equiv \langle \{i \mid b = \text{ownerOfID}(i)\} \rangle$
- $s_0^b.\text{secrets}$  contains an entry  $\langle \langle d, s \rangle, \langle \text{Secrets}^{b,p} \rangle \rangle$  for each  $p \in \text{AS} \cup \text{C} \cup \text{RS}$  and every domain  $d \in \text{dom}(p)$  (and nothing else), i.e.,
$$s_0^b.\text{secrets} \equiv \left\langle \left\{ \langle \langle d, s \rangle, \langle \text{Secrets}^{b,p} \rangle \rangle \mid \exists p, d : p \in \text{AS} \cup \text{C} \cup \text{RS} \wedge d \in \text{dom}(p) \right\} \right\rangle$$
- $s_0^b.\text{keyMapping} \equiv \langle \{ \langle d, \text{pub}(\text{tlskey}(d)) \rangle \mid d \in \text{Doms} \} \rangle$

### F. Resources

We model the management of resources as follows: We assume that each resource is managed by at most one AS. We also assume that resources are identified by URLs at the RS. Thus, when getting a request to such a resource URL, the RS has to

- 1) identify the AS that is managing the resource, and
- 2) identify the identity for which the access token was issued.

If the access token is a structured JWT, the RS retrieves the identity from the subject field. Otherwise, the identity is retrieved from the introspection response.

For identifying the AS, we first define the URL paths of resources managed by a RS, and then define a mapping from these paths to AS.

*Definition 8.* For each  $rs \in \text{RS}$ , let  $\text{resourceURLPath}^{rs} \subseteq \mathbb{S}$  be a finite set of strings. These are the URL paths identifying the resources managed by the RS.<sup>5</sup>

<sup>5</sup>A resource is managed by the RS if and only if  $\text{resourceID} \in \text{resourceURLPath}^{rs}$ .

*Definition 9.* For each  $rs \in RS$ , let  $\text{supportedAuthorizationServer}^{rs} \subseteq AS$  be a finite set of ASs. These are the ASs supported by the RS.

*Definition 10.* For each  $rs \in RS$ , let  $\text{authorizationServerOfResource}^{rs} : \text{resourceURLPath}^{rs} \rightarrow \text{supportedAuthorizationServer}^{rs}$  be a mapping that assigns an AS to each resource URL path suffix of resources managed by the RS.

If the access token is valid and the resource is managed by an AS supported by the RS, the RS model responds with a fresh nonce that it stores under the identity of the resource owner and the path under which it returns the resource. By using fresh nonces, the RS does not return a nonce twice – even for requests for the same path and the same resource owner (identified via token introspection or the sub claim in the access token). Without this, the authorization property would need to exclude the case that the resource owner granted some malicious client access to a resource at some point.

### G. Client Registration

ASs and clients have to establish some relationship with each other before starting a FAPI 2.0 flow. Such a relationship can be established out of band, e.g., via manual configuration. While FAPI 2.0 also supports the use of dynamic client registration, our model assumes an out of band registration, captured by the following definitions.

*Definition 11.* A *client information dictionary* is a dictionary of the form  $[\text{client\_id}: \text{clientId}, \text{client\_type}: \text{clientType}, \text{mtls\_skey}: \text{mtlsSkey}, \text{jwt\_skey}: \text{jwtSkey}]$  where  $\text{clientId} \in \mathcal{T}_{\mathcal{N}}$ ,  $\text{clientType} \in \{\text{mTLS\_mTLS}, \text{mTLS\_DPoP}, \text{pkjwt\_mTLS}, \text{pkjwt\_DPoP}\}$ ,  $\text{mtlsSkey} \in K_{\text{tls}} \cup \{\diamond\}$ , and  $\text{jwtSkey} \in K_{\text{sign}} \cup \{\diamond\}$ . We further require  $\text{jwtSkey} \neq \diamond$  if and only if  $\text{clientType} \in \{\text{pkjwt\_mTLS}, \text{pkjwt\_DPoP}, \text{mTLS\_DPoP}\}$ , as well as  $\text{mtlsSkey} \neq \diamond$  if and only if  $\text{clientType} \in \{\text{mTLS\_mTLS}, \text{pkjwt\_mTLS}, \text{mTLS\_DPoP}\}$ . Let  $\text{ClientInfos}$  be the set of all client information dictionaries.

*Definition 12.* Let  $\text{clientInfo}: \text{Doms} \times \text{Doms} \rightarrow \text{ClientInfos}$  be a (partial) mapping from AS and client domains to client information dictionaries, assigning a client information dictionary to some of the possible AS–client pairs with the following restrictions: 1) There are no two clients with the same  $\text{clientId}$  at the same AS, formalized as  $\forall \text{as} \in AS, d_{\text{as}} \in \text{dom}(\text{as}), d_{\text{client}}, d'_{\text{client}} \in \text{Doms}: \text{clientInfo}(d_{\text{as}}, d_{\text{client}}).\text{client\_id} \equiv \text{clientInfo}(d_{\text{as}}, d'_{\text{client}}).\text{client\_id} \Rightarrow d_{\text{client}} \equiv d'_{\text{client}}$ , 2) TLS keys are assigned according to  $\text{tlskey}$ , formalized as  $\text{clientInfo}(d_{\text{as}}, d_{\text{client}}).\text{mtlsSkey} \in \{\diamond, \text{tlskey}(d_{\text{client}})\}$ , and 3) signing keys are assigned according to  $\text{signkey}$ , formally expressed as  $\text{clientInfo}(d_{\text{as}}, d_{\text{client}}).\text{jwtSkey} \in \{\diamond, \text{signkey}(\text{dom}^{-1}(d_{\text{client}}))\}$ . Note that while this definition requires the client to use the same key to sign JWTs and DPoP proofs for all ASs, it allows the client to use a different client type for each AS. mTLS keys are different for each of the client’s domains.

*Definition 13.* A client  $c \in C$  has the client identifier  $\text{clientId}$  at an authorization server  $\text{as} \in AS$  if  $\exists d_{\text{client}} \in \text{dom}(c), d_{\text{as}} \in \text{dom}(\text{as})$  such that  $\text{clientInfo}(d_{\text{as}}, d_{\text{client}}).\text{client\_id} \equiv \text{clientId}$ .

*Definition 14.* Let  $\text{clientInfoAS}: AS \rightarrow \mathcal{T}_{\mathcal{N}}$  be a (partial) mapping from an AS to a dictionary. The keys of this dictionary are client IDs and the values AS *client information* dictionaries. We define  $\text{clientInfoAS}$  by  $\text{as} \mapsto \{\{\langle \text{cli}.\text{client\_id}, \text{as\_cli}(\text{cli}) \rangle \mid \exists d_{\text{client}} \in \text{Doms}, d_{\text{as}} \in \text{dom}(\text{as}): \text{clientInfo}(d_{\text{as}}, d_{\text{client}}) \equiv \text{cli}\}\}$  where  $\text{as\_cli}(\text{cli}) := [\text{client\_type}: \text{cli}.\text{client\_type}] + \langle \rangle$

$$\begin{cases} [\text{mtls\_key}: \text{pub}(\text{cli}.\text{mtlsSkey})] & \text{if } \text{cli}.\text{client\_type} \equiv \text{mTLS\_mTLS} \\ [\text{mtls\_key}: \text{pub}(\text{cli}.\text{mtlsSkey}), \text{jwt\_key}: \text{pub}(\text{cli}.\text{jwtSkey})] & \text{if } \text{cli}.\text{client\_type} \in \{\text{pkjwt\_mTLS}, \text{mTLS\_DPoP}\} \\ [\text{jwt\_key}: \text{pub}(\text{cli}.\text{jwtSkey})] & \text{if } \text{cli}.\text{client\_type} \equiv \text{pkjwt\_DPoP} \end{cases}$$

Note: In  $\exists d_{\text{client}} \in \text{Doms}, d_{\text{as}} \in \text{dom}(\text{as}): \text{clientInfo}(d_{\text{as}}, d_{\text{client}}) \equiv \text{cli}$ , we refer to values  $d_{\text{client}}$  and  $d_{\text{as}}$  for which  $\text{clientInfo}(d_{\text{as}}, d_{\text{client}})$  is defined.

*Definition 15.* Let  $\text{clientInfoClient}: C \rightarrow [\text{Doms} \times \mathcal{T}_{\mathcal{N}}]$  be a (partial) mapping from a client to a dictionary. The keys of this dictionary are AS domains and the values simple dictionaries, containing client type and client ID.  $\text{clientInfoClient}$  is defined as  $c \mapsto \{\{\{d_{\text{as}}: [\text{client\_id}: \text{cli}.\text{client\_id}, \text{client\_type}: \text{cli}.\text{client\_type}] \mid \exists d_{\text{client}} \in \text{dom}(c), d_{\text{as}} \in \text{Doms}: \text{clientInfo}(d_{\text{as}}, d_{\text{client}}) \equiv \text{cli}\}\}\}$ . Note: In  $\exists d_{\text{client}} \in \text{dom}(c), d_{\text{as}} \in \text{Doms}: \text{clientInfo}(d_{\text{as}}, d_{\text{client}}) \equiv \text{cli}$ , we refer to values  $d_{\text{client}}$  and  $d_{\text{as}}$  for which  $\text{clientInfo}(d_{\text{as}}, d_{\text{client}})$  is defined.

### H. Modeling mTLS

*OAuth 2.0 Mutual TLS for Client Authentication and Certificate Bound Access Tokens* (mTLS) [11] provides a method for both client authentication and token binding. Note that both mechanisms may be used independently of each other.

OAuth 2.0 Mutual TLS Client Authentication makes use of *TLS client authentication*<sup>6</sup>, which the client can use for client authentication at the pushed authorization request and token endpoints (in Step [5] and Step [14] of Figure 1). In TLS client authentication, not only the server authenticates to the client (as is common for TLS), but the client also authenticates to the server.

<sup>6</sup>As noted in Section 7.2 of [11], this extension supports all TLS versions with certificate-based client authentication.

To this end, the client proves that it knows the private key belonging to a certificate that is either (a) self-signed and pre-configured at the respective AS or that is (b) issued for the respective client id by a predefined certificate authority within a public key infrastructure (PKI).

Token binding means binding an access token to a client such that only this client is able to use the access token at the RS. To achieve this, the AS associates the access token with the certificate used by the client for the TLS connection to the token endpoint. In the TLS connection to the RS (in Step [18] of Figure 1), the client then authenticates using the same certificate. The RS accepts the access token only if the client certificate is the one associated with the access token.<sup>7</sup>

The WIM models TLS at a high level of abstraction. An HTTP request is encrypted with the public key of the recipient and contains a symmetric key, which is used for encrypting the HTTP response. Furthermore, the model contains no certificates or public key infrastructures but uses a function that maps domains to their public key.

We model mTLS similarly to [38]. An overview of the mTLS model is shown in Figure 6. The basic idea is that the server sends a nonce encrypted with the public key of the client. The client proves possession of the private key by decrypting this message. In Step [1], the client sends its client identifier to the AS. The AS then looks up the public key associated with the client identifier, chooses a nonce, and encrypts it with the public key. As depicted in Step [2], the server additionally includes its public key. When the client decrypts the message, it checks if the public key belongs to the server it wants to send the original message to. This prevents man-in-the-middle attacks, as only the honest client can decrypt the response and as the public key of the server cannot be changed by an attacker. In Step [3], the client sends the original request with the decrypted nonce. When the server receives this message, it knows that the nonce was decrypted by the honest client (as only the client knows the corresponding private key) and that the client had chosen to send the nonce to the server (due to the public key included in the response). Therefore, the server can conclude that the message was sent by the honest client.

In effect, this resembles the behavior of the TLS handshake, as the verification of the client certificate in TLS is done by signing all handshake messages [68, Section 7.4.8], which also includes information about the server certificate, which means that the signature cannot be reused for another server. Instead of signing a sequence that contains information about the receiver, in our model, the client checks the sender of the nonce, and only sends the decrypted nonce to the creator of the nonce. In other words, a nonce decrypted by an honest server that gets decrypted by the honest client is never sent to the attacker.

As explained above, the client uses the same certificate it used for the token request when sending the access token to the RS. While the RS has to check the possession of corresponding private keys, the validity of the certificate was already checked at the AS and can be ignored by the RS. Therefore, in our model of FAPI 2.0, the client does not send its client id to the RS, but its public key, and the RS encrypts the message with this public key.

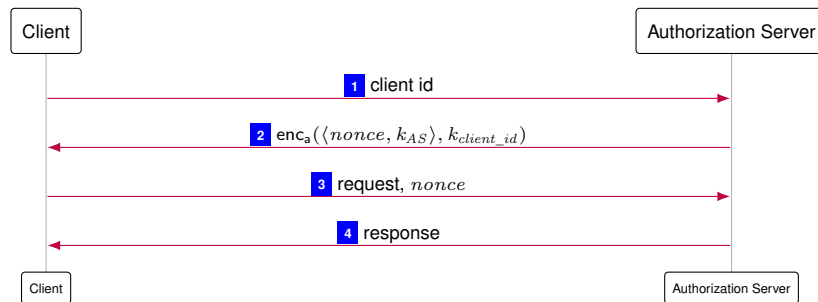


Figure 6. Overview of mTLS model

All messages are sent by the generic HTTPS server model (Appendix G-L), which means that each request is encrypted asymmetrically, and the responses are encrypted symmetrically with a key that was included in the request. For completeness, Figure 7 shows the complete messages, i.e., with the encryption used for transmitting the messages.

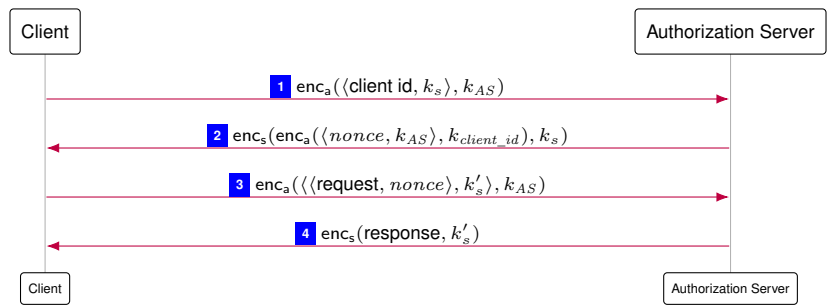
### I. Additional HTTP Headers

In order to model DPoP, we extend the list of headers of Definition 42 with the following header:

$(\text{DPoP}, p)$  where  $p \in \mathcal{T}_{\mathcal{N}}$  is (for honest senders) a DPoP proof (i.e., a signed JWT).

<sup>7</sup>The RS can read this information either directly from the access token if the access token is a signed document, or uses token introspection to retrieve the data from the AS.





**Figure 7.** Detailed view on mTLS model

## J. Clients

A client  $c \in \mathcal{C}$  is a Web server modeled as an atomic DY process  $(I^c, Z^c, R^c, s_0^c)$  with the addresses  $I^c := \text{addr}(c)$ . Next, we define the set  $Z^c$  of states of  $c$  and the initial state  $s_0^c$  of  $c$ .

*Definition 16.* A state  $s \in Z^c$  of a client  $c$  is a term of the form  $\langle \text{DNSAddress}, \text{pendingDNS}, \text{pendingRequests}, \text{corrupt}, \text{keyMapping}, \text{tlskeys}, \text{sessions}, \text{oauthConfigCache}, \text{jwksCache}, \text{asAccounts}, \text{mtlsCache}, \text{jwt}, \text{resourceASMapping}, \text{dpopNonces} \rangle$  with  $\text{DNSAddress} \in \text{IPs}$ ,  $\text{pendingDNS} \in [\mathcal{N} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{pendingRequests} \in [\mathcal{N} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{corrupt} \in \mathcal{T}_{\mathcal{N}}$ ,  $\text{keyMapping} \in [\text{Doms} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{tlskeys} \in [\text{Doms} \times K_{\text{TLS}}]$  (all former components as in [Definition 71](#)),  $\text{sessions} \in [\mathcal{N} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{oauthConfigCache} \in [\mathcal{T}_{\mathcal{N}} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{jwksCache} \in [\mathcal{T}_{\mathcal{N}} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{asAccounts} \in [\mathcal{T}_{\mathcal{N}} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{mtlsCache} \in \mathcal{T}_{\mathcal{N}}$ ,  $\text{jwt} \in K_{\text{sign}}$ ,  $\text{resourceASMapping} \in [\text{Doms} \times [\mathbb{S} \times \text{Doms}]]$ , and  $\text{dpopNonces} \in [\text{Doms} \times \mathcal{T}_{\mathcal{N}}]$ .

An initial state  $s_0^c$  of  $c$  is a state of  $c$  with

- $s_0^c.\text{DNSAddress} \in \text{IPs}$ ,
- $s_0^c.\text{pendingDNS} \equiv \langle \rangle$ ,
- $s_0^c.\text{pendingRequests} \equiv \langle \rangle$ ,
- $s_0^c.\text{corrupt} \equiv \perp$ ,
- $s_0^c.\text{keyMapping}$  being the same as the keymapping for browsers,
- $s_0^c.\text{tlskeys} \equiv \text{tlskeys}^c$  (see [Appendix D-C](#)),
- $s_0^c.\text{sessions} \equiv \langle \rangle$ ,
- $s_0^c.\text{oauthConfigCache} \equiv \langle \rangle$ ,
- $s_0^c.\text{jwksCache} \equiv \langle \rangle$ ,
- $s_0^c.\text{asAccounts} \equiv \text{clientInfoClient}(c)$  (see [Definition 15](#)),
- $s_0^c.\text{mtlsCache} \equiv \langle \rangle$ ,
- $s_0^c.\text{jwt} \equiv \text{signkey}(c)$  (see [Appendix D-C](#)),
- $s_0^c.\text{resourceASMapping}[\text{domRS}][\text{resourceID}] \in \text{dom}(\text{authorizationServerOfResource}^{rs}(\text{resourceID}))$ ,  $\forall rs \in \text{RS}$  and  $\forall \text{domRS} \in \text{dom}(rs)$  and  $\forall \text{resourceID} \in \text{resourceURLPath}^{rs}$  (a domain of the AS managing the resource stored at  $rs$  identified by  $\text{resourceID}$ ), and
- $s_0^c.\text{dpopNonces} \equiv \langle \rangle$ .

We now specify the relation  $R^c$ : This relation is based on the model of generic HTTPS servers (see [Appendix G-L](#)). Hence we only need to specify algorithms that differ from or do not exist in the generic server model. These algorithms are defined in [Algorithms 2–7](#). Note that in several places throughout these algorithms, we use placeholders of the form  $\nu_x$  to generate “fresh” nonces as described in the communication model (see [Definition 27](#)).

The script that is used by the client on its index page is specified in [Algorithm 8](#). This script uses the  $\text{GETURL}(tree, \text{docnonce})$  function to extract the current URL of a document. We define this function as follows: It searches for the document with the identifier  $\text{docnonce}$  in the (cleaned) tree  $tree$  of the browser’s windows and documents. It then returns the URL  $u$  of that document. If no document with nonce  $\text{docnonce}$  is found in the tree  $tree$ ,  $\diamond$  is returned.

---

**Algorithm 1** Relation of a Client  $R^c$  – Processing HTTPS Requests

---

```
1: function PROCESS_HTTPS_REQUEST( $m, k, a, f, s'$ ) → Process an incoming HTTPS request. Other message types are handled
   in separate functions.  $m$  is the incoming message,  $k$  is the encryption key for the response,  $a$  is the receiver,  $f$  the sender of the message.  $s'$ 
   is the current state of the atomic DY process  $c$ .
2:   if  $m.path \equiv /$  then → Serve index page (start flow).
3:     let  $m' := enc_s(\langle\langle HTTPResp, m.nononce, 200, headers, \langle script\_client\_index, \langle \rangle \rangle \rangle, k)$  → Reply with  $script\_client\_index$ .
4:     stop  $\langle\langle f, a, m' \rangle\rangle, s'$ 
5:   else if  $m.path \equiv /startLogin \wedge m.method \equiv POST$  then → Start a new FAPI 2.0 flow (probably from  $script\_client\_index$ )
6:     if  $m.headers[Origin] \neq \langle m.host, S \rangle$  then
7:       stop → Check the Origin header for CSRF protection to prevent attacker from starting a flow in the background (as this would
         trivially violate the session integrity property).
8:     let  $selectedAS := m.body$ 
9:     let  $sessionId := \nu_1$  → Session id is a freshly chosen nonce.
10:    let  $s'.sessions[sessionId] := [startRequest: [message: m, key: k, receiver: a, sender: f], selected\_AS: selectedAS]$ 
11:    call PREPARE_AND_SEND_PAR( $sessionId, a, s'$ ) → Start an authorization flow with the AS (see Algorithm 6)
12:  else if  $m.path \equiv /redirect\_ep$  then → User is being redirected after authentication to the AS.
13:    let  $sessionId := m.headers[Cookie][\langle \_Host, sessionId \rangle]$ 
14:    if  $sessionId \notin s'.sessions$  then
15:      stop
16:    let  $session := s'.sessions[sessionId]$  → Retrieve session data.
17:    let  $selectedAS := session[selected\_AS]$ 
18:    if  $m.parameters[iss] \neq selectedAS$  then → Check issuer parameter ([75]).
19:      stop
      → Store browser's request for use in CHECK_ID_TOKEN (Algorithm 5) and PROCESS_HTTPS_RESPONSE (Algorithm 2)
20:    let  $s'.sessions[sessionId][redirectEpRequest] := [message: m, key: k, receiver: a, sender: f]$ 
21:    call SEND_TOKEN_REQUEST( $sessionId, m.parameters[code], a, s'$ ) → Retrieve a token from AS's token endpoint.
22:  stop → Unknown endpoint or malformed request.
```

---

---

**Algorithm 2** Relation of a Client  $R^c$  – Processing HTTPS Responses

---

```
1: function PROCESS_HTTPS_RESPONSE( $m, reference, request, a, f, s'$ )
2:   if  $reference[responseTo] \equiv \text{MTLS}$  then  $\rightarrow$  Client received an mTLS nonce (see Appendix D-H)
3:     let  $m_{dec}, k'$  such that  $m_{dec} \equiv \text{dec}_a(m.\text{body}, k') \wedge \langle \text{dom}, k' \rangle \in s'.\text{tlskeys}$  if possible; otherwise stop
4:     let  $\text{mtlsNonce}, \text{serverPubKey}$  such that  $m_{dec} \equiv \langle \text{mtlsNonce}, \text{serverPubKey} \rangle$  if possible; otherwise stop
5:     if  $\text{serverPubKey} \equiv s'.\text{keyMapping}[request]$  then  $\rightarrow$  Verify sender of mTLS nonce
6:       let  $\text{clientId} := reference[\text{client\_id}]$   $\rightarrow$  Note: If  $\text{clientId} \notin reference[\text{client\_id}] \equiv \langle \rangle$ 
7:       let  $\text{pubKey} := reference[\text{pub\_key}]$   $\rightarrow$  See note for client ID above
8:       let  $s'.\text{mtlsCache} := s'.\text{mtlsCache} + \langle \rangle (\text{request}.\text{host}, \text{clientId}, \text{pubKey}, \text{mtlsNonce})$ 
9:     stop  $\langle \rangle, s'$ 
10:  let  $\text{sessionId} := reference[\text{session}]$ 
11:  let  $\text{session} := s'.\text{sessions}[\text{sessionId}]$ 
12:  let  $\text{selectedAS} := \text{session}[\text{selected\_AS}]$ 
13:   $\rightarrow$  Note: PREPARE_AND_SEND_PAR issues CONFIG and JWKS requests as required – once these get a response, we continue the
14:  PAR preparation by calling PREPARE_AND_SEND_PAR again.
15:  if  $reference[responseTo] \equiv \text{CONFIG}$  then
16:    if  $m.\text{body}[\text{issuer}] \neq \text{selectedAS}$  then  $\rightarrow$  Verify issuer identifier according to [48, Sec. 3.3]
17:    stop
18:    let  $s'.\text{oauthConfigCache}[\text{selectedAS}] := m.\text{body}$ 
19:    call PREPARE_AND_SEND_PAR( $\text{sessionId}, a, s'$ )
20:  else if  $reference[responseTo] \equiv \text{JWKS}$  then
21:    let  $s'.\text{jwtCache}[\text{selectedAS}] := m.\text{body}$ 
22:    call PREPARE_AND_SEND_PAR( $\text{sessionId}, a, s'$ )
23:  else if  $reference[responseTo] \equiv \text{PAR}$  then
24:    let  $\text{requestUri} := m.\text{body}[\text{request\_uri}]$ 
25:    let  $s'.\text{sessions}[\text{sessionId}][\text{request\_uri}] := \text{requestUri}$ 
26:    let  $\text{clientId} := \text{session}[\text{client\_id}]$ 
27:    let  $\text{request} := \text{session}[\text{startRequest}]$ 
28:     $\rightarrow$  In the following, we construct the response to the initial request by some browser
29:    let  $\text{authEndpoint} := s'.\text{oauthConfigCache}[\text{selectedAS}][\text{auth\_ep}]$ 
30:     $\rightarrow$  The authorization endpoint URL may include query components, which must be retained while also ensuring that no parameter
31:    appears more than once ([41, Sec. 3.1]). However, following [57, Sec. 4] and [73, Sec. 5] closely could introduce duplicates. We
32:    opted to overwrite  $\text{client\_id}$  and  $\text{request\_uri}$  parameters if present.
33:    let  $\text{authEndpoint}.\text{parameters}[\text{client\_id}] := \text{clientId}$ 
34:    let  $\text{authEndpoint}.\text{parameters}[\text{request\_uri}] := \text{requestUri}$ 
35:    let  $\text{headers} := [\text{Location}: \text{authEndpoint}]$ 
36:    let  $\text{headers}[\text{Set-Cookie}] := [(\_Host, \text{sessionId}): \langle \text{sessionId}, \top, \top, \top \rangle]$ 
37:    let  $\text{response} := \text{enc}_s(\langle \text{HTTPResp}, \text{request}[\text{message}].\text{nonce}, 303, \text{headers}, \langle \rangle \rangle, \text{request}[\text{key}])$ 
38:    let  $\text{leakAuthZReq} \leftarrow \{ \top, \perp \}$   $\rightarrow$  We assume that the authorization request, in particular  $\text{request\_uri}$  and  $\text{client\_id}$ , may leak
39:    to the attacker, see [30] and Section IV-B2.
40:  if  $\text{leakAuthZReq} \equiv \top$  then
41:    let  $\text{leak} := \langle \text{LEAK}, \text{authEndpoint} \rangle$ 
42:    let  $\text{leakAddress} \leftarrow \text{IPs}$ 
43:    stop  $\langle \langle \text{request}[\text{sender}], \text{request}[\text{receiver}], \text{response} \rangle, \langle \text{leakAddress}, \text{request}[\text{receiver}], \text{leak} \rangle \rangle, s'$ 
44:  else
45:    stop  $\langle \langle \text{request}[\text{sender}], \text{request}[\text{receiver}], \text{response} \rangle \rangle, s'$ 
46:  else if  $reference[responseTo] \equiv \text{TOKEN}$  then
47:    let  $\text{useAccessTokenNow} := \top$ 
48:    if  $\text{session}[\text{scope}] \equiv \text{openid}$  then  $\rightarrow$  Non-deterministically decide whether to use the AT or check the ID token (if requested)
49:      let  $\text{useAccessTokenNow} \leftarrow \{ \top, \perp \}$ 
50:    if  $\text{useAccessTokenNow} \equiv \top$  then
51:      call USE_ACCESS_TOKEN( $reference[\text{session}], m.\text{body}[\text{access\_token}], \text{request}.\text{host}, a, s'$ )
52:    let  $\text{selectedAsTokenEp} := s'.\text{oauthConfigCache}[\text{selectedAS}][\text{token\_ep}]$ 
53:    if  $\text{request}.\text{host} \neq \text{selectedAsTokenEp}.\text{host}$  then
54:      stop  $\rightarrow$  Verify sender of HTTPS response is the expected AS (see [70, Sec. 3.1.3.7])
55:    call CHECK_ID_TOKEN( $reference[\text{session}], m.\text{body}[\text{id\_token}], s'$ )
56:  else if  $reference[responseTo] \equiv \text{RESOURCE\_USAGE}$  then
57:     $\rightarrow$  Reply to browser's request to the client's redirect endpoint (with the retrieved resource as payload)
58:    let  $\text{resource} := m.\text{body}[\text{resource}]$ 
59:    let  $s'.\text{sessions}[\text{sessionId}][\text{resource}] := \text{resource}$   $\rightarrow$  Store received resource
60:    let  $s'.\text{sessions}[\text{sessionId}][\text{resourceServer}] := \text{request}.\text{host}$   $\rightarrow$  Store the domain of the RS
61:    let  $\text{request} := \text{session}[\text{redirectEpRequest}]$   $\rightarrow$  Data on browser's request to client's redirect endpoint
62:    let  $m' := \text{enc}_s(\langle \text{HTTPResp}, \text{request}[\text{message}].\text{nonce}, 200, \langle \rangle, \text{resource} \rangle, \text{request}[\text{key}])$ 
63:    stop  $\langle \langle \text{request}[\text{sender}], \text{request}[\text{receiver}], m' \rangle \rangle, s'$ 
64:   $\rightarrow$  Algorithm continues on next page.
```

---

---

```

56: else if reference[responseTo]  $\equiv$  DPOP_NONCE then
57:   let dpopNonce := m.body[nonce]
58:   let rsDomain := request.host
59:   let s'.dpopNonces[rsDomain] := s'.dpopNonces[rsDomain] + $\diamond$  dpopNonce
60:   stop  $\langle \rangle$ , s'
61: stop

```

---

### Algorithm 3 Relation of a Client $R^c$ – Request to token endpoint.

---

```

1: function SEND_TOKEN_REQUEST(sessionId, code, a, s')
2:   let session := s'.sessions[sessionId]
3:   let pkceVerifier := session[code_verifier]
4:   let selectedAS := session[selected_AS]
5:   let headers := []
6:   let body := [grant_type: authorization_code, code: code, redirect_uri: session[redirect_uri]]
7:   let body[code_verifier] := pkceVerifier  $\rightarrow$  add PKCE Code Verifier ([72], Section 4.5)
8:   let clientId := s'.asAccounts[selectedAS][client_id]
9:   let clientType := s'.asAccounts[selectedAS][client_type]
10:  let oauthConfig := s'.oauthConfigCache[selectedAS]
11:  let tokenEndpoint := oauthConfig[token_ep]
     $\rightarrow$  Client Authentication:
12:  if clientType  $\in$  {mTLS_mTLS, mTLS_DPoP} then  $\rightarrow$  mTLS client authentication
13:    let body[client_id] := clientId  $\rightarrow$  [11] mandates client_id when using mTLS authentication
14:    let mtlsNonce such that  $\langle$ tokenEndpoint.host, clientId,  $\langle \rangle$ , mtlsNonce $\rangle \in$  s'.mtlsCache if possible; otherwise stop
15:    let authData := [TLS_AuthN: mtlsNonce]
16:    let s'.mtlsCache := s'.mtlsCache - $\diamond$   $\langle$ tokenEndpoint.host, clientId,  $\langle \rangle$ , mtlsNonce $\rangle$ 
17:  else if clientType  $\in$  {pkjwt_mTLS, pkjwt_DPoP} then  $\rightarrow$  private_key_jwt client authentication
18:    let jwt := [iss: clientId, sub: clientId, aud: selectedAS]
19:    let jws := sig(jwt, s'.jwk)
20:    let authData := [client_assertion: jws]
21:  else
22:    stop  $\rightarrow$  Invalid client type
     $\rightarrow$  Sender Constraining:
23:  if clientType  $\equiv$  mTLS_mTLS then  $\rightarrow$  mTLS sender constraining (same nonce as for mTLS authN)
24:    let mtlsNonce := authData[TLS_AuthN]
25:    let body[TLS_binding] := mtlsNonce
26:  else if clientType  $\equiv$  pkjwt_mTLS then  $\rightarrow$  mTLS sender constraining (fresh mTLS nonce)
27:    let mtlsNonce such that  $\langle$ tokenEndpoint.host, clientId,  $\langle \rangle$ , mtlsNonce $\rangle \in$  s'.mtlsCache if possible; otherwise stop
28:    let s'.mtlsCache := s'.mtlsCache - $\diamond$   $\langle$ tokenEndpoint.host, clientId,  $\langle \rangle$ , mtlsNonce $\rangle$ 
29:    let body[TLS_binding] := mtlsNonce
30:  else  $\rightarrow$  Sender constraining using DPoP
31:    let privKey := s'.jwk  $\rightarrow$  get private key
32:    let htu := tokenEndpoint
33:    let htu[parameters] :=  $\langle \rangle$   $\rightarrow$  [39, Sec. 4.2]: without query
34:    let htu[fragment] :=  $\perp$   $\rightarrow$  [39, Sec. 4.2]: without fragment
35:    let dpopJwt := [headers: [jwk: pub(privKey)]]
36:    let dpopJwt[payload] := [htm: POST, htu: htu]
37:    let dpopProof := sig(dpopJwt, privKey)
38:    let headers[DPoP] := dpopProof  $\rightarrow$  add DPoP header; the dpopJwt can be extracted with the extractmsg() function
39:  let body := body + $\diamond$  authData
40:  let message := (HTTPReq, v2, POST, tokenEndpoint.host, tokenEndpoint.path, tokenEndpoint.parameters, headers, body)
41:  call HTTPS_SIMPLE_SEND(responseTo: TOKEN, session: sessionId], message, a, s')

```

---

---

**Algorithm 4** Relation of a Client  $R^c$  – Using the access token.

---

```
1: function USE_ACCESS_TOKEN(sessionId, token, tokenEPDomain, a, s')
2:   let session := s'.sessions[sessionId]
3:   let selectedAS := session[selected_AS]
4:   let rsDomain ← Doms → This domain may or may not belong to a “real” RS. If it belongs to the attacker, this request leaks the
   → Note: All paths except the mTLS and DPoP preparation endpoints are resource paths at the RS.
   access token (but no mTLS nonce, nor a DPoP proof for an honest server).
5:   let resourceID ←  $\mathbb{S}$  such that resourceID  $\notin$  {/MTLS-prepare, /DPoP-nonce}
6:   let url := ⟨URL, S, rsDomain, resourceID, ⟨⟩,  $\perp$ ⟩
7:   if s'.resourceASMapping[rsDomain][resourceID]  $\neq$  tokenEPDomain then
8:     stop → The AS from which the client received the AT is not managing the resource
   → The access token is sender-constrained, so the client must add a corresponding key proof.
9:   let clientType := s'.asAccounts[selectedAS][client_type]
10:  let clientId := s'.asAccounts[selectedAS][client_id]
11:  let body := []
12:  if clientType  $\in$  {mTLS_mTLS, pkjwt_mTLS} then → mTLS sender constraining
13:    let mtlsNonce such that (rsDomain, ⟨⟩, pubKey, mtlsNonce)  $\in$  s'.mtlsCache if possible; otherwise stop
14:    let body[TLS_binding] := mtlsNonce → This nonce is not necessarily associated with the same of the client’s keys as the access
   token. In such a case, the RS will reject this request and the client has to try again.
15:    let headers := [Authorization: [Bearer: token]] → FAPI 2.0 mandates to send access token in header
16:    let s'.mtlsCache := s'.mtlsCache  $-$  ( $\langle \rangle$  (rsDomain, ⟨⟩, pubKey, mtlsNonce))
17:  else if clientType  $\in$  {mTLS_DPoP, pkjwt_DPoP} then → DPoP sender constraining
18:    let privKey := s'.jwk → get private key
19:    let dpopNonce such that dpopNonce  $\in$  s'.dpopNonces[rsDomain] if possible; otherwise stop
20:    let s'.dpopNonces[rsDomain] := s'.dpopNonces[rsDomain]  $-$  ( $\langle \rangle$  dpopNonce)
21:    let htu := url
22:    let htu[parameters] := ⟨⟩ → [39, Sec. 4.2]: without query
23:    let htu[fragment] :=  $\perp$  → [39, Sec. 4.2]: without fragment
24:    let dpopJwt := [headers: [jwk: pub(privKey)]]
25:    let dpopJwt[payload] := [htm: POST, htu: htu, ath: hash(token), nonce: dpopNonce]
26:    let dpopProof := sig(dpopJwt, privKey)
27:    let headers := [Authorization: [DPoP: token]] → See [39, Sec. 7.1]
28:    let headers[DPoP] := dpopProof → add DPoP header; the dpopJwt can be extracted with the extractmsg() function
29:    let message := ⟨HTTPReq,  $\nu_3$ , POST, url.domain, url.path, ⟨⟩, headers, body⟩
30:    call HTTPS_SIMPLE_SEND([responseTo: RESOURCE_USAGE, session: sessionId], message, a, s')
```

---

---

**Algorithm 5** Relation of a Client  $R^c$  – Check ID Token and log user in at *c*.

---

```
1: function CHECK_ID_TOKEN(sessionId, idToken, s') → Check ID Token validity and create service session.
2:   let session := s'.sessions[sessionId] → Retrieve session data.
3:   let selectedAS := session[selected_AS]
4:   let oauthConfig := s'.oauthConfigCache[selectedAS] → Retrieve configuration for user-selected AS.
5:   let clientInfo := s'.asAccounts[selectedAS] → Retrieve client info used at that AS.
6:   let juks := s'.juksCache[selectedAS] → Retrieve signature verification key for AS.
7:   let data := extractmsg(idToken) → Extract contents of signed ID Token.
   → The following ID token checks are mandated by [70, Sec. 3.1.3.7]. Note that OIDC allows clients to skip ID token signature verification
   if the ID token is received directly from the AS (which it is here). Hence, we do not check the token’s signature (see also Line 47 of
   Algorithm 2).
8:   if data[iss]  $\neq$  selectedAS then
9:     stop → Check the issuer; note that previous checks ensure oauthConfig[issuer]  $\equiv$  selectedAS
10:  if data[aud]  $\neq$  clientInfo[client_id] then
11:    stop → Check the audience against own client id.
12:  if nonce  $\in$  session  $\wedge$  data[nonce]  $\neq$  session[nonce] then
13:    stop → If a nonce was used, check its value.
14:  let s'.sessions[sessionId][loggedInAs] := ⟨selectedAS, data[sub]⟩ → User is now logged in. Store user identity and issuer of ID
   token.
15:  let s'.sessions[sessionId][serviceSessionId] :=  $\nu_4$  → Choose a new service session id.
16:  let request := session[redirectEpRequest] → Retrieve stored meta data of the request from the browser to the redir. end-
   point in order to respond to it now. The request’s meta data was stored in
   PROCESS_HTTPS_REQUEST (Algorithm 1).
17:  let headers[Set-Cookie] := [serviceSessionId: [ $\nu_4$ , T, T, T]] → Create a cookie containing the service session id, effectively
   logging the user identified by data[sub] in at this client.
18:  let m' := encs(⟨HTTPResp, request[message].nonce, 200, headers, ok), request[key])
19:  stop ⟨⟨request[sender], request[receiver], m'⟩⟩, s'
```

---

---

**Algorithm 6** Relation of a Client  $R^c$  – Prepare and send pushed authorization request.

---

```
1: function PREPARE_AND_SEND_PAR( $sessionId, a, s'$ )
2:   let  $redirectUris := \{ \langle URL, S, d, /redirect\_ep, \langle \rangle, \perp \rangle \mid d \in \text{dom}(c) \}$   $\rightarrow$  Set of redirect URIs for all domains of  $c$ .
3:   let  $redirectUri \leftarrow redirectUris$   $\rightarrow$  Select a (potentially) different redirect URI for each authorization request
4:   let  $session := s'.sessions[sessionId]$ 
5:   let  $selectedAS := session[selected\_AS]$   $\rightarrow$  AS selected by the user at the beginning of the flow.
    $\rightarrow$  Check whether the client needs to fetch AS metadata first and do so if required.
6:   if  $selectedAS \notin s'.oauthConfigCache$  then
7:     let  $path \leftarrow \{ /.well\_known/openid-configuration, /.well\_known/oauth-authorization-server \}$ 
8:     let  $message := \langle HTTPReq, \nu_5, GET, selectedAS, path, \langle \rangle, \langle \rangle \rangle$ 
9:     call HTTPS_SIMPLE_SEND( $[responseTo: CONFIG, session: sessionId], message, a, s'$ )
10:  let  $oauthConfig := s'.oauthConfigCache[selectedAS]$ 
    $\rightarrow$  Check whether the client needs to fetch the AS's signature verification key first and do so if required.
11:  if  $selectedAS \notin s'.jwksCache$  then
12:    let  $url := oauthConfig[jwks\_uri]$ 
13:    let  $message := \langle HTTPReq, \nu_5, GET, url.host, url.path, url.parameters, \langle \rangle, \langle \rangle \rangle$ 
14:    call HTTPS_SIMPLE_SEND( $[responseTo: JWKS, session: sessionId], message, a, s'$ )
    $\rightarrow$  Construct pushed authorization request
15:  let  $parEndpoint := oauthConfig[par\_ep]$ 
16:  let  $clientId := s'.asAccounts[selectedAS][client\_id]$ 
17:  let  $clientType := s'.asAccounts[selectedAS][client\_type]$ 
18:  if  $clientType \in \{ mTLS\_mTLS, mTLS\_DPoP \}$  then  $\rightarrow$  mTLS client authentication
19:    let  $mtlsNonce$  such that  $(parEndpoint.host, clientId, \langle \rangle, mtlsNonce) \in s'.mtlsCache$  if possible; otherwise stop
20:    let  $authData := [TLS\_AuthN: mtlsNonce]$ 
21:    let  $s'.mtlsCache := s'.mtlsCache - \langle \rangle (parEndpoint.host, clientId, \langle \rangle, mtlsNonce)$ 
22:  else if  $clientType \in \{ pkjwt\_mTLS, pkjwt\_DPoP \}$  then  $\rightarrow$  private_key_jwt client authentication
23:    let  $jwt := [iss: clientId, sub: clientId, aud: selectedAS]$ 
24:    let  $jws := sig(jwt, s'.jwk)$ 
25:    let  $authData := [client\_assertion: jws]$ 
26:  let  $pkceVerifier := \nu_{pkce}$   $\rightarrow$  Fresh random value
27:  let  $pkceChallenge := hash(pkceVerifier)$ 
28:  let  $parData := [response\_type: code, code\_challenge\_method: S256, client\_id: clientId,$ 
    $\hookrightarrow$   $redirect\_uri: redirectUri, code\_challenge: pkceChallenge]$ 
29:  let  $useOidc \leftarrow \{ \top, \perp \}$   $\rightarrow$  Use of OIDC is optional
30:  if  $useOidc \equiv \top$  then
31:    let  $parData[scope] := openid$ 
32:  let  $s'.sessions[sessionId] := s'.sessions[sessionId] + \langle \rangle parData$ 
33:  let  $parData := parData + \langle \rangle authData$ 
34:  let  $s'.sessions[sessionId][code\_verifier] := pkceVerifier$   $\rightarrow$  Store PKCE randomness in state
35:  let  $authzReq := \langle HTTPReq, \nu_{parNonce}, POST, parEndpoint.host, parEndpoint.path, parEndpoint.parameters, \langle \rangle, parData \rangle$ 
36:  call HTTPS_SIMPLE_SEND( $[responseTo: PAR, session: sessionId], authzReq, a, s'$ )
```

---

---

**Algorithm 7** Relation of a Client  $R^c$  – Handle trigger events.

---

```
1: function PROCESS_TRIGGER( $a, s'$ )
2:   let  $action \leftarrow \{ \text{MTLS\_PREPARE\_AS}, \text{MTLS\_PREPARE\_RS}, \text{MTLS\_PREPARE\_MISCONFIGURED\_TOKEN\_EP},$ 
    $\hookrightarrow \text{GET\_DPOP\_NONCE} \}$ 
3:   switch  $action$  do
4:     case  $\text{MTLS\_PREPARE\_AS}$ 
5:       let  $server \leftarrow \text{Doms}$  such that  $server \in s'.\text{asAccounts}$  if possible; otherwise stop
6:       let  $asAcc := s'.\text{asAccounts}[server]$ 
7:       let  $clientId := asAcc[\text{client\_id}]$ 
8:       let  $body := [\text{client\_id}: clientId]$ 
9:       let  $message := \langle \text{HTTPReq}, \nu_{\text{mtls}}, \text{GET}, server, /\text{MTLS-prepare}, \langle \rangle, \langle \rangle, body \rangle$ 
10:      call  $\text{HTTPS\_SIMPLE\_SEND}([\text{responseTo}: \text{MTLS}, \text{client\_id}: clientId], message, a, s')$ 
11:     case  $\text{MTLS\_PREPARE\_RS}$ 
12:       let  $resourceServer \leftarrow \text{Doms}$   $\rightarrow$  Note: This may or may not be a “real” RS.
13:       let  $domainAndKey \leftarrow s'.\text{tlskeys}$ 
14:       let  $pubKey := \text{pub}(domainAndKey.2)$ 
15:       let  $body := [\text{pub\_key}: pubKey]$ 
16:       let  $message := \langle \text{HTTPReq}, \nu_{\text{mtls}}, \text{GET}, resourceServer, /\text{MTLS-prepare}, \langle \rangle, \langle \rangle, body \rangle$ 
17:       call  $\text{HTTPS\_SIMPLE\_SEND}([\text{responseTo}: \text{MTLS}, \text{pub\_key}: pubKey], message, a, s')$ 
18:     case  $\text{MTLS\_PREPARE\_MISCONFIGURED\_TOKEN\_EP}$ 
    $\rightarrow$  This case allows the client to retrieve mTLS nonces from attacker-controlled servers and subsequently make requests to such
   servers. Without this case, the model would not capture attacks in which the client talks to attacker-controlled endpoints
   protected by mTLS.
19:     let  $server \leftarrow \text{Doms}$  such that  $server \in s'.\text{asAccounts}$  if possible; otherwise stop
20:     let  $asAcc := s'.\text{asAccounts}[server]$ 
21:     let  $clientId := asAcc[\text{client\_id}]$ 
22:     let  $host \leftarrow \text{Doms}$   $\rightarrow$  Non-deterministically choose the domain instead of sending to the correct AS
23:     let  $body := [\text{client\_id}: clientId]$ 
24:     let  $message := \langle \text{HTTPReq}, \nu_{\text{mtls}}, \text{GET}, host, /\text{MTLS-prepare}, \langle \rangle, \langle \rangle, body \rangle$ 
25:     call  $\text{HTTPS\_SIMPLE\_SEND}([\text{responseTo}: \text{MTLS}, \text{client\_id}: clientId], message, a, s')$ 
26:     case  $\text{GET\_DPOP\_NONCE}$ 
    $\rightarrow$  Our client uses DPoP server-provided nonces at the RS. The RS model offers a special endpoint to retrieve nonces.
27:     let  $resourceServer \leftarrow \text{Doms}$   $\rightarrow$  Note: This may or may not be a “real” RS.
28:     let  $message := \langle \text{HTTPReq}, \nu_{\text{DPoPreq}}, \text{GET}, resourceServer, /\text{DPoP-nonce}, \langle \rangle, \langle \rangle, \langle \rangle \rangle$ 
29:     call  $\text{HTTPS\_SIMPLE\_SEND}([\text{responseTo}: \text{DPOP\_NONCE}], message, a, s')$ 
30:   stop
```

---

---

**Algorithm 8** Relation of  $script\_client\_index$ 

---

**Input:**  $(tree, docnonce, scriptstate, scriptinputs, cookies, localStorage, sessionStorage, ids, secrets)$   $\rightarrow$  **Script that models the index page of a client.** Users can initiate the login flow or follow arbitrary links. The script receives various information about the current browser state, filtered according to the access rules (same origin policy and others) in the browser.

```
1: let  $switch \leftarrow \{ \text{auth}, \text{link} \}$   $\rightarrow$  Non-deterministically decide whether to start a login flow or to follow some link.
2: if  $switch \equiv \text{auth}$  then  $\rightarrow$  Start login flow.
3:   let  $url := \text{GETURL}(tree, docnonce)$   $\rightarrow$  Retrieve URL of current document.
4:   let  $id \leftarrow ids$   $\rightarrow$  Retrieve one of user’s identities.
5:   let  $as := id.\text{domain}$   $\rightarrow$  Extract domain of AS from chosen  $id$ .
6:   let  $url' := \langle \text{URL}, \text{S}, url.\text{host}, /startLogin, \langle \rangle, \perp \rangle$   $\rightarrow$  Assemble request URL.
7:   let  $command := \langle \text{FORM}, url', \text{POST}, as, \perp \rangle$   $\rightarrow$  Post a form including the selected AS to the client.
8:   stop  $(s, cookies, localStorage, sessionStorage, command)$   $\rightarrow$  Finish script’s run and instruct the browser to execute the command
   (i.e., to POST the form).
9: else  $\rightarrow$  Follow (random) link to facilitate referrer-based attacks.
10: let  $protocol \leftarrow \{ \text{P}, \text{S} \}$   $\rightarrow$  Non-deterministically select protocol (HTTP or HTTPS).
11: let  $host \leftarrow \text{Doms}$   $\rightarrow$  Non-det. select host.
12: let  $path \leftarrow \text{S}$   $\rightarrow$  Non-det. select path.
13: let  $fragment \leftarrow \text{S}$   $\rightarrow$  Non-det. select fragment part.
14: let  $parameters \leftarrow [\text{S} \times \text{S}]$   $\rightarrow$  Non-det. select parameters.
15: let  $url := \langle \text{URL}, protocol, host, path, parameters, fragment \rangle$   $\rightarrow$  Assemble request URL.
16: let  $command := \langle \text{HREF}, url, \perp, \perp \rangle$   $\rightarrow$  Follow link to the selected URL.
17: stop  $(s, cookies, localStorage, sessionStorage, command)$   $\rightarrow$  Finish script’s run and instruct the browser to execute the command
   (follow link).
```

---



### K. Authorization Servers

An authorization server  $as \in AS$  is a Web server modeled as an atomic process  $(I^{as}, Z^{as}, R^{as}, s_0^{as})$  with the addresses  $I^{as} := \text{addr}(as)$ . Next, we define the set  $Z^{as}$  of states of  $as$  and the initial state  $s_0^{as}$  of  $as$ .

*Definition 17.* A state  $s \in Z^{as}$  of an AS  $as$  is a term of the form  $\langle \text{DNSAddress}, \text{pendingDNS}, \text{pendingRequests}, \text{corrupt}, \text{keyMapping}, \text{tlskeys}, \text{jwt}, \text{registrationRequests}, \text{clients}, \text{records}, \text{authorizationRequests}, \text{mtlsRequests}, \text{rsCredentials} \rangle$  with:  $\text{DNSAddress} \in \text{IPs}$ ,  $\text{pendingDNS} \in [\mathcal{N} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{pendingRequests} \in \mathcal{T}_{\mathcal{N}}$ ,  $\text{corrupt} \in \mathcal{T}_{\mathcal{N}}$ ,  $\text{keyMapping} \in [\text{Doms} \times \mathcal{T}_{\mathcal{N}}]$  (as in [Definition 71](#)),  $\text{tlskeys} \in [\text{Doms} \times K_{\text{TLS}}]$ ,  $\text{jwt} \in K_{\text{sign}}$ ,  $\text{registrationRequests} \in \mathcal{T}_{\mathcal{N}}$ ,  $\text{clients} \in [\mathcal{T}_{\mathcal{N}} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{records} \in \mathcal{T}_{\mathcal{N}}$ ,  $\text{authorizationRequests} \in [\mathcal{T}_{\mathcal{N}} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{mtlsRequests} \in [\mathcal{T}_{\mathcal{N}} \times \mathcal{T}_{\mathcal{N}}]$ , and  $\text{rsCredentials} \in \mathcal{T}_{\mathcal{N}}$ .

An initial state  $s_0^{as}$  of  $as$  is a state of  $as$  with

- $s_0^{as}.\text{DNSAddress} \in \text{IPs}$ ,
- $s_0^{as}.\text{pendingDNS} \equiv \langle \rangle$ ,
- $s_0^{as}.\text{pendingRequests} \equiv \langle \rangle$ ,
- $s_0^{as}.\text{corrupt} \equiv \perp$ ,
- $s_0^{as}.\text{keyMapping}$  being the same as the keymapping for browsers,
- $s_0^{as}.\text{tlskeys} \equiv \text{tlskeys}^{as}$ ,
- $s_0^{as}.\text{jwt} \equiv \text{signkey}(as)$  (see [Appendix D-C](#)),
- $s_0^{as}.\text{registrationRequests} \equiv \langle \rangle$ ,
- $s_0^{as}.\text{clients} \equiv \text{clientInfoAS}(as)$  (see [Definition 14](#)),
- $s_0^{as}.\text{records} \equiv \langle \rangle$ ,
- $s_0^{as}.\text{authorizationRequests} \equiv \langle \rangle$ ,
- $s_0^{as}.\text{mtlsRequests} \equiv \langle \rangle$ , and
- $s_0^{as}.\text{rsCredentials} \equiv \text{rsCreds}$  where  $\text{rsCreds}$  is a sequence and  $\forall c: c \in \text{rsCreds} \Leftrightarrow (\exists d \in \text{dom}(as), rs \in \text{Doms}: c \equiv \text{secretOfRS}(d, rs))$ .

We now specify the relation  $R^{as}$ : This relation is based on the model of generic HTTPS servers (see [Appendix G-L](#)). We specify algorithms that differ from or do not exist in the generic server model in [Algorithms 9](#) to [10](#). [Algorithm 11](#) shows the script `script_as_form` that is used by ASs.

---

**Algorithm 9** Relation of AS  $R^{AS}$  – Processing HTTPS Requests

---

```
1: function PROCESS_HTTPS_REQUEST( $m, k, a, f, s'$ )
2:   if  $m.path \equiv /.well-known/openid-configuration \vee$ 
    $\hookrightarrow m.path \equiv /.well-known/oauth-authorization-server$  then  $\rightarrow$  We model both OIDD and RFC 8414.
3:     let  $metaData := [issuer: m.host]$ 
4:     let  $metaData[auth\_ep] := \langle URL, S, m.host, /auth, \langle \rangle, \perp \rangle$ 
5:     let  $metaData[token\_ep] := \langle URL, S, m.host, /token, \langle \rangle, \perp \rangle$ 
6:     let  $metaData[par\_ep] := \langle URL, S, m.host, /par, \langle \rangle, \perp \rangle$ 
7:     let  $metaData[introspec\_ep] := \langle URL, S, m.host, /introspect, \langle \rangle, \perp \rangle$ 
8:     let  $metaData[jwks\_uri] := \langle URL, S, m.host, /jwks, \langle \rangle, \perp \rangle$ 
9:     let  $m' := enc_s(\langle \langle HTTPResp, m.nonce, 200, \langle \rangle, metaData \rangle, k \rangle)$ 
10:    stop  $\langle \langle f, a, m' \rangle \rangle, s'$ 
11:  else if  $m.path \equiv /jwks$  then
12:    let  $m' := enc_s(\langle \langle HTTPResp, m.nonce, 200, \langle \rangle, pub(s'.jwk) \rangle, k \rangle)$ 
13:    stop  $\langle \langle f, a, m' \rangle \rangle, s'$ 
14:  else if  $m.path \equiv /auth$  then  $\rightarrow$  Authorization endpoint: Reply with login page.
15:    if  $m.method \equiv GET$  then
16:      let  $data := m.parameters$ 
17:    else if  $m.method \equiv POST$  then
18:      let  $data := m.body$ 
19:    let  $requestUri := data[request\_uri]$ 
20:    if  $requestUri \equiv \langle \rangle$  then
21:      stop  $\rightarrow$  FAPI 2.0 mandates PAR, therefore a request URI is required
22:    let  $authzRecord := s'.authorizationRequests[requestUri]$ 
23:    let  $clientId := data[client\_id]$ 
24:    if  $authzRecord[client\_id] \neq clientId$  then  $\rightarrow$  Check binding of request URI to client
25:      stop
26:    if  $clientId \notin s'.clients$  then
27:      stop  $\rightarrow$  Unknown client
28:    let  $s'.authorizationRequests[requestUri][auth2\_reference] := \nu_5$ 
29:    let  $m' := enc_s(\langle \langle HTTPResp, m.nonce, 200, \langle \langle ReferrerPolicy, origin \rangle \rangle, \langle script\_as\_form, [auth2\_reference: \nu_5] \rangle \rangle, k \rangle)$ 
30:    stop  $\langle \langle f, a, m' \rangle \rangle, s'$ 
31:  else if  $m.path \equiv /auth2 \wedge m.method \equiv POST \wedge m.headers[Origin] \equiv \langle m.host, S \rangle$  then  $\rightarrow$  Second step of authorization
32:    let  $identity := m.body[identity]$ 
33:    let  $password := m.body[password]$ 
34:    if  $identity.domain \notin dom(as)$  then
35:      stop  $\rightarrow$  This AS does not manage  $identity$ 
36:    if  $password \neq secretOfID(identity)$  then
37:      stop  $\rightarrow$  Invalid user credentials
38:    let  $auth2Reference := m.body[auth2\_reference]$ 
39:    let  $requestUri$  such that  $s'.authorizationRequests[requestUri][auth2\_reference] \equiv auth2Reference$ 
    $\hookrightarrow$  if possible; otherwise stop
40:    let  $authzRecord := s'.authorizationRequests[requestUri]$ 
41:    let  $authzRecord[subject] := identity$ 
42:    let  $authzRecord[issuer] := m.host$ 
43:    let  $authzRecord[code] := \nu_1$   $\rightarrow$  Generate a fresh, random authorization code
44:    let  $s'.records := s'.records + \langle \rangle authzRecord$ 
45:    let  $responseData := [code: authzRecord[code]]$ 
46:    if  $authzRecord[state] \neq \langle \rangle$  then
47:      let  $responseData[state] := authzRecord[state]$ 
48:    let  $redirectUri := authzRecord[redirect\_uri]$ 
49:    let  $redirectUri.parameters := redirectUri.parameters \cup responseData$ 
50:    let  $redirectUri.parameters[iss] := authzRecord[issuer]$ 
51:    let  $m' := enc_s(\langle \langle HTTPResp, m.nonce, 303, \langle \langle Location, redirectUri \rangle \rangle, \langle \rangle \rangle, k \rangle)$ 
52:    stop  $\langle \langle f, a, m' \rangle \rangle, s'$ 
53:  else if  $m.path \equiv /par \wedge m.method \equiv POST$  then  $\rightarrow$  Pushed Authorization Request
54:    if  $m.body[response\_type] \neq code \vee m.body[code\_challenge\_method] \neq S256$  then
55:      stop
56:    let  $authnResult := AUTHENTICATE_CLIENT(m, s')$   $\rightarrow$  Stops in case of errors/failed authentication
57:    let  $clientId := authnResult.1$ 
58:    let  $s' := authnResult.2$ 
59:    let  $mtlsInfo := authnResult.3$ 
60:    if  $clientId \neq m.body[client\_id]$  then
61:      stop  $\rightarrow$  Key used in client authentication is not registered for  $m.body[client\_id]$ 
62:    let  $redirectUri := m.body[redirect\_uri]$   $\rightarrow$  Clients are required to send  $redirect\_uri$  with each request
```

$\rightarrow$  Algorithm continues on next page.

---

---

```

63:   if redirectUri ≡ ⟨⟩ then
64:     stop
65:   if redirectUri.protocol ≠ S then
66:     stop
67:   let codeChallenge := m.body[code_challenge] → PKCE challenge
68:   if codeChallenge ≡ ⟨⟩ then
69:     stop → Missing PKCE challenge
70:   let requestUri := νA → Choose random URI
71:   let authzRecord := [client_id: clientId]
72:   let authzRecord[state] := m.body[state]
73:   let authzRecord[scope] := m.body[scope]
74:   if nonce ∈ m.body then
75:     let authzRecord[nonce] := m.body[nonce]
76:   let authzRecord[redirect_uri] := redirectUri
77:   let authzRecord[code_challenge] := codeChallenge
78:   let s'.authorizationRequests[requestUri] := authzRecord → Store data linked to requestUri
79:   let m' := encs((HTTPResp, m.nonce, 201, ⟨⟩, [request_uri: requestUri]), k)
80:   stop ⟨⟨f, a, m'⟩⟩, s'
81: else if m.path ≡ /token ∧ m.method ≡ POST then
82:   if m.body[grant_type] ≠ authorization_code then
83:     stop
84:   let authnResult := AUTHENTICATE_CLIENT(m, s') → Stops in case of errors/failed authentication
85:   let clientId := authnResult.1
86:   let s' := authnResult.2
87:   let mtlInfo := authnResult.3
88:   let code := m.body[code]
89:   let codeVerifier := m.body[code_verifier]
90:   if code ≡ ⟨⟩ ∨ codeVerifier ≡ ⟨⟩ then
91:     stop → Missing code or code_verifier
92:   let record, ptr such that record ≡ s'.records.ptr ∧ record[code] ≡ code ∧ code ≠ ⊥ ∧ ptr ∈ ℕ if possible; otherwise stop
93:   if record[client_id] ≠ clientId then
94:     stop
95:   if record[code_challenge] ≠ hash(codeVerifier) ∨ record[redirect_uri] ≠ m.body[redirect_uri] then
96:     stop → PKCE verification failed or URI mismatch
97:   let clientType := s'.clients[clientId][client_type]
98:   if clientType ≡ pkjwt_DPoP ∨ clientType ≡ mTLS_DPoP then → DPoP token binding
99:     let tokenType := DPoP
100:    let dpopProof := m.headers[DPoP]
101:    let dpopJwt := extractmsg(dpopProof)
102:    let verificationKey := dpopJwt[headers][jwk]
103:    if checksig(dpopProof, verificationKey) ≠ ⊤ ∨ verificationKey ≡ ⟨⟩ then
104:      stop → Invalid DPoP signature (or empty jwk header)
105:    let dpopClaims := dpopJwt[payload]
106:    let reqUri := ⟨URL, S, m.host, m.path, ⟨⟩, ⊥⟩
107:    if dpopClaims[htm] ≠ m.method ∨ dpopClaims[htu] ≠ reqUri then
108:      stop → DPoP claims do not match corresponding message
109:    let cnfContent := [jkt: hash(verificationKey)]
110:  else if clientType ≡ pkjwt_mTLS ∨ clientType ≡ mTLS_mTLS then → mTLS token binding
111:    let tokenType := Bearer
112:    let mtlNonce := m.body[TLS_binding]
113:    if clientType ≡ mTLS_mTLS then → Client used mTLS authentication, reuse data from authentication
114:      if mtlNonce ≠ mtlInfo.1 then
115:        stop → Client tried to use different mTLS key for authentication and token binding
116:      else → Client did not use mTLS authentication
117:        let mtlInfo such that mtlInfo ∈ s'.mtlsRequests[clientId] ∧ mtlInfo.1 ≡ mtlNonce if possible; otherwise stop
118:        let s'.mtlsRequests[clientId] := s'.mtlsRequests[clientId] - ⟨⟩ mtlInfo
119:        let mTlsKey := mtlInfo.2 → mTLS public key of client
120:        let cnfContent := [x5t#S256: hash(mTlsKey)]
121:  else
122:    stop → Client used neither DPoP nor mTLS
123:  let s'.records.ptr[code] := ⊥ → Invalidate code
124:  let atType ← {JWT, opaque} → The AS chooses randomly whether it issues a structured or an opaque access token

```

→ Algorithm continues on next page.

---

---

```

125:   if atType  $\equiv$  JWT then  $\rightarrow$  Structured access token
126:     let accessTokenContent := [cnf: cnfContent, sub: record[subject]]
127:     let accessToken := sig(accessTokenContent, s'.jwk)
128:   else  $\rightarrow$  Opaque access token
129:     let accessToken :=  $\nu_2$   $\rightarrow$  Fresh random value
130:   let s'.records.ptr[access_token] := accessToken  $\rightarrow$  Store for token introspection
131:   let s'.records.ptr[cnf] := cnfContent  $\rightarrow$  Store for token introspection
132:   let body := [access_token: accessToken, token_type: tokenType]
133:   if record[scope]  $\equiv$  openid then  $\rightarrow$  Client requested ID token
134:     let idTokenBody := [iss: record[issuer]]
135:     let idTokenBody[sub] := record[subject]
136:     let idTokenBody[aud] := record[client_id]
137:     if nonce  $\in$  record then
138:       let idTokenBody[nonce] := record[nonce]
139:     let idToken := sig(idTokenBody, s'.jwk)
140:     let body[id_token] := idToken
141:   let m' := encs(⟨HTTPResp, m.nonce, 200, ⟨⟩, body⟩, k)
142:   stop ⟨⟨f, a, m'⟩⟩, s'
143: else if m.path  $\equiv$  /introspect  $\wedge$  m.method  $\equiv$  POST  $\wedge$  token  $\in$  m.body then
144:   let rsCredentials such that ⟨Basic, rsCredentials⟩  $\equiv$  m.headers[Authorization] if possible; otherwise stop
145:   if rsCredentials  $\notin$  s'.rsCredentials then
146:     stop  $\rightarrow$  RS authentication failed
147:   let token := m.body[token]
148:   let record such that record  $\in$  s'.records  $\wedge$  record[access_token]  $\equiv$  token if possible; otherwise let record :=  $\diamond$ 
149:   if record  $\equiv$   $\diamond$  then  $\rightarrow$  Unknown token
150:     let m' := encs(⟨HTTPResp, m.nonce, 200, ⟨⟩, [active:  $\perp$ ], k)
151:   else  $\rightarrow$  token was issued by this AS
152:     let body := [active:  $\top$ , cnf: record[cnf], sub: record[subject]]  $\rightarrow$  cnf claim contains hash of token binding key
153:     let m' := encs(⟨HTTPResp, m.nonce, 200, ⟨⟩, body⟩, k)
154:   stop ⟨⟨f, a, m'⟩⟩, s'
155: else if m.path  $\equiv$  /MTLS-prepare then  $\rightarrow$  See Appendix D-H
156:   let clientId := m.body[client_id]
157:   let mtlsNonce :=  $\nu_3$ 
158:   let clientKey := s'.clients[clientId][mtls_key]
159:   if clientKey  $\equiv$  ⟨⟩  $\vee$  clientKey  $\equiv$  pub( $\diamond$ ) then
160:     stop  $\rightarrow$  Client has no mTLS key
161:   let s'.mtlsRequests[clientId] := s'.mtlsRequests[clientId] +  $\langle \rangle$  ⟨mtlsNonce, clientKey⟩
162:   let m' := encs(⟨HTTPResp, m.nonce, 200, ⟨⟩, enca(⟨mtlsNonce, s'.keyMapping[m.host], clientKey⟩), k)
163:   stop ⟨⟨f, a, m'⟩⟩, s'
164: stop  $\rightarrow$  Request was malformed or sent to non-existing endpoint.

```

---

---

**Algorithm 10** Relation of AS  $R^{\text{as}}$  – Client Authentication

---

```
1: function AUTHENTICATE_CLIENT( $m, s'$ )  $\rightarrow$  Check client authentication in message  $m$ . Stops the current processing step in case of errors or failed authentication.
2:   if  $\text{client\_assertion} \in m.\text{body}$  then  $\rightarrow$  private_key_jwt client authentication
3:     let  $\text{jwts} := m.\text{body}[\text{client\_assertion}]$ 
4:     let  $\text{clientId}, \text{verificationKey}$  such that  $\text{verificationKey} \equiv s'.\text{clients}[\text{clientId}][\text{jwt\_key}] \wedge$   

 $\hookrightarrow \text{checksig}(\text{jwts}, \text{verificationKey}) \equiv \top$  if possible; otherwise stop
5:     if  $\text{verificationKey} \equiv \langle \rangle \vee \text{verificationKey} \equiv \text{pub}(\diamond)$  then
6:       stop  $\rightarrow$  Client has no jwt key
7:     let  $\text{clientInfo} := s'.\text{clients}[\text{clientId}]$ 
8:     let  $\text{clientType} := \text{clientInfo}[\text{client\_type}]$ 
9:     if  $\text{clientType} \not\equiv \text{pkjwt\_mTLS} \wedge \text{clientType} \not\equiv \text{pkjwt\_DPoP}$  then
10:      stop  $\rightarrow$  Client authentication type mismatch
11:     let  $\text{jwt} := \text{extractmsg}(\text{jwts})$ 
12:     if  $\text{jwt}[\text{iss}] \not\equiv \text{clientId} \vee \text{jwt}[\text{sub}] \not\equiv \text{clientId}$  then
13:       stop
14:     if  $\text{jwt}[\text{aud}] \not\equiv \langle \text{URL}, \text{S}, m.\text{host}, /token, \langle \rangle, \perp \rangle \wedge \text{jwt}[\text{aud}] \not\equiv m.\text{host}$   $\rightarrow$  issuer in AS metadata is just the host part  

 $\hookrightarrow \wedge \text{jwt}[\text{aud}] \not\equiv \langle \text{URL}, \text{S}, m.\text{host}, /par, \langle \rangle, \perp \rangle$  then
15:       stop  $\rightarrow$  aud claim value is neither token, nor PAR endpoint nor AS issuer identifier
16:   else if  $\text{TLS\_AuthN} \in m.\text{body}$  then  $\rightarrow$  mTLS client authentication
17:     let  $\text{clientId} := m.\text{body}[\text{client\_id}]$   $\rightarrow$  [11] mandates client_id when using mTLS authentication
18:     let  $\text{mtlsNonce} := m.\text{body}[\text{TLS\_AuthN}]$ 
19:     let  $\text{mtlsInfo}$  such that  $\text{mtlsInfo} \in s'.\text{mtlsRequests}[\text{clientId}] \wedge \text{mtlsInfo}.1 \equiv \text{mtlsNonce}$  if possible; otherwise stop
20:     let  $\text{clientInfo} := s'.\text{clients}[\text{clientId}]$ 
21:     let  $\text{clientType} := \text{clientInfo}[\text{client\_type}]$ 
22:     if  $\text{clientType} \not\equiv \text{mTLS\_mTLS} \wedge \text{clientType} \not\equiv \text{mTLS\_DPoP}$  then
23:       stop  $\rightarrow$  Client authentication type mismatch
24:     let  $s'.\text{mtlsRequests}[\text{clientId}] := s'.\text{mtlsRequests}[\text{clientId}] - \langle \rangle \text{mtlsInfo}$ 
25:   else
26:     stop  $\rightarrow$  Unsupported client (authentication) type
27:   if  $\text{clientType} \equiv \text{mTLS\_mTLS} \vee \text{clientType} \equiv \text{mTLS\_DPoP}$  then
28:     return  $\langle \text{clientId}, s', \text{mtlsInfo} \rangle$ 
29:   else
30:     return  $\langle \text{clientId}, s', \perp \rangle$   $\rightarrow$  private_key_jwt client authentication, i.e., no mTLS info
```

---

---

**Algorithm 11** Relation of  $\text{script\_as\_form}$ : A login page for the user.

---

**Input:**  $\langle \text{tree}, \text{docnonce}, \text{scriptstate}, \text{scriptinputs}, \text{cookies}, \text{localStorage}, \text{sessionStorage}, \text{ids}, \text{secrets} \rangle$

```
1: let  $\text{url} := \text{GETURL}(\text{tree}, \text{docnonce})$ 
2: let  $\text{url}' := \langle \text{URL}, \text{S}, \text{url}.\text{host}, /auth2, \langle \rangle, \perp \rangle$ 
3: let  $\text{formData} := \text{scriptstate}$ 
4: let  $\text{identity} \leftarrow \text{ids}$ 
5: let  $\text{secret} \leftarrow \text{secrets}$ 
6: let  $\text{formData}[\text{identity}] := \text{identity}$ 
7: let  $\text{formData}[\text{password}] := \text{secret}$ 
8: let  $\text{command} := \langle \text{FORM}, \text{url}', \text{POST}, \text{formData}, \perp \rangle$ 
9: stop  $\langle s, \text{cookies}, \text{localStorage}, \text{sessionStorage}, \text{command} \rangle$ 
```

---

### L. Resource Servers

A resource server  $rs \in \text{RS}$  is a Web server modeled as an atomic process  $(I^{rs}, Z^{rs}, R^{rs}, s_0^{rs})$  with the addresses  $I^{rs} := \text{addr}(rs)$ . The set of states  $Z^{rs}$  and the initial state  $s_0^{rs}$  of  $rs$  are defined in the following.

*Definition 18.* A state  $s \in Z^{rs}$  of a resource server  $rs$  is a term of the form  $\langle \text{DNSaddress}, \text{pendingDNS}, \text{pendingRequests}, \text{corrupt}, \text{keyMapping}, \text{tlskeys}, \text{mtlsRequests}, \text{pendingResponses}, \text{resourceNonces}, \text{ids}, \text{asInfo}, \text{resourceASMapping}, \text{dpopNonces} \rangle$  with  $\text{DNSaddress} \in \text{IPs}$ ,  $\text{pendingDNS} \in [\mathcal{N} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{pendingRequests} \in [\mathcal{N} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{corrupt} \in \mathcal{T}_{\mathcal{N}}$ ,  $\text{keyMapping} \in [\text{Doms} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{tlskeys} \in [\text{Doms} \times K_{\text{TLS}}]$  (all former components as in [Definition 71](#)),  $\text{mtlsRequests} \in \mathcal{T}_{\mathcal{N}}$ ,  $\text{pendingResponses} \in \mathcal{T}_{\mathcal{N}}$ ,  $\text{resourceNonces} \in [\text{ID} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{ids} \subset \text{ID}$ ,  $\text{asInfo} \in [\text{Doms} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{resourceASMapping} \in [\text{resourceURLPath}^{rs} \times \mathcal{T}_{\mathcal{N}}]$ , and  $\text{dpopNonces} \in \mathcal{T}_{\mathcal{N}}$ .

An initial state  $s_0^{rs}$  of  $rs$  is a state of  $rs$  with

- $s_0^{rs}.\text{DNSaddress} \in \text{IPs}$ ,
- $s_0^{rs}.\text{pendingDNS} \equiv \langle \rangle$ ,
- $s_0^{rs}.\text{pendingRequests} \equiv \langle \rangle$ ,
- $s_0^{rs}.\text{corrupt} \equiv \perp$ ,
- $s_0^{rs}.\text{keyMapping}$  being the same as the keymapping for browsers,
- $s_0^{rs}.\text{tlskeys} \equiv \text{tlskeys}^{rs}$ ,
- $s_0^{rs}.\text{mtlsRequests} \equiv \langle \rangle$ ,
- $s_0^{rs}.\text{pendingResponses} \equiv \langle \rangle$ ,
- $s_0^{rs}.\text{resourceNonces}$  being a dictionary where the RS stores the resource nonces for each identity and resource id pair, initialized as  $s_0^{rs}.\text{resourceNonces}[id][\text{resourceID}] := \langle \rangle$ ,  $\forall id \in \langle \rangle s_0^{rs}.\text{ids}, \forall \text{resourceID} \in \mathbb{S}$
- $s_0^{rs}.\text{ids} \subset \langle \rangle \langle \text{ID} \rangle$  such that  $\forall id \in s_0^{rs}.\text{ids} : \text{governor}(id) \in \text{supportedAuthorizationServer}^{rs}$ , i.e., the RS manages only resources of identities that are governed by one of the AS supported by the RS,
- and for each domain of a supported AS  $\text{dom}_{as} \in \text{supportedAuthorizationServerDoms}^{rs}$ , let  $s_0^{rs}.\text{asInfo}$  contain a dictionary entry with the following values:
  - $s_0^{rs}.\text{asInfo}[\text{dom}_{as}][\text{as\_introspect\_ep}] \equiv \langle \text{URL}, \text{S}, \text{dom}_{as}, /introspect, \langle \rangle, \perp \rangle$  (the URL of the introspection endpoint of the AS)
  - $s_0^{rs}.\text{asInfo}[\text{dom}_{as}][\text{as\_key}] \equiv \text{signkey}(\text{dom}^{-1}(\text{dom}_{as}))$  being the verification key for the AS
  - $s_0^{rs}.\text{asInfo}[\text{dom}_{as}][\text{rs\_credentials}]$  being a sequence s.t.  $\forall c: c \in \langle \rangle s_0^{rs}.\text{asInfo}[\text{dom}_{as}][\text{rs\_credentials}] \Leftrightarrow (\exists \text{rsDom} \in \text{dom}(rs) : c \equiv \text{secretOfRS}(\text{dom}_{as}, \text{rsDom}))$ , i.e., the secrets used by the RS for authenticating at the AS.
- $s_0^{rs}.\text{resourceASMapping}[\text{resourceID}] \in \text{dom}(\text{authorizationServerOfResource}^{rs}(\text{resourceID}))$ ,  $\forall \text{resourceID} \in \text{resourceURLPath}^{rs}$  (a domain of the AS managing the resource identified by  $\text{resourceID}$ ),
- $s_0^{rs}.\text{dpopNonces} \equiv \langle \rangle$

The relation  $R^{rs}$  is again based on the generic HTTPS server model (see [Appendix G-L](#)), for which the algorithms used for processing HTTP requests and responses are defined in [Algorithm 12](#) and [Algorithm 13](#).

---

**Algorithm 12** Relation of RS  $R^{rs}$  – Processing HTTPS Requests

---

```
1: function PROCESS_HTTPS_REQUEST( $m, k, a, f, s'$ )
2:   if  $m.path \equiv /MTLS-prepare$  then
3:     let  $mtlsNonce := \nu_1$ 
4:     let  $clientKey := m.body[pub\_key]$   $\rightarrow$  Certificate is not required to be checked [11, Section 4.2]
5:     let  $s'.mtlsRequests := s'.mtlsRequests + \langle \langle mtlNonce, clientKey \rangle \rangle$ 
6:     let  $m' := enc_s(\langle \langle HTTPResp, m.nononce, 200, \langle \rangle, enc_a(\langle \langle mtlNonce, s'.keyMapping[m.host] \rangle, clientKey) \rangle \rangle, k)$ 
7:     stop  $\langle \langle f, a, m' \rangle \rangle, s'$ 
8:   else if  $m.path \equiv /DPoP-nonce$  then
9:     let  $freshDpopNonce := \nu_{dpop}$ 
10:    let  $s'.dpopNonces := s'.dpopNonces + \langle \langle freshDpopNonce \rangle \rangle$ 
11:    let  $m' := enc_s(\langle \langle HTTPResp, m.nononce, 200, \langle \rangle, [nonce: freshDpopNonce] \rangle \rangle, k)$ 
12:    stop  $\langle \langle f, a, m' \rangle \rangle, s'$ 
13:   else
14:     let  $resourceID := m.path$ 
15:     let  $responsibleAS := s'.resourceASMapping[resourceID]$ 
16:     if  $responsibleAS \equiv \langle \rangle$  then
17:       stop  $\rightarrow$  Resource is not managed by any of the supported ASs
18:     let  $asInfo := s'.asInfo[responsibleAS]$ 
19:     if Authorization  $\in m.headers$  then
20:       let  $authnScheme := m.headers[Authorization].1$ 
21:       let  $accessToken := m.headers[Authorization].2$ 
22:       if  $authnScheme \equiv Bearer$  then  $\rightarrow$  mTLS sender constraining
23:         let  $mtlsNonce := m.body[TLS\_binding]$ 
24:         let  $mtlsInfo$  such that  $mtlsInfo \in \langle \rangle s'.mtlsRequests \wedge mtlInfo.1 \equiv mtlNonce$  if possible; otherwise stop
25:         let  $s'.mtlsRequests := s'.mtlsRequests - \langle \langle mtlInfo \rangle \rangle$ 
26:         let  $mtlsKey := mtlInfo.2$ 
27:         let  $cnfValue := [x5t\#S256: hash(mTlsKey)]$ 
28:       else if  $authnScheme \equiv DPoP$  then  $\rightarrow$  DPoP sender constraining
29:         let  $dpopProof := m.headers[DPoP]$ 
30:         let  $dpopJwt := extractmsg(dpopProof)$ 
31:         let  $verificationKey := dpopJwt[headers][jwk]$ 
32:         if  $checksig(dpopProof, verificationKey) \neq \top \vee verificationKey \equiv \langle \rangle$  then
33:           stop  $\rightarrow$  Invalid DPoP signature (or empty jwk header)
34:         let  $dpopClaims := dpopJwt[payload]$ 
35:         let  $reqUri := \langle URL, S, m.host, m.path, \langle \rangle, \perp \rangle$ 
36:         if  $dpopClaims[htm] \neq m.method \vee dpopClaims[htu] \neq reqUri$  then
37:           stop  $\rightarrow$  DPoP claims do not match corresponding message
38:         if  $dpopClaims[nonce] \notin s'.dpopNonces$  then
39:           stop  $\rightarrow$  Invalid DPoP nonce
40:         if  $dpopClaims[ath] \neq hash(accessToken)$  then
41:           stop  $\rightarrow$  Invalid access token hash
42:         let  $s'.dpopNonces := s'.dpopNonces - \langle \langle dpopClaims[nonce] \rangle \rangle$ 
43:         let  $cnfValue := [jkt: hash(verificationKey)]$ 
44:       else
45:         stop  $\rightarrow$  Wrong Authorization header value
46:       let  $resource := \nu_4$   $\rightarrow$  Generate a fresh resource nonce
47:       let  $accessTokenContent$  such that  $accessTokenContent \equiv extractmsg(accessToken)$ 
48:          $\hookrightarrow$  if possible; otherwise let  $accessTokenContent := \diamond$ 
49:       if  $accessTokenContent \equiv \diamond$  then  $\rightarrow$  Not a structured AT, do Token Introspection
50:        $\rightarrow$  Store values for the pending request (needed when the RS gets the introspection response)
51:       let  $requestId := \nu_2$ 
52:       let  $s'.pendingResponses[requestId] := [expectedCNF: cnfValue, requestingClient: f,$ 
53:          $\hookrightarrow$   $resourceID: resourceID, originalRequest: m, originalRequestKey: k, resource: resource]$ 
54:       let  $url := asInfo[as\_introspect\_ep]$ 
55:       let  $rsCred \leftarrow asInfo[rs\_credentials]$   $\rightarrow$  Secret for authenticating at the AS (see also [69, Sec. 2.1])
56:       let  $headers := [Authorization: (Basic, rsCred)]$ 
57:       let  $body := [token: accessToken]$ 
58:       let  $message := \langle HTTPReq, \nu_3, POST, url.domain, url.path, url.parameters, headers, body \rangle$ 
59:       call HTTPS_SIMPLE_SEND( $[responseTo: TOKENINTROSPECTION, requestId: requestId], message, a, s'$ )
60:     else  $\rightarrow$  Check structured AT
61:       if  $cnfValue.1 \neq accessTokenContent[cnf].1 \vee cnfValue.2 \neq accessTokenContent[cnf].2$  then
62:         stop  $\rightarrow$  AT is bound to a different key
63:       if  $checksig(accessToken, asInfo[as\_key]) \neq \top$  then
64:         stop  $\rightarrow$  Verification of AT signature failed
```

$\rightarrow$  Algorithm continues on next page.

---

---

```

62:   let id := accessTokenContent[sub]
63:   if id  $\notin$   $s'.ids$  then
64:     stop → RS does not manage resources of this RO
        → Token binding successfully checked, the RS gives access to a resource of the identity
65:     let  $s'.resourceNonces[id][resourceID] := s'.resourceNonces[id][resourceID] + \langle \rangle resource$ 
66:     let  $m' := enc_s(\langle HTTPResp, m.nonce, 200, \langle \rangle, [resource:resource] \rangle, k)$ 
        → Leak resource request. Note that we only leak the application-layer message, and in particular, not the mTLS nonce.
67:     let leakingMessage :=  $\langle HTTPReq, \nu_{RRleak}, POST, m.domain, m.path, m.parameters, m.headers, [] \rangle$ 
68:     let leakAddress  $\leftarrow$  IPs
69:     stop  $\langle \langle f, a, m' \rangle, \langle leakAddress, a, \langle LEAK, leakingMessage \rangle \rangle \rangle, s'$ 
70:   else
71:     stop → Expected AT in Authorization header as mandated by FAPI 2.0

```

---

### Algorithm 13 Relation of a Resource Server $R^{rs}$ – Processing HTTPS Responses

---

```

1: function PROCESS_HTTPS_RESPONSE( $m, reference, request, key, a, f, s'$ )
2:   if  $reference[responseTo] \equiv \text{TOKENINTROSPECTION}$  then
3:     let pendingRequestInfo :=  $s'.pendingResponses[reference[requestId]]$ 
4:     let clientAddress :=  $pendingRequestInfo[requestingClient]$ 
5:     let expectedCNF :=  $pendingRequestInfo[expectedCNF]$ 
6:     let origReq :=  $pendingRequestInfo[originalRequest]$ 
7:     let originalRequestKey :=  $pendingRequestInfo[originalRequestKey]$ 
8:     let resourceID :=  $pendingRequestInfo[resourceID]$ 
9:     let resource :=  $pendingRequestInfo[resource]$ 
10:    if  $m.body[active] \neq \top$  then
11:      stop → Access token was invalid
12:    let responseCNF :=  $m.body[cnf]$ 
13:    if  $responseCNF.1 \neq expectedCNF.1 \vee responseCNF.2 \neq expectedCNF.2$  then
14:      stop → Access token was bound to a different key
15:    let id :=  $m.body[sub]$ 
16:    if id  $\notin$   $s'.ids$  then
17:      stop → RS does not manage resources of this RO
        → Token binding successfully checked, the RS gives access to a resource of the identity
18:      let  $s'.resourceNonces[id][resourceID] := s'.resourceNonces[id][resourceID] + \langle \rangle resource$ 
19:      let  $m' := enc_s(\langle HTTPResp, origReq.nonce, 200, \langle \rangle, [resource:resource] \rangle, originalRequestKey)$ 
        → Leak resource request. Note that we only leak the application-layer message, and in particular, not the mTLS nonce.
20:      let leakingMessage :=  $\langle HTTPReq, \nu_{RRleak}, POST, origReq.domain, origReq.path, origReq.parameters, origReq.headers, [] \rangle$ 
21:      let leakAddress  $\leftarrow$  IPs
22:      stop  $\langle \langle f, a, m' \rangle, \langle leakAddress, a, \langle LEAK, leakingMessage \rangle \rangle \rangle, s'$ 

```

---



APPENDIX E  
FORMAL SESSION INTEGRITY PROPERTIES

In [Section IV-E](#), we present our authorization and authentication properties. Here, we give the formal definition of the session integrity properties.

A. *Session Integrity*

On a high-level view, the two session integrity properties state that (1) an honest user, after logging in, is indeed logged in under their own account and not under the account of an attacker, and (2) similarly, that an honest user is accessing their own resources and not the resources of the attacker.

In addition to the authorization and authentication properties, it is important that the integrity of user sessions is not compromised. This is captured by two different session integrity properties.

More precisely, the first one, session integrity for authorization, ensures that an honest user should never use resources of the attacker. The second property, session integrity for authentication, captures that an honest user should never be logged in under the identity of the attacker.

We first define notations for the processing steps that represent important events during a run of a FAPI 2.0 web system.

*Definition 19 (User is logged in).* For a run  $\rho$  of a FAPI 2.0 web system with network attacker  $\mathcal{F}\mathcal{A}\mathcal{P}\mathcal{I}$  we say that a browser  $b$  was authenticated to a client  $c$  using an authorization server  $as$  and an identity  $id$  in a login session identified by a nonce  $lsid$  in processing step  $Q$  in  $\rho$  with

$$Q = (S, E, N) \xrightarrow{c \rightarrow E_{\text{out}}} (S', E', N')$$

and some event  $\langle y, y', m \rangle \in E_{\text{out}}$  such that  $m$  is an HTTPS response to an HTTPS request sent by  $b$  to  $c$  and we have that in the headers of  $m$  there is a header of the form  $\langle \text{Set-Cookie}, [\text{serviceSessionId}: \langle ssid, \top, \top, \top \rangle] \rangle$  for some nonce  $ssid$  such that  $S(c).\text{sessions}[lsid][\text{serviceSessionId}] \equiv ssid$  and  $S(c).\text{sessions}[lsid][\text{loggedInAs}] \equiv \langle d, id \rangle$  with  $d \in \text{dom}(as)$ . We then write  $\text{loggedIn}_\rho^Q(b, c, id, as, lsid)$ .

*Definition 20 (User started a login flow).* For a run  $\rho$  of a FAPI 2.0 web system with network attacker  $\mathcal{F}\mathcal{A}\mathcal{P}\mathcal{I}$  we say that the user of the browser  $b$  started a login session identified by a nonce  $lsid$  at the client  $c$  in a processing step  $Q$  in  $\rho$  if (1) in that processing step, the browser  $b$  was triggered, selected a document loaded from an origin of  $c$ , executed the script  $\text{script\_client\_index}$  in that document, and in that script, executed the Line 8 of Algorithm 8, and (2)  $c$  sends an HTTPS response corresponding to the HTTPS request sent by  $b$  in  $Q$  and in that response, there is a header of the form  $\langle \text{Set-Cookie}, [\text{__Host\_sessionId}: \langle lsid, \top, \top, \top \rangle] \rangle$ . We then write  $\text{started}_\rho^Q(b, c, lsid)$ .

*Definition 21 (User authenticated at an AS).* For a run  $\rho$  of a FAPI 2.0 web system with network attacker  $\mathcal{F}\mathcal{A}\mathcal{P}\mathcal{I}$  we say that the user of the browser  $b$  authenticated to an authorization server  $as$  using an identity  $id$  for a login session identified by a nonce  $lsid$  at the client  $c$  if there is a processing step  $Q = (S, E, N) \rightarrow (S', E', N')$  in  $\rho$  in which the browser  $b$  was triggered, selected a document loaded from an origin of  $as$ , executed the script  $\text{script\_as\_form}$  in that document, and in that script, (1) in Line 4 of Algorithm 11, selected the identity  $id$ , and (2) we have that

- the  $\text{scriptstate}$  of that document, when triggered in  $Q$ , contains a nonce  $\text{auth2Reference}$  such that  $\text{scriptstate}[\text{auth2\_reference}] \equiv \text{auth2Reference}$ , and
- there is a nonce  $\text{requestUri}$  such that  $S(as).\text{authorizationRequests}[\text{requestUri}][\text{auth2\_reference}] \equiv \text{auth2Reference}$ , and
- $S(c).\text{sessions}[lsid][\text{request\_uri}] \equiv \text{requestUri}$ .

We then write  $\text{authenticated}_\rho^Q(b, c, id, as, lsid)$ .

*Definition 22 (Resource Access).* For a run  $\rho$  of a FAPI 2.0 web system with network attacker  $\mathcal{F}\mathcal{A}\mathcal{P}\mathcal{I}$  we say that a browser  $b \in \mathcal{B}$  gets access to a resource of identity  $u$  stored at resource server  $rs$  managed by authorization server  $as$  through the session of client  $c$  identified by the nonce  $lsid$  in a processing step  $Q = (S, E, N) \rightarrow (S', E', N')$  in  $\rho$  if  $c$  executes Line 55 of Algorithm 2 in  $Q$ , includes the resource  $r$  in the body of the HTTPS response that is sent out there, and it holds true that

- 1)  $r \in \langle S'(rs).\text{resourceNonces}[u][\text{resourceId}] \rangle$  and  $as = \text{authorizationServerOfResource}^{rs}(\text{resourceID})$  (for some value  $\text{resourceId} \in \mathcal{T}_{\mathcal{N}}$ ),
- 2)  $\langle \langle \text{__Host\_sessionId}, \langle lsid, y, z, z' \rangle \rangle \in \langle S'(b).\text{cookies}[d] \rangle$  for  $d \in \text{dom}(c)$ ,  $y, z, z' \in \mathcal{T}_{\mathcal{N}}$ ,
- 3)  $S'(c).\text{sessions}[lsid][\text{resourceServer}] \in \text{dom}(rs)$ .
- 4) the request to which the client is responding contains a Cookie header with the cookie  $\langle \text{__Host\_sessionId} \rangle$  with the value  $lsid$

We then write  $\text{accessesResource}_\rho^Q(b, r, u, c, rs, as, lsid)$ .

*Definition 23 (Client Leaked Authorization Request).* Let  $\mathcal{F}API$  be a FAPI 2.0 web system with network attacker. For a run  $\rho$  of  $\mathcal{F}API$  with a processing step  $Q$ , a client  $c \in \mathcal{C}$ , a browser  $b$ , an authorization server  $as \in \mathcal{AS}$ , an identity  $id$ , a login session id  $lsid$ , and  $\text{loggedIn}_\rho^Q(b, c, id, as, lsid)$ , we say that  $c$  leaked the authorization request for  $lsid$ , if there is a processing step  $Q' = (S, E, N) \xrightarrow[c \rightarrow E_{\text{out}}]{} (S', E', N')$  in  $\rho$  prior to  $Q$  such that in  $Q'$ ,  $c$  executes Line 36 of Algorithm 2 and there is a nonce  $requestUri$  and an event  $\langle x, y, m \rangle \in E_{\text{out}}$  with  $m.1 \equiv \text{LEAK}$  and  $m.2.parameters[request\_uri] \equiv requestUri$  such that  $S'(c).sessions[lsid][request\_uri] \equiv requestUri$ .

#### Session Integrity Property for Authentication

This security property captures that (a) a user should only be logged in when the user actually expressed the wish to start a FAPI flow before, and (b) if a user expressed the wish to start a FAPI flow using some honest authorization server and a specific identity, then user is not logged in under a different identity.

*Definition 24 (Session Integrity for Authentication).* Let  $\mathcal{F}API$  be a FAPI 2.0 web system with network attacker. We say that  $\mathcal{F}API$  is secure w.r.t. session integrity for authentication iff for every run  $\rho$  of  $\mathcal{F}API$ , every processing step  $Q = (S, E, N) \rightarrow (S', E', N')$  in  $\rho$ , every browser  $b$  that is honest in  $S$ , every  $as \in \mathcal{AS}$ , every identity  $id$ , every client  $c \in \mathcal{C}$  that is honest in  $S$ , every nonce  $lsid$ , and  $\text{loggedIn}_\rho^Q(b, c, id, as, lsid)$  and  $c$  did not leak the authorization request for  $lsid$  (see Definition 23), we have that (1) there exists a processing step  $Q'$  in  $\rho$  (before  $Q$ ) such that  $\text{started}_\rho^{Q'}(b, c, lsid)$ , and (2) if  $as$  is honest in  $S$ , then there exists a processing step  $Q''$  in  $\rho$  (before  $Q$ ) such that  $\text{authenticated}_\rho^{Q''}(b, c, id, as, lsid)$ .

#### Session Integrity Property for Authorization

This security property captures that (a) a user should only access resources when the user actually expressed the wish to start a FAPI flow before, and (b) if a user expressed the wish to start a FAPI flow using some honest authorization server and a specific identity, then the user is not using resources of a different identity. We note that for this, we require that the resource server which the client uses is honest, as otherwise, the attacker can trivially return any resource.

*Definition 25 (Session Integrity for Authorization).* Let  $\mathcal{F}API$  be a FAPI 2.0 web system with network attacker. We say that  $\mathcal{F}API$  is secure w.r.t. session integrity for authorization iff for every run  $\rho$  of  $\mathcal{F}API$ , every processing step  $Q = (S, E, N) \rightarrow (S', E', N')$  in  $\rho$ , every browser  $b$  that is honest in  $S$ , every  $as \in \mathcal{AS}$ , every identity  $u$ , every client  $c \in \mathcal{C}$  that is honest in  $S$ , every  $rs \in \mathcal{RS}$  that is honest in  $S$ , every nonce  $r$ , every nonce  $lsid$ , we have that if  $\text{accessesResource}_\rho^Q(b, r, u, c, rs, as, lsid)$  and  $c$  did not leak the authorization request for  $lsid$  (see Definition 23), then (1) there exists a processing step  $Q'$  in  $\rho$  (before  $Q$ ) such that  $\text{started}_\rho^{Q'}(b, c, lsid)$ , and (2) if  $as$  is honest in  $S$ , then there exists a processing step  $Q''$  in  $\rho$  (before  $Q$ ) such that  $\text{authenticated}_\rho^{Q''}(b, c, u, as, lsid)$ .

By *session integrity* we denote the conjunction of both properties.

## APPENDIX F PROOFS

### A. Helper Lemmas

*Lemma 1 (Host of HTTP Request).* For any run  $\rho$  of a FAPI web system  $\mathcal{F}API$  with a network attacker, every configuration  $(S, E, N)$  in  $\rho$  and every process  $p \in \mathcal{C} \cup \mathcal{AS} \cup \mathcal{RS}$  that is honest in  $S$  it holds true that if the generic HTTPS server calls  $\text{PROCESS\_HTTPS\_REQUEST}(m_{dec}, k, a, f, s)$  in Algorithm 31, then  $m_{dec}.host \in \text{dom}(p)$ , for all values of  $k, a, f$  and  $s$ .

PROOF.  $\text{PROCESS\_HTTPS\_REQUEST}$  is called only in Line 9 of Algorithm 31. The input message  $m$  is an asymmetrically encrypted ciphertext. Intuitively, such a message is only decrypted if the process knows the private TLS key, where the private key used to decrypt is chosen (non-deterministically) according to the host of the decrypted message.

More formally, when  $\text{PROCESS\_HTTPS\_REQUEST}$  is called, the **stop** in Line 8 is not called. Therefore, it holds true that

$$\begin{aligned} & \exists inDomain, k' : \langle inDomain, k' \rangle \in S(p).tlskeys \wedge m_{dec}.host \equiv inDomain \\ & \Rightarrow \exists inDomain, k' : \langle inDomain, k' \rangle \in tlskeys^p \wedge m_{dec}.host \equiv inDomain \\ & \stackrel{\text{Def. (Appendix D-C)}}{\Rightarrow} \exists inDomain, k' : \langle inDomain, k' \rangle \in \{ \langle d, tlskey(d) \rangle \mid d \in \text{dom}(p) \} \wedge m_{dec}.host \equiv inDomain \end{aligned}$$

From this, it follows directly that  $m_{dec}.host \in \text{dom}(p)$ .

The first implication holds true due to  $S(p).tlskeys \equiv s_0^p.tlskeys \equiv tlskeys^p$ , as this sequence is never changed by any honest process  $p \in \mathcal{C} \cup \mathcal{AS} \cup \mathcal{RS}$  and due to the definitions of the initial states of clients, authorization servers, and resource servers (Definition 16, Definition 17, Definition 18). ■

*Lemma 2 (Client's Signing Key Does Not Leak).* For any run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of a FAPI web system  $\mathcal{F}\mathcal{A}\mathcal{P}\mathcal{I}$  with a network attacker, every configuration  $(S, E, N)$  in  $\rho$ , every client  $c \in \mathcal{C}$  that is honest in  $S$ , and every process  $p$  with  $p \neq c$ , all of the following hold true:

- $\text{signkey}(c) \notin d_\emptyset(S(p))$
- $\text{signkey}(c) \equiv s_0^c.\text{jwt}$
- $\text{signkey}(c) \equiv S(c).\text{jwt}$

PROOF.  $\text{signkey}(c) \equiv s_0^c.\text{jwt}$  immediately follows from [Definition 16](#).  $\text{signkey}(c) \equiv S(c).\text{jwt}$  follows from [Definition 16](#) and by induction over the processing steps: state subterm  $\text{jwt}$  of a client is never changed.

The only places in which an honest client accesses the  $\text{jwt}$  state subterm are: [Line 19](#) of [Algorithm 3](#), [Line 31](#) of [Algorithm 3](#), [Line 18](#) of [Algorithm 4](#), and [Line 24](#) of [Algorithm 6](#).

In [Line 19](#) of [Algorithm 3](#) and [Line 24](#) of [Algorithm 6](#), the  $\text{jwt}$  state subterm is only used in a  $\text{sig}(\cdot, \cdot)$  term constructor as signature key, i.e., cannot be extracted from the respective terms. Thus, it does not matter where these terms are stored or sent to. We conclude that these two usages of the  $\text{jwt}$  state subterm do not leak  $\text{signkey}(c)$  to any other process, in particular  $p$ .

In [Line 31](#) of [Algorithm 3](#) and [Line 18](#) of [Algorithm 4](#), the value of the  $\text{jwt}$  state subterm is stored in a variable  $\text{privKey}$ , which is then used in two places each:

- 1) In a  $\text{pub}(\cdot)$  term constructor ([Line 35](#) of [Algorithm 3](#) and [Line 24](#) of [Algorithm 4](#)). The  $\text{privKey}$  value cannot be extracted from these terms. Thus, it does not matter where these terms are stored or sent to.
- 2) In a  $\text{sig}(\cdot, \cdot)$  term constructor as signature key ([Line 37](#) of [Algorithm 3](#) and [Line 26](#) of [Algorithm 4](#)), i.e., cannot be extracted from the respective terms. Thus, it does not matter where these terms are stored or sent to.

By definition of  $\text{signkey}$  in [Appendix D-C](#) and the initial states of authorization servers ([Definition 17](#)), clients ([Definition 16](#)), browsers ([Definition 7](#)), and resource servers ([Definition 18](#)), we have that no other process initially knows  $\text{signkey}(c)$ .

We conclude that  $\text{signkey}(c) \notin d_\emptyset(S(p))$ . ■

*Lemma 3 (Client's TLS Keys Does Not Leak).* For any run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of a FAPI web system  $\mathcal{F}\mathcal{A}\mathcal{P}\mathcal{I}$  with a network attacker, every configuration  $(S, E, N)$  in  $\rho$ , every client  $c \in \mathcal{C}$  that is honest in  $S$ , every domain  $d_c \in \text{dom}(c)$ , and every process  $p$  with  $p \neq c$ , all of the following hold true:

- 1)  $\text{tlskey}(d_c) \notin d_\emptyset(S(p))$
- 2)  $\langle d_c, \text{tlskey}(d_c) \rangle \in^\diamond S^0(c).\text{tlskeys}$
- 3)  $\langle d_c, \text{tlskey}(d_c) \rangle \in^\diamond S(c).\text{tlskeys}$

PROOF. With [Definition 16](#),  $\langle d_c, \text{tlskey}(d_c) \rangle \in^\diamond S^0(c).\text{tlskeys}$  is equivalent to  $\langle d_c, \text{tlskey}(d_c) \rangle \in^\diamond \text{tlskeys}^c$ . This, in turn follows immediately from the definition of  $\text{tlskeys}^c$  in [Appendix D-C](#). Building on this, it is easy to check that the client never changes the contents of its  $\text{tlskeys}$  state subterm, i.e., we have  $\langle d_c, \text{tlskey}(d_c) \rangle \in^\diamond S(c).\text{tlskeys}$ .

The only places in which an honest client accesses any value in its  $\text{tlskeys}$  state subterm are:

**Line 3 of Algorithm 2** Here, the value is only used to decrypt a message (i.e., cannot leak).

**Line 13 of Algorithm 7** Here, the client only uses the value to create a public key. As the equational theory does not allow extraction of private keys from public keys, it does not matter where that public key is stored or sent to.

**Line 7 of Algorithm 31** Here, the value is once again only used to decrypt a message.

By definition of  $\text{tlskey}$ ,  $\text{tlskeys}^p$  in [Appendix D-C](#) and the initial states of authorization servers ([Definition 17](#)), clients ([Definition 16](#)), browsers ([Definition 7](#)), and resource servers ([Definition 18](#)), we have that no other process initially knows  $\text{tlskey}(d_c)$ .

We conclude that  $\text{tlskey}(d_c) \notin d_\emptyset(S(p))$ . ■

*Lemma 4 (Generic Server - Correctness of Reference and Request).* For any run  $\rho$  of a FAPI web system  $\mathcal{F}\mathcal{A}\mathcal{P}\mathcal{I}$  with a network attacker, every processing step  $P = (S^P, E^P, N^P) \rightarrow (S^{P'}, E^{P'}, N^{P'})$  in  $\rho$ , every  $p \in \mathcal{C} \cup \text{AS} \cup \text{RS}$  being honest in  $S$ , it holds true that if  $p$  calls `PROCESS_HTTPS_RESPONSE` in  $P$  with *reference* being the second and *request* being the third input argument, then there exists a previous processing step in which  $p$  calls `HTTPS_SIMPLE_SEND` with *reference* being the first and *request* being the second input argument.

PROOF. Let  $p \in \mathcal{C} \cup \text{AS} \cup \text{RS}$  be honest in  $S^P$ .  $p$  calls the `PROCESS_HTTPS_RESPONSE` function only in the generic HTTPS server algorithm in [Line 26](#) of [Algorithm 31](#). The values *reference* and *request* are taken from  $S^P(p).\text{pendingRequests}$  in [Line 19](#) of [Algorithm 31](#). Thus,  $p$  added these values to  $\text{pendingRequests}$  in a previous processing step  $O = (S^O, E^O, N^O) \rightarrow (S^{O'}, E^{O'}, N^{O'})$  by executing [Line 15](#) of [Algorithm 31](#), as this is the only location where a client, authorization server, or resource server adds entries to  $\text{pendingRequests}$  and as  $\text{pendingRequests}$  is initially empty (see [Definitions 16, 17, and 18](#)). In  $O$ , the process  $p$  takes both values from  $S^O(p).\text{pendingDNS}$  in [Line 13](#) and [Line 14](#) of [Algorithm 31](#). Initially  $\text{pendingDNS}$  is empty (as  $p$  is a client, an authorization server, or a resource server), and  $p$  adds values to  $\text{pendingDNS}$  only in [Line 2](#) of [Algorithm 26](#),

where the reference and request values are the input arguments of `HTTPS_SIMPLE_SEND`. Thus, in some processing step prior to  $O$ ,  $p$  called `HTTPS_SIMPLE_SEND` with *reference* being the first and *request* being the second input argument. ■

*Lemma 5 (Code used in Token Request was received at Redirection Endpoint).* For any run  $\rho$  of a FAPI web system  $\mathcal{F}API$  with a network attacker, every processing step

$$P = (S, E, N) \xrightarrow[c \rightarrow E_{out}^P]{e_{in}^P \rightarrow c} (S', E', N')$$

in  $\rho$  with  $c \in \mathcal{C}$  being honest in  $S$ , it holds true that if Algorithm 2 (`PROCESS_HTTPS_RESPONSE`) is called in  $P$  with *reference* being the second and *request* being the third input argument, and if  $reference[responseTo] \equiv \text{TOKEN}$ , then there is a previous configuration  $(S^{L'}, E^{L'}, N^{L'})$  such that  $request.body[code] \equiv S^{L'}(c).sessions[reference[session]][redirectEpRequest][message].parameters[code]$ .

PROOF. As shown in Lemma 4, there exists a processing step  $L = (S^L, E^L, N^L) \rightarrow (S^{L'}, E^{L'}, N^{L'})$  prior to  $P$  in which  $c$  called `HTTPS_SIMPLE_SEND` with the same reference and request values. The only line in which a client calls `HTTPS_SIMPLE_SEND` with  $reference[responseTo] \equiv \text{TOKEN}$  is Line 41 of Algorithm 3 (`SEND_TOKEN_REQUEST`). The code included in the request is the input parameter of `SEND_TOKEN_REQUEST` (see Lines 6, 39, and 40 of Algorithm 3).

`SEND_TOKEN_REQUEST` is called only in Line 21 of Algorithm 1, i.e., at the redirection endpoint (`/redirect_ep`) of the client, and the code is taken from the parameters of the redirection request. The redirection request is stored into  $S^{L'}(c).sessions[sessionId][redirectionEpRequest][message]$  in Line 20 of Algorithm 1, with  $sessionId \equiv reference[session]$ . ■

*Lemma 6 (Authorization Server's Signing Key Does Not Leak).* For any run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of a FAPI web system  $\mathcal{F}API$  with a network attacker, every configuration  $Q = (S, E, N)$  in  $\rho$ , every authorization server  $as \in \text{AS}$  that is honest in  $S$ , every term  $t$  with  $checksig(t, \text{pub}(\text{signkey}(as))) \equiv \top$ , and every process  $p$  with  $p \neq as$ , all of the following hold true:

- $\text{signkey}(as) \notin d_0(S(p))$
- $\text{signkey}(as) \equiv s_0^{as}.jwk$
- $\text{signkey}(as) \equiv S(as).jwk$
- if  $t$  is known (Definition 77) to  $p$  in  $Q$ , then  $t$  was created (Definition 75) by  $as$  in a processing step  $s_e$  prior to  $Q$  in  $\rho$

PROOF.  $\text{signkey}(as) \equiv s_0^{as}.jwk$  immediately follows from Definition 17.  $\text{signkey}(as) \equiv S(as).jwk$  follows from Definition 17 and by induction over the processing steps: state subterm `jwk` of a client is never changed.

By Definitions 16, 17, 18, 60, and Appendix D-C, we have that no process (except  $as$ ) initially knows  $\text{signkey}(as)$ , i.e.,  $\text{signkey}(as) \notin d_0(S^0(p))$ .

The only places in which an honest authorization server accesses the `jwk` state subterm are: Line 11 of Algorithm 9, Line 127 of Algorithm 9, and Line 139 of Algorithm 9.

In Line 127 of Algorithm 9 and Line 139 of Algorithm 9, the `jwk` state subterm is only used in a  $\text{sig}(\cdot, \cdot)$  term constructor as signature key, i.e., cannot be extracted from the respective terms. Thus, it does not matter where these terms are stored or sent to. We conclude that these two usages of the `jwk` state subterm do not leak  $\text{signkey}(as)$  to any other process, in particular  $p$ .

In Line 11 of Algorithm 9, the `jwk` state subterm is only used in a  $\text{pub}(\cdot)$  term constructor, i.e., cannot be extracted from the constructed term. Thus, it does not matter where such a term is stored or sent to.

We conclude that  $\text{signkey}(as) \notin d_0(S(p))$ .

To complete the proof, we now have to show that any term  $t$  with  $checksig(t, \text{pub}(\text{signkey}(as))) \equiv \top$  known to  $p$  in  $Q$  was created by  $as$  in a processing step  $s_e$  prior to  $Q$  in  $\rho$ :

By Definitions 16, 17, 18, 7, and Appendix D-C, we have that no process (except  $as$ ) initially knows such a term  $t$ , i.e.,  $t \notin d_0(S^0(p))$ . Together with Definition 51 and Definition 75, this implies that  $t$  can only be known to  $p$  in some configuration  $Q'$  if  $t$  was contained in some event  $e$  “received” by  $p$  at an earlier point in  $\rho$  (i.e.,  $e$  was the input event in a processing step in  $\rho$  with  $p$ ). Since such an  $e$  is not part of  $E^0$  (Definition 70),  $e$  must have been emitted by some process in a processing step  $s_e$  prior to  $Q'$  in  $\rho$ . Definition 51 and Definition 72 imply that  $p$  (or any other process  $\neq as$ ) cannot have emitted  $e$  in  $s_e$  (i.e., cannot have created  $t$  in  $s_e$ ).

Therefore,  $as$  must have emitted  $e$  and hence created  $t$  in  $s_e$ , i.e., prior to  $Q$  in  $\rho$ . ■

*Lemma 7 (mTLS Nonce created by AS does not Leak).* For every run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of a FAPI web system  $\mathcal{F}API$  with a network attacker, every configuration  $(S, E, N)$  in  $\rho$ , every authorization server  $as \in \text{AS}$  that is honest in  $S^n$ , every client  $c \in \mathcal{C}$  that is honest in  $S^n$  with client id  $clientId$  at  $as$ , every  $i \in \mathbb{N}$  with  $0 \leq i \leq |S(as).mtlsRequests[clientId]|$ , and every process  $p$  with  $as \neq p \neq c$  it holds true that  $mtlsNonce := S(as).mtlsRequests[clientId].i.1$  does not leak to  $p$ , i.e.,  $mtlsNonce \notin d_0(S^n(p))$ .

PROOF. Initially, the `mtlsRequests` subterm of the authorization server's state is empty, i.e.,  $S^0(as).mtlsRequests \equiv \langle \rangle$  (Definition 17). An authorization server only adds values to the `mtlsRequests` subterm in Line 161 of Algorithm 9, where the mTLS nonce is chosen as a fresh nonce (Line 157 of Algorithm 9).

Let  $(S^i, E^i, N^i) \rightarrow (S^{i'}, E^{i'}, N^{i'})$  be the processing step in which the nonce is chosen (note that  $(S^i, E^i, N^i)$  is prior to  $(S, E, N)$  in  $\rho$ ). In the same processing step, the authorization server sends out the nonce in Line 163 of Algorithm 9, asymmetrically encrypted with the public key

$$\begin{aligned}
& \text{clientKey} \\
& \equiv S^i(as).clients[\text{clientId}][\text{mtls\_key}] && \text{(Line 158, Algorithm 9)} \\
& \equiv s_0^{as}.clients[\text{clientId}][\text{mtls\_key}] && \text{(value is never changed)} \\
& \equiv \text{clientInfoAS}(as)[\text{clientId}][\text{mtls\_key}] && \text{(Definition 17)} \\
& \equiv \langle \{ \langle \text{cli.client\_id}, \text{as\_cli}(\text{cli}) \rangle \mid d_c \in \text{Doms}, d_{as} \in \text{dom}(as) \} \rangle [\text{clientId}][\text{mtls\_key}] && \text{(Definition 14)}
\end{aligned}$$

with  $\text{cli} = \text{clientInfo}(d_{as}, d_c)$  and  $\text{as\_cli}$  as in Definition 14. As  $\text{clientId}$  is the identifier of  $c$  at  $as$  and as no two clients have the same identifier at an authorization server (see Definition 12), it follows that  $d_c \in \text{dom}(c)$ .

As there is a dictionary entry in  $S^i(as).clients[\text{clientId}]$  with the key `mtls_key` (Line 159, Algorithm 9), it follows that

$$\begin{aligned}
& \text{clientKey} \\
& \equiv \text{pub}(\text{clientInfo}(d_{as}, d_c).tlsSKey) && (d_{as} \in \text{dom}(as), d_c \in \text{dom}(c); \text{ see also Definition 17}) \\
& \equiv \text{pub}(\text{tlskey}(d_c)) && (d_c \in \text{dom}(c))
\end{aligned}$$

The corresponding private key is  $\text{tlskey}(d_c) \in \text{tlskeys}^c$ , which is only known to  $c$  (Lemma 3). (The `mtlsNonce` saved in `mtlsRequests` is not sent in any other place).

This implies that the encrypted nonce can only be decrypted by  $c$ .  $c$  decrypts messages only in Line 3 of Algorithm 2. (The only other place where a message is decrypted asymmetrically by  $c$  is in the generic HTTPS server (Line 7 of Algorithm 31), where the process would stop due to the requirement that the decrypted message must begin with `HTTPReq`).

We also note that the encrypted message created by the authorization server containing the nonce also contains a public TLS key of  $as$ . (This holds true due to Lemma 1).

After decrypting the mTLS nonce and public TLS key of  $as$  in Line 3 of Algorithm 2, the client stores the sequence  $\langle \text{request.host}, \text{clientId}, \text{pubKey}, \text{mtlsNonce} \rangle$  into the `mtlsCache` subterm of its state in Line 8 of Algorithm 2, where  $\text{clientId}, \text{pubKey} \in \mathcal{T}_{\mathcal{N}}$  and, in particular,

- $\text{request.host}$  is a domain of  $as$  (see Line 5, Algorithm 2)
- $\text{mtlsNonce}$  is the mTLS nonce chosen by  $as$ .

Thus, the nonce is stored at the client together with a domain of the authorization server. After storing the values, the client stops in Line 9 of Algorithm 2 without creating an event and without storing the nonce in any other place.

$c$  sends mTLS nonces only to domains of  $as$ . The client accesses values stored in the `mtlsCache` subterm of its state only in the following places:

#### Case 1: Algorithm 3

In this algorithm, the client accesses the `mtlsCache` subterm only in Line 14 and Line 27.

In both cases, the sequence containing the nonce is removed from the `mtlsCache` subterm (Lines 16 and 28), and the mTLS nonce is sent by calling the `HTTPS_SIMPLE_SEND` function. The HTTP request that is passed to `HTTPS_SIMPLE_SEND` in Line 41 contains the retrieved mTLS nonces only in the body, under the dictionary key `TLS_AuthN` (Line 15, Line 39) or `TLS_binding` (Line 25, Line 29, Line 39).

In all cases, the domain stored in the sequence that is retrieved from the `mtlsCache` subterm of the client state (i.e., the first entry of the sequence) is the host of the HTTPS request that the client constructs (see Lines 14, 27).

#### Case 2: Algorithm 4

Here, the client accesses the `mtlsCache` state subterm only in Line 13. As in the first case, the sequence from which the mTLS nonce is chosen is removed from the `mtlsCache` subterm (Line 16 of Algorithm 4). The nonce is sent in the body of an HTTP request, using the dictionary key `TLS_binding` (see Line 14) by calling `HTTPS_SIMPLE_SEND` in Line 30. The request is sent to the same domain that is stored in the sequence containing the mTLS nonce.

#### Case 3: Algorithm 6

Here, Line 19 is the last line in which the client accesses the `mtlsCache` state subterm. As in the previous cases, the client removes the corresponding sequence from the `mtlsCache` subterm (Line 21). The client creates an HTTPS request which contains the mTLS nonce in the body under the key `TLS_AuthN` (Lines 20, 33, and 35). Again, the request is sent to the same domain that is stored in the sequence containing the mTLS nonce (see Line 35).

In all cases, the HTTP request is sent to the domain stored in the first entry of the sequence containing the mTLS nonce (stored in the `mtlsCache` subterm). Let  $req_{c \rightarrow as}$  be the HTTP request that the client sends by calling `HTTPS_SIMPLE_SEND`.

`HTTPS_SIMPLE_SEND` stores the request  $req_{c \rightarrow as}$  (which contains the mTLS nonce) in the `pendingDNS` state subterm of  $c$ , see Line 2 of Algorithm 26 and then stops with the DNS request (which does not contain the nonce) in Line 3 of Algorithm 26. Thus, after finishing this processing step, the client stores the mTLS nonce only in its `pendingDNS` state subterm.

The client accesses the `pendingDNS` state subterm only within the else case in Line 10 of Algorithm 31, i.e., when it receives the DNS response. There, it either stops without a new event and without changing its state in Line 12 of Algorithm 31, or creates a new `pendingRequests` entry containing the request  $req_{c \rightarrow as}$  (and thus, also the mTLS nonce) in Line 15 of Algorithm 31. In this case, the client removes the request from the `pendingDNS` state subterm in Line 17 of Algorithm 31, i.e., regarding the client state, the mTLS nonce is only contained in the newly created `pendingRequests` entry. The client finishes the processing step by encrypting  $req_{c \rightarrow as}$  with the key of the domain that was stored along with the mTLS nonce, i.e., a key of  $as$ , see Lines 16 and 18 of Algorithm 31.

**$as$  does not leak mTLS nonce contained in request.**

As the client encrypts  $req_{c \rightarrow as}$  asymmetrically with a key of  $as$ , it follows that only  $as$  can decrypt the HTTPS request (Lemma 27).

The authorization server only decrypts terms in the generic HTTPS server algorithms. More specifically, this request is decrypted (only) in Line 7 of Algorithm 31, as this is the only place where an authorization server decrypts a message asymmetrically, and then used as a function argument for `PROCESS_HTTPS_REQUEST` which is modeled in Algorithm 9.

In Algorithm 9, none of the endpoints except for the `PAR` (Line 53) and token endpoint (Line 81) reads, stores, or sends out a value stored in the body of the request under the `TLS_AuthN` or `TLS_binding` key.

The `PAR` and token endpoints pass the HTTP request to the `AUTHENTICATE_CLIENT` helper function (Algorithm 10), which removes an entry from the `mtlsRequests` state subterm and returns this entry; the `/par` endpoint code does not use this value. The token endpoint uses this value for token binding (Lines 110–120), but the nonce is not added to any state subterm and not sent out in a network message. Thus, the endpoints of the authorization server do not store the mTLS nonces contained in requests in any state subterm and do not send them out in any network message.

**$c$  does not leak mTLS nonce in request after getting the response.** When the client receives the HTTPS response to  $req_{c \rightarrow as}$ , the generic HTTPS server removes the message from the `pendingRequests` state subterm and calls `PROCESS_HTTPS_RESPONSE` with the request as the third function argument. Algorithm 2 does not store a nonce contained in the body of the request and does not create new network messages containing such a nonce.

Summing up, the client sends the mTLS nonce created by the authorization server only back to that authorization server. As an honest authorization server never sends out such a nonce received in a request, we conclude that the nonce never leaks to any other process, in particular not to  $p$ . ■

*Lemma 8 (mTLS Nonce created by RS does not Leak).* For every run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of a FAPI web system  $\mathcal{F}_{API}$  with a network attacker, every configuration  $(S, E, N)$  in  $\rho$ , every resource server  $rs \in RS$  that is honest in  $S^n$ , every client  $c \in C$  that is honest in  $S^n$ , every domain  $d_c \in \text{dom}(c)$ , every  $i \in \mathbb{N}$  with  $0 \leq i \leq |S(rs).\text{mtlsRequests}|$  and with  $S(rs).\text{mtlsRequests}.i.2 \equiv \text{pub}(\text{tlskey}(d_c))$ , and every process  $p$  with  $rs \neq p \neq c$  it holds true that  $\text{mtlsNonce} := S(rs).\text{mtlsRequests}.i.1$  does not leak to  $p$ , i.e.,  $\text{mtlsNonce} \notin d_0(S^n(p))$ .

**PROOF.** This proof is similar to the proof of Lemma 7:

Initially, the `mtlsRequests` subterm of the resource server's state is empty, i.e.,  $S^0(rs).\text{mtlsRequests} \equiv \langle \rangle$  (Definition 18). A resource server only adds values to the `mtlsRequests` subterm in Line 5 of Algorithm 12, where the mTLS nonce (the first value of the sequence that is added to `mtlsRequests`) is a fresh nonce (Line 3 of Algorithm 12).

Let  $(S^i, E^i, N^i) \rightarrow (S^{i'}, E^{i'}, N^{i'})$  be the processing step in which the nonce is chosen (note that  $(S^i, E^i, N^i)$  is prior to  $(S, E, N)$  in  $\rho$ ). In the same processing step, the resource server sends out the nonce in Line 7 of Algorithm 12, asymmetrically encrypted with the public key  $\text{pub}(\text{tlskey}(d_c))$  (precondition of the lemma, see also Line 5 and Line 6 of Algorithm 12).

The corresponding private key is  $\text{tlskey}(d_c) \in \text{tlskeys}^c$ , which is only known to  $c$  (Lemma 3). (The  $\text{mtlsNonce}$  saved in `mtlsRequests` is not sent in any other place).

This implies that the encrypted nonce can only be decrypted by  $c$ .  $c$  decrypts messages only in Line 3 of Algorithm 2. (The only other place where a message is decrypted asymmetrically by  $c$  is in the generic HTTPS server (Line 7 of Algorithm 31), where the process would stop due to the requirement that the decrypted message must begin with `HTTPReq`).

We also note that the encrypted message created by the resource server containing the nonce also contains a public TLS key of  $rs$ . (This holds true due to Lemma 1).

After decrypting the mTLS nonce and public TLS key of  $rs$  in Line 3 of Algorithm 2, the client stores the sequence  $\langle \text{request.host}, \text{clientId}, \text{pubKey}, \text{mtlsNonce} \rangle$  into the `mtlsCache` subterm of its state, where  $\text{clientId}, \text{pubKey} \in \mathcal{T}_{\mathcal{N}}$  and, in particular,

- $\text{request.host}$  is a domain of  $rs$  (see Line 5, Algorithm 2)
- $\text{mtlsNonce}$  is the mTLS nonce chosen by  $rs$ .

Thus, the nonce is stored at the client together with a domain of the resource server. After storing the values, the client stops in Line 9 of Algorithm 2 without creating an event and without storing the nonce in any other place.

*c* sends mTLS nonces only to domains of *rs*. The client accesses values stored in the `mtlsCache` subterm of its state only in the following places:

### Case 1: Algorithm 3

In this algorithm, the client accesses the `mtlsCache` subterm only in Line 14 and Line 27.

In both cases, the sequence containing the nonce is removed from the `mtlsCache` subterm (Lines 16 and 28), and the mTLS nonce is sent by calling the `HTTPS_SIMPLE_SEND` function. The HTTP request that is passed to `HTTPS_SIMPLE_SEND` in Line 41 contains the retrieved mTLS nonces only in the body, under the dictionary key `TLS_AuthN` (Line 15, Line 39) or `TLS_binding` (Line 25, Line 29, Line 39).

In all cases, the domain stored in the sequence that is retrieved from the `mtlsCache` subterm of the client state (i.e., the first entry of the sequence) is the host of the HTTPS request that the client constructs (see Lines 14, 27).

### Case 2: Algorithm 4

Here, the client accesses the `mtlsCache` state subterm only in Line 13. As in the first case, the sequence from which the mTLS nonce is chosen is removed from the `mtlsCache` subterm (Line 16 of Algorithm 4). The nonce is sent in the body of an HTTP request, using the dictionary key `TLS_binding` (see Line 14) by calling `HTTPS_SIMPLE_SEND` in Line 30. The request is sent to the same domain that is stored in the sequence containing the mTLS nonce.

### Case 3: Algorithm 6

Here, Line 19 is the last line in which the client accesses the `mtlsCache` state subterm. As in the previous cases, the client removes the corresponding sequence from the `mtlsCache` subterm (Line 21). The client creates an HTTPS request which contains the mTLS nonce in the body under the key `TLS_AuthN` (Lines 20, 33, and 35). Again, the request is sent to the same domain that is stored in the sequence containing the mTLS nonce (see Line 35).

In all cases, the HTTP request is sent to the domain stored in the first entry of the sequence containing the mTLS nonce (stored in the `mtlsCache` subterm). Let  $req_{c \rightarrow rs}$  be the request that the client sends by calling `HTTPS_SIMPLE_SEND`.

`HTTPS_SIMPLE_SEND` stores the request  $req_{c \rightarrow rs}$  (which contains the mTLS nonce) in the `pendingDNS` state subterm of *c*, see Line 2 of Algorithm 26, and then stops with the DNS request (which does not contain the nonce) in Line 3 of Algorithm 26. Thus, after finishing this processing step, the client stores the mTLS nonce only in its `pendingDNS` state subterm.

The client accesses the `pendingDNS` state subterm only within the else case in Line 10 of Algorithm 31, i.e., when it receives the DNS response. There, it either stops without a new event and without changing its state in Line 12 of Algorithm 31, or creates a new `pendingRequests` entry containing the request  $req_{c \rightarrow rs}$  (and thus, also the mTLS nonce) in Line 15 of Algorithm 31. In this case, the client removes the request from the `pendingDNS` state subterm in Line 17 of Algorithm 31, i.e., regarding the client state, the mTLS nonce is only contained in the newly created `pendingRequests` entry. The client finishes the processing step by encrypting  $req_{c \rightarrow rs}$  with the key of the domain that was stored along with the mTLS nonce, i.e., a key of *rs*, see Lines 16 and 18 of Algorithm 31.

*rs* does not leak mTLS nonce contained in request. As the HTTP request  $req_{c \rightarrow rs}$  is encrypted asymmetrically with a key of *rs*, it follows that only the resource server can decrypt the request. The resource server only decrypts terms in the generic HTTPS server algorithms. More specifically, this request is decrypted (only) in Line 7 of Algorithm 31, as this is the only place where an resource server decrypts a message asymmetrically, and then used as a function argument for `PROCESS_HTTPS_REQUEST` which is modeled in Algorithm 12.

In Algorithm 12, the `/MTLS-prepare` and `/DPoP-nonce` endpoints (Line 2 and Line 8 of Algorithm 12) do not read, store, or send out a value stored in the body of the request under the `TLS_AuthN` or `TLS_binding` key.

The last endpoint starting at Line 13 of Algorithm 12 accesses values stored in the body of the request under the `TLS_binding` key in Line 23. This value is not added to any state subterm and not sent out in a network message if Line 57 of Algorithm 12 is executed. If Line 48 of Algorithm 12 is true, then the whole request (including the `TLS_binding` value in the request body) is stored in the `pendingResponses` subterm of the resource server's state. However, the resource server never stores the body of requests stored in `pendingResponses` into any other subterm of its state and does not send out any value contained in the body.

*c* does not leak mTLS nonce in request after getting the response. When receiving the HTTPS response to  $req_{c \rightarrow rs}$ , the generic HTTPS server removes the message from the `pendingRequests` state subterm and calls `PROCESS_HTTPS_RESPONSE` with the request as the third function argument. Algorithm 2 does not store a nonce contained in the body of the request and does not create new network messages containing such a nonce.

Summing up, the client sends the mTLS nonce created by the resource server only back to that resource server. As an honest resource server never sends out such a nonce received in a request, we conclude that the nonce never leaks to any other process, in particular not to *p*. ■

*Lemma 9 (JWS client assertion created by client does not leak).* For any run  $\rho$  of a FAPI web system  $\mathcal{F}API$  with a network attacker, every configuration  $(S, E, N)$  in  $\rho$ , every authorization server  $as \in AS$  that is honest in  $S$ , every client  $c \in C$  that is honest in  $S$  with client identifier  $clientId$  at  $as$ , every term  $t$  with

- $checksig(t, pub(signkey(c))) \equiv \top$
- $extractmsg(t)[iss] \equiv clientId$
- $extractmsg(t)[sub] \equiv clientId$
- $extractmsg(t)[aud].host \in dom(as)$  or  $extractmsg(t)[aud] \in dom(as)$

and every process  $p$  with  $as \neq p \neq c$  it holds true that  $t \notin d_\emptyset(S(p))$ .

PROOF. The private signing key  $signkey(c)$  is part of the clients initial state ( $s_c^0.jwk$ , see Definition 16). Initially, no other process and no waiting event contains the key except as a public key  $pub(signkey(c))$ , which cannot be used to retrieve the private key.

The client  $c$  never leaks the private key: The client uses the private key only for creating signatures or creating public keys (from which the private key cannot be extracted) see Algorithm 3 (Lines 19, 35, and 37), Algorithm 4 (Lines 24 and 26), and Algorithm 6 (Line 24).

Thus, only  $c$  can create a term  $t$ , i.e., the attacker cannot create  $t$  itself by signing a dictionary with the corresponding  $iss$ ,  $sub$ , and  $aud$  values. In the following, we show that such a term created by  $c$  does not leak to the attacker.

The client signs dictionaries with  $aud$  dictionary key only in two locations:

**Case 1: Line 19 of Algorithm 3** The signature created in Line 19 of Algorithm 3 is added to the body of an HTTP request (Lines 20, 39, and 40). The client sends that HTTP request (the token request) to the token endpoint it has cached for the AS identified by the issuer identifier in  $extractmsg(t)[aud]$  (i.e.,  $selectedAS$  in the context of Algorithm 3). From Lemma 16, we know that this token endpoint is a URL of the selected AS, i.e., the token request is sent to and encrypted for the party to which the domain  $selectedAS$  belongs (see the call of `HTTPS_SIMPLE_SEND` in Line 41 (using `responseTo: TOKEN` in the first function argument). This party is an honest AS by the preconditions of this lemma.

**Case 2: Line 24 of Algorithm 6** As in the previous case, the signature created in Line 24 of Algorithm 6 is added to the body of an HTTP request (Lines 25, 33, and 35). Similar to the first case, this request (the PAR request) is encrypted for and sent to the PAR endpoint  $c$  has cached for the party to which  $extractmsg(t)[aud]$  belongs. Analogous to the first case, we can apply Lemma 16 to conclude that this party must be an honest AS (and the request is stored by  $c$  with `responseTo: PAR`).

When the client receives a corresponding HTTPS response, the generic HTTPS server decrypts the message and calls `PROCESS_HTTPS_RESPONSE`. The original request (containing the client assertion) is used as the third function argument. The instantiation of `PROCESS_HTTPS_RESPONSE` (Algorithm 2) does not access the body of the request when processing `TOKEN` or `PAR` responses and hence cannot leak it in any way.

When processing the HTTPS request in Algorithm 9, the authorization server does not store the client assertion and does not create a network message containing the client assertion.

Overall, we conclude that no other process can derive a client assertion created by an honest client for an honest authorization server. ■

*Lemma 10 (Client Authentication).* For any run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of a FAPI web system  $\mathcal{F}API$  with a network attacker, every authorization server  $as \in AS$ , every client  $c \in C$ , every processing step  $Q$  in  $\rho$

$$(S, E, N) \xrightarrow[as \rightarrow E_{out}]{e_{in} \rightarrow as} (S', E', N')$$

with  $c$  and  $as$  being honest in  $S'$  and every client identifier  $clientId$  of  $c$  at  $as$  it holds true that:

If  $e_{in} \equiv \langle x, y, enc_a(\langle m, k \rangle, k') \rangle$  (for some  $x, y, k, k'$ ) with  $m$  being an HTTP request such that all of the following hold true, then  $c$  created  $m$  (Definition 75):

- $client\_id \in m.body \Rightarrow m.body[client\_id] \equiv clientId$  and
- $client\_assertion \in m.body \Rightarrow extractmsg(m.body[client\_assertion])[iss] \equiv clientId$  and
- $m.path \equiv /par \vee m.path \equiv /token$  and
- $E_{out}$  is not empty.

PROOF. We first note that authorization servers do not emit any HTTP(S) requests.

$as$  executes either Line 80 or Line 142 of Algorithm 9. We first show that in processing step  $Q$ ,  $as$  executes either Line 80 or Line 142 of Algorithm 9: When processing  $e_{in}$ , the generic HTTPS server calls `PROCESS_HTTPS_REQUEST` in Line 9 of Algorithm 31, as the checks done in Line 2, 10, 19 and 27 complete successfully and the authorization server's instantiation of `PROCESS_OTHER` (Algorithm 30) does not create an output event, i.e., if  $as$  ended up in Algorithm 30,  $E_{out}$  would be empty. Furthermore, the `stop` in Line 8 of Algorithm 31 is not executed as this `stop` does not emit an output event.

Thus, Algorithm 31 calls `PROCESS_HTTPS_REQUEST` in processing step  $Q$ . The authorization server's instantiation of `PROCESS_HTTPS_REQUEST` is defined in Algorithm 9.



If  $m.path \equiv /par$  (with  $m$  as in the statement of the lemma), then the PAR endpoint starting in Line 53 of Algorithm 9 is executed. No **stop** except for the last (unconditional) stop in Line 80 emits an event.

Analogously, if  $m.path \equiv /token$  (with  $m$  as in the statement of the lemma), then the (unconditional) **stop** in Line 142 was executed in the processing step, as no other **stop** emits events.

**HTTP request contains values that only  $c$  and  $as$  know.** Non-empty  $E_{out}$  implies that the checks done in the AUTHENTICATE\_CLIENT helper function (Algorithm 10, called in Line 56 or 84 of Algorithm 9) did not lead to a **stop**.

In both cases, Algorithm 10 is called with the HTTP request  $m$  and  $S(as)$  as input arguments. As Line 26 of Algorithm 10 is not executed (because  $E_{out}$  is not empty), it follows that  $client\_assertion \in m.body$  or  $TLS\_AuthN \in m.body$ .

**Case 1:  $client\_assertion \in m.body$ :** As  $extractmsg(m.body[client\_assertion])[iss] \equiv clientId$  (lemma precondition), it holds true that the verification key used for verifying the signature in Line 4 of Algorithm 10 is

$$\begin{aligned} & S(as).clients[clientId][jwk\_key] && \text{(Line 4 and 12, Algorithm 10)} \\ \equiv & s_0^{as}.clients[clientId][jwk\_key] && \text{(value is never changed)} \\ \equiv & clientInfoAS(as)[clientId][jwk\_key] && \text{(Definition 17)} \\ \equiv & \langle \{ \langle cli.client\_id, as\_cli(cli) \rangle \mid d_c \in \text{Doms}, d_{as} \in \text{dom}(as) \} \rangle [clientId][jwk\_key] && \text{(Definition 14)} \end{aligned}$$

with  $cli = clientInfo(d_{as}, d_c)$  and  $as\_cli$  as in Definition 14. As  $clientId$  is the identifier of  $c$  at  $as$  and as no two clients have the same identifier at an authorization server (see Definition 12), it follows that  $d_c \in \text{dom}(c)$ .

As there is a dictionary entry in  $S(as).clients[clientId]$  with the key  $jwk\_key$  (otherwise,  $as$  would **stop** with empty  $E_{out}$  in Line 5 of Algorithm 10), it follows that the verification key is

$$\begin{aligned} & \text{pub}(clientInfo(d_{as}, d_c).jwtSkey) && (d_{as} \in \text{dom}(as), d_c \in \text{dom}(c); \text{ see also Definition 14}) \\ \equiv & \text{pub}(\text{signkey}(c)) && \text{(Definition 12)} \end{aligned}$$

The corresponding private key  $\text{signkey}(c)$  is only known to  $c$  (Lemma 2).

In the following, let  $cli\_assertion = extractmsg(m.body[client\_assertion])$ . As  $cli\_assertion[iss] \equiv clientId$  (precondition of the lemma),  $cli\_assertion[sub] \equiv clientId$  (Line 12 of Algorithm 10), and  $cli\_assertion[aud].host \in \text{dom}(as)$  or  $cli\_assertion[aud] \in \text{dom}(as)$  (Line 14 of Algorithm 10 and as the host of the request is a domain of the authorization server as shown in Lemma 1), we can apply Lemma 9.

Thus, for all processes  $p$ , it holds true that  $m.body[client\_assertion] \notin d_\emptyset(S'(p))$  if  $as \neq p \neq c$ , i.e., only  $c$  and  $as$  can derive  $m.body[client\_assertion]$ . As authorization servers do not create HTTP requests, it follows that  $m$  was *created* by  $c$ .

**Case 2:  $TLS\_AuthN \in m.body$ :** From Lines 17–19 of Algorithm 10 it follows that

$$\exists i \in \mathbb{N}. S(as).mtlsRequests[m.body[client\_id]].i.1 \equiv m.body[TLS\_AuthN]$$

Note that  $client\_id \in m.body$  as otherwise, the **stop** in Line 23 of Algorithm 10 will be executed.

Now, we can apply Lemma 7 with  $\rho'$  ( $\rho'$  being the trace prefix of  $\rho$  up to and including  $(S', E', N')$ ).

Thus, for all processes  $p$ , it holds true that  $m.body[TLS\_AuthN] \notin d_\emptyset(S'(p))$  if  $as \neq p \neq c$ , i.e., only  $c$  and  $as$  can derive  $m.body[TLS\_AuthN]$ . As authorization servers do not create HTTP requests, we conclude that  $m$  was *created* by  $c$ . ■

**Lemma 11 (mTLS Keys of Clients stored at Authorization Servers).** For every authorization server  $as \in \text{AS}$  and every term  $clientId \in \mathcal{T}_{\mathcal{C}}$  it holds true that if  $s_0^{as}.clients[clientId][mtls\_key]$  is not  $\text{pub}(\diamond)$  and not  $\langle \rangle$ , then  $\exists c \in \mathbb{C}, d \in \text{dom}(c)$  such that  $s_0^{as}.clients[clientId][mtls\_key] \equiv \text{pub}(\text{tlskey}(d))$ .

PROOF.

$$\begin{aligned} & s_0^{as}.clients[clientId][mtls\_key] && (1) \\ \equiv & clientInfoAS(as)[clientId][mtls\_key] && (2) \\ \equiv & \langle \{ \langle cli.client\_id, as\_cli(cli) \rangle \mid \exists d_c \in \text{Doms}, d_{as} \in \text{dom}(as): clientInfo(d_{as}, d_c) \equiv cli \} \rangle [clientId][mtls\_key] && (3) \\ \equiv & as\_cli(clientInfo(d_{as}, d_c))[mtls\_key] \text{ (with } d_c \in \text{Doms}, d_{as} \in \text{dom}(as), clientInfo(d_{as}, d_c)[client\_id] \equiv clientId) && (4) \\ \equiv & \text{pub}(clientInfo(d_{as}, d_c).mtlsSKey) && (5) \end{aligned}$$

Notes on the above equivalences:

- (2) As per Definition 17.
- (3) See Definition 14.

(5) The value is not  $\langle \rangle$  (precondition); with  $d_c \in \text{Doms}$ ,  $d_{as} \in \text{dom}(as)$ ,  $\text{clientInfo}(d_{as}, d_c)[\text{client\_id}] \equiv \text{clientId}$ ; see [Definition 14](#) for the definition of  $\text{as\_cli}$ .

As this value is not  $\text{pub}(\diamond)$  (precondition of the lemma), it follows from [Definition 12](#) that the value is equal to  $\text{pub}(\text{tlskey}(d_c))$  and  $d_c$  is a domain of a client  $c \in \mathcal{C}$  (as  $\text{clientInfos}$  is defined for domains of clients, see [Definition 12](#)).

*Lemma 12 (DPoP proof secrecy (RS)).* For any run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of a FAPI web system  $\mathcal{F}\mathcal{A}\mathcal{P}\mathcal{I}$  with a network attacker, every configuration  $(S, E, N)$  in  $\rho$ , every resource server  $rs \in \text{RS}$  that is honest in  $S$ , every client  $c \in \mathcal{C}$  that is honest in  $S$ , every term  $t$  with

- $\text{checksig}(t, \text{pub}(\text{signkey}(c))) \equiv \top$
- $\text{extractmsg}(t)[\text{payload}][\text{htu}].\text{host} \in \text{dom}(rs)$ ,
- $\text{ath} \in \langle \rangle \text{extractmsg}(t)[\text{payload}]$ ,
- $\text{extractmsg}(t)[\text{payload}][\text{nonce}] \in S(rs).\text{dpopNonces}$

and every process  $p$  with  $rs \neq p \neq c$  it holds true that  $t \notin d_\emptyset(S(p))$ .

**PROOF.** From [Lemma 2](#), we have that only  $c$  can create a term  $t$ , i.e., the attacker cannot create  $t$  itself by signing a dictionary with the corresponding  $\text{payload}$  value. In the following, we show that such a term created by  $c$  does not leak to the attacker.

The client signs dictionaries with a  $\text{payload}$  dictionary key only in two locations: In [Line 37](#) of [Algorithm 3](#), where the  $\text{payload}$  dictionary does not contain an  $\text{ath}$  value (see [Line 36](#) of [Algorithm 3](#)), and in [Line 26](#) of [Algorithm 4](#).

In [Line 26](#) of [Algorithm 4](#), the client sends the term  $t$  to  $\text{extractmsg}(t)[\text{payload}][\text{htu}].\text{host}$  via `HTTPS_SIMPLE_SEND` (using `responseTo: RESOURCE_USAGE` in the first function argument), see [Lines 21, 25, 29, and 30](#) of [Algorithm 4](#). The client does not store  $t$  in any other subterm except for those needed by `HTTPS_SIMPLE_SEND`. The term  $t$  is added (only) to the headers of the HTTP request using the DPoP dictionary key, see [Line 28](#) of [Algorithm 4](#). The client also adds an `Authorization` header containing a dictionary with a DPoP dictionary key, and in particular, no other value, see [Lines 27 and 29](#) of [Algorithm 4](#).

We note that the generic part of the client model (which takes care of DNS resolution and sending the actual HTTPS request after the `HTTPS_SIMPLE_SEND` call) does not send out or use  $t$  in any way – except for the sending of the actual request, which is encrypted for the domain  $\text{extractmsg}(t)[\text{payload}][\text{htu}].\text{host}$ , i.e., for  $rs$ , which can only be decrypted by  $rs$  ([Lemma 27](#)).

When the client receives the HTTPS response to this request, the generic HTTPS server decrypts the message and calls `PROCESS_HTTPS_RESPONSE`. The original request (containing the signed term) is used as the third function argument. The instantiation of `PROCESS_HTTPS_RESPONSE` ([Algorithm 2](#)) does not access the headers of the request when processing `RESOURCE_USAGE` responses.

When processing the HTTPS request created by the client in [Algorithm 12](#), the resource server does not access the DPoP header (in particular, it does not add the term to its state and does not create a network message containing the value) in the `/MTLS-prepare` and `/DPoP-nonce` endpoints ([Lines 2 and 8](#) of [Algorithm 12](#)). For all other path values ([Line 13](#) of [Algorithm 12](#)), the resource server first checks whether the resource identified by the path is managed by a supported authorization server. If this is not the case, then the resource server stops without changing the state and without emitting events ([Line 17](#) of [Algorithm 12](#)). Otherwise, the resource server will eventually invalidate the nonce value stored in the DPoP proof in [Line 42](#) of [Algorithm 12](#) (by removing it from the `dpopNonces` subterm of the resource server's state), as the request contains an `Authorization` header containing a dictionary with the DPoP keyword (see [Lines 19 and 28](#) of [Algorithm 12](#)). The **stops** before the removal of the nonce from the state of the resource server do not modify the state of the resource server and do not lead to new events.

We note that the `dpopNonces` state subterm of the resource server does not contain any value twice, as the resource server only adds fresh nonces to the state subterm, see the endpoint in [Line 8](#) of [Algorithm 12](#). Thus, the nonce is not contained in `dpopNonces` after [Line 42](#) of [Algorithm 12](#) is executed, and the resource server it does not add it back to the `dpopNonces` state subterm afterwards.

Thus, if the resource server does not finish with a **stop** without any arguments, it holds true that  $\text{extractmsg}(t)[\text{payload}][\text{nonce}]$  is not contained in the `dpopNonces` subterm of the new resource server's state, as it always stops with the updated state. (If it finishes with a **stop** without any arguments, then  $t$  will not leak, as there is no change in any state and no new event).

Overall, we conclude that no other process can derive a signed term  $t$  (as in the statement of the lemma) created by an honest client for an honest resource server. ■

*Lemma 13 (Access Token can only be used by Honest Client).* For

- every run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of  $\mathcal{F}\mathcal{A}\mathcal{P}\mathcal{I}$  with a network attacker,
- every resource server  $rs \in \text{RS}$  that is honest in  $S^n$ ,
- every identity  $id \in \langle \rangle s_0^{rs}.\text{ids}$ ,
- every processing step in  $\rho$

$$Q = (S^Q, E^Q, N^Q) \xrightarrow[\text{rs} \rightarrow E_{\text{out}}^Q]{e_{\text{in}}^Q \rightarrow \text{rs}} (S^{Q'}, E^{Q'}, N^{Q'})$$

- every  $resourceID \in \mathbb{S}$  with  $as = \text{authorizationServerOfResource}^{rs}(resourceID)$  being honest in  $S^Q$ ,

it holds true that:

If  $\exists r, x, y, k, m_{\text{resp}}. \langle x, y, \text{enc}_s(m_{\text{resp}}, k) \rangle \in \langle \rangle E_{\text{out}}^Q$  such that  $m_{\text{resp}}$  is an HTTP response,  $r := m_{\text{resp}}.\text{body}[\text{resource}]$ , and  $r \in \langle \rangle S^{Q'}(rs).\text{resourceNonce}[id][resourceID]$ , then

- 1) There exists a processing step

$$P = (S^P, E^P, N^P) \xrightarrow[\text{rs} \rightarrow E_{\text{out}}^P]{e_{\text{in}}^P \rightarrow \text{rs}} (S^{P'}, E^{P'}, N^{P'})$$

such that

- a) either  $P = Q$  or  $P$  prior to  $Q$  in  $\rho$ , and
  - b)  $e_{\text{in}}^P$  is an event  $\langle x, y, \text{enc}_a(\langle m_{\text{req}}, k_1 \rangle, k_2) \rangle$  for some  $x, y, k_1$ , and  $k_2$  where  $m_{\text{req}} \in \mathcal{T}_{\mathcal{N}}$  is an HTTP request which contains a term (access token)  $t$  in its Authorization header, i.e.,  $t \equiv m_{\text{req}}.\text{headers}[\text{Authorization}].2$ , and
  - c)  $r$  is a fresh nonce generated in  $P$  at the resource endpoint of  $rs$  in Line 46 of Algorithm 12.
- 2)  $t$  is bound to a key  $k \in \mathcal{T}_{\mathcal{N}}$ ,  $as$ , and  $id$  in  $S^Q$  (see Definition 1).
  - 3) If there exists a client  $c \in \mathbb{C}$  such that  $k \equiv \text{pub}(\text{signkey}(c))$  or  $k \equiv \text{pub}(\text{tlskey}(d_c))$  and  $d_c \in \text{dom}(c)$ , and if  $c$  is honest in  $S^n$ , then  $e_{\text{in}}^P$  was created by  $c$ .

PROOF. An honest resource server sends HTTPS responses with a resource dictionary key only in Line 69 of Algorithm 12 and Line 22 of Algorithm 13.

### Case 1: Line 69 of Algorithm 12

**First Postcondition** In the same processing step, i.e.,  $P = Q$ , the resource server received an HTTPS request with an access token and generated the resource:

$e_{\text{in}}^Q$  is an event containing an HTTPS request, as Algorithm 12 is only called by the generic HTTPS server in Line 9 of Algorithm 31. As the check done in Line 7 of Algorithm 31 was true and the stop in Line 8 was not executed, it follows that the input event of Algorithm 31 was an event containing an HTTPS request  $m_{\text{req}}$  (as in the first statement of the post-condition of the lemma).

$m_{\text{req}}$  contains an Authorization header (Line 19 of Algorithm 12).

The resource that is sent out in Line 69 of Algorithm 12 is a freshly chosen nonce generated in the same processing step in Line 46 of Algorithm 12 (see also Line 66 of Algorithm 12). This concludes the proof of the first post-condition.

**Second Postcondition** As Line 69 of Algorithm 12 is in the second case of Line 48 of Algorithm 12, it follows that  $\text{extractmsg}(m.\text{headers}[\text{Authorization}].2)$  is a structured access token (see Lines 21 and 47).

The access token is signed by  $\text{authorizationServerOfResource}^{rs}(resourceID)$ : The value of  $responsibleAS$  (in Line 15) is equal to

$$\begin{aligned} & S^Q(rs).\text{resourceASMapping}[resourceID] && \text{(Line 15 of Algorithm 12)} \\ \equiv & s_0^{rs}.\text{resourceASMapping}[resourceID] && \text{(value is never changed)} \\ \in & \text{dom}(\text{authorizationServerOfResource}^{rs}(resourceID)) && \text{(Definition 18)} \end{aligned}$$

As required by the precondition of the lemma,  $as = \text{authorizationServerOfResource}^{rs}(resourceID)$  is honest in  $S^Q$ . The signature of the access token is checked in Line 60 of Algorithm 12 using the verification key

$$\begin{aligned} & asInfo[as\_key] \\ \equiv & S^Q(rs).\text{asInfo}[responsibleAS][as\_key] && (responsibleAS \in \text{dom}(as), \text{Line 18}) \\ \equiv & s_0^{rs}.\text{asInfo}[responsibleAS][as\_key] && \text{(value is never changed)} \\ \equiv & \text{signkey}(\text{dom}^{-1}(responsibleAS)) && \text{(Definition 18)} \\ \equiv & \text{signkey}(as) \end{aligned}$$

The authorization server  $as$  only uses this key in the following locations:

- Line 13 of Algorithm 9: Endpoint returning public key
- Line 127 of Algorithm 9: Signing access token
- Line 139 of Algorithm 9: Signing ID token

As ID tokens created by an authorization server do not contain a  $\text{cnf}$  claim (see Lines 134-139 of Algorithm 9), it follows that  $\text{extractmsg}(m.\text{headers}[\text{Authorization}].2)$  is an access token created by  $as$  in Line 127 of Algorithm 9.

Let  $O = (S^O, E^O, N^O) \xrightarrow[as \rightarrow E_{out}^O]{e_{in}^O \rightarrow as} (S^{O'}, E^{O'}, N^{O'})$  be the processing step in which the authorization server created and signed the access token. After finishing the processing step,  $as$  stores the access token in  $S^{O'}(as).records.i[access\_token]$ , for some natural number  $i$  (as Line 130 of Algorithm 9 was executed by the authorization server). Note: we know that  $i$  is a natural number and not a “longer” pointer due to the last condition in Line 92 of Algorithm 9.

The structured access token contains a value  $extractmsg(m.headers[Authorization].2)[sub] \in {}^\diamond S^Q(rs).ids$  (Line 47, 62, and 63 of Algorithm 12). This identity is used as a dictionary key for storing the resource (see Line 65 of Algorithm 12). The ids stored at the resource server are never changed, i.e.,  $S^Q(rs).ids \equiv s_0^{rs}.ids$ . When creating the access token, the authorization server takes this value from  $S^O(as).records.i[sub]$  with the same  $i$  as above (Line 92 and 126 of Algorithm 9). As the remaining lines of the token endpoint do not change this value, it follows that  $S^O(as).records.i[sub] \equiv S^{O'}(as).records.i[sub]$ .

From the successful check of Line 58 of Algorithm 12 (as we assume that the resource server returns a resource in Line 69), it follows that either

- $accessTokenContent[cnf].1 \equiv x5t\#S256$  or
- $accessTokenContent[cnf].1 \equiv jkt$ ,

as  $cnfValue$  is set in Line 27 or Line 43 of Algorithm 12.

The authorization server sets the  $cnf$  value of access tokens only in Line 126 of Algorithm 9. The value is determined either in Line 109 or Line 120 of Algorithm 9, and the authorization server stores the  $cnf$  value into the same record as the  $access\_token$  and  $sub$  values, see Line 131 of Algorithm 9, i.e.,  $S^{O'}(as).records.i[cnf]$  is either  $[jkt: hash(k)]$  or  $[x5t\#S256: hash(k)]$ , for some value  $k$ .

As authorization servers do not remove sequences from their  $records$  state subterm, it follows that the access token is bound to some term  $k \in \mathcal{T}_{\mathcal{N}}$ , the authorization server  $as$ , and  $id$  in  $S^Q$ , by which we conclude the proof of the second postcondition for this case.

### Third Postcondition

Let  $c \in C$  be honest in  $S^n$ .

**Case 1.3.1:**  $k \equiv pub(signkey(c))$  As already shown, the authorization server determines the  $cnf$  value either in Line 109 or Line 120 of Algorithm 9.

- Line 109 of Algorithm 9: The structured access token contains the value  $accessTokenContent[cnf].1 \equiv jkt$ . Thus, the resource server executed Line 43 of Algorithm 12 (in the processing step  $P$ ), i.e., the request  $e_{in}^P$  contains a DPoP proof  $dpopProof$  (in the header of the request, see Line 29 of Algorithm 12). All preconditions of Lemma 12 are true:
  - $checksig(dpopProof, pub(signkey(c))) \equiv \top$  (see Line 32 of Algorithm 12)
  - $extractmsg(dpopProof)[payload][htu].host \in dom(rs)$  (see Line 36 of Algorithm 12 and Lemma 1)
  - $ath \in {}^\diamond extractmsg(dpopProof)[payload]$ , (see Line 40 of Algorithm 12)
  - $extractmsg(dpopProof)[payload][nonce] \in S(rs).dpopNonces$  (see Line 38 of Algorithm 12)

Thus, we can apply Lemma 12 and conclude that  $dpopProof$  can only be known by  $c$  and  $rs$ . The only places where a resource server sends a request are Lines 55 and 67 of Algorithm 12. In the first case, the request in question is a token introspection request whose Authorization header uses the Basic scheme. Processing of such a request by the resource server would lead to an empty  $E_{out}^Q$  in Line 45 of Algorithm 12. In the latter case, the resource server leaks the resource request – but only after invalidating the mTLS nonce (Lines 24f. of Algorithm 12) or DPoP nonce (Line 42 of Algorithm 12), i.e., processing this request again would lead to an empty  $E_{out}^Q$  in Line 24 of Algorithm 12, or Line 38 of Algorithm 12. Hence, resource servers do not send requests with valid DPoP or mTLS nonces to themselves and it follows that only  $c$  could have created the request  $e_{in}^P$ .

- Line 120 of Algorithm 9: The structured access token contains the value  $accessTokenContent[cnf].1 \equiv x5t\#S256$ .

The value  $accessTokenContent[cnf].2$  was chosen by the authorization server in Line 120 of Algorithm 9 (as this is the only location where the authorization server sets the  $cnf$  value with the  $x5t\#S256$  dictionary key), where it was set to  $hash(mTlsKey)$ . The value  $mTlsKey$  is set to  $mtlsInfo.2$  in Line 119 of Algorithm 9. The sequence  $mtlsInfo$  is chosen in Line 84 or Line 117 of Algorithm 9. In both cases,  $mTlsKey$  is set to  $s_0^{as}.clients[clientId][mtls\_key]$ :

- Line 84 of Algorithm 9:  $mtlsInfo$  is the third entry of the return value of AUTHENTICATE\_CLIENT. As shown above, AUTHENTICATE\_CLIENT (Algorithm 10) determines the client identifier  $clientId$  from the HTTP request and also determines the type of the client (see Lines 7, 8, 20, 21). As shown previously, the type of the client is either  $pkjwt\_mTLS$  or  $mTLS\_mTLS$ . Thus, the body of the request does not contain a

value `client_assertion`, as otherwise, the **stop** in Line 10 of Algorithm 10 would have prevented the authorization server to issue the access token. In particular, the **return** in Line 28 was executed and the third return value was taken from  $S^O(as).mTLSRequests[clientId]$  (Line 19 of Algorithm 10). Initially, the `mTLSRequests` subterm of the authorization server's state is empty (see Definition 17). The authorization server adds values to `mTLSRequests` only in Line 161 of Algorithm 9. The second sequence entry is  $s_0^{as}.clients[clientId][mTLS\_key]$  (due to Line 158 of Algorithm 9 and because the `clients` subterm of the authorization server's state is never modified).

- Line 117 of Algorithm 9:  $mTLSInfo$  is taken from  $S^O(as).mTLSRequests[clientId]$ . As shown in the previous case, the second sequence entry of  $mTLSInfo$  is equal to  $s_0^{as}.clients[clientId][mTLS\_key]$ .

When adding values to `mTLSRequests` in Line 161 of Algorithm 9, the authorization server ensures that the value of  $s_0^{as}.clients[clientId][mTLS\_key]$  is not  $\langle \rangle$  and not  $\text{pub}(\diamond)$  (Line 159 of Algorithm 9). Thus, we can apply Lemma 11 and conclude that there exists a client  $c \in \mathbb{C}$  and  $s_0^{as}.clients[clientId][mTLS\_key] \equiv \text{pub}(\text{tlskey}(d_c))$  with  $d_c \in \text{dom}(c)$ .

Thus, it is not possible that this value is a public signature verification key.

**Case 1.3.2:**  $k \equiv \text{pub}(\text{tlskey}(d_c))$  and  $d_c \in \text{dom}(c)$  As in the previous case, the authorization server determined the `cnf` value either in Line 109 or Line 120 of Algorithm 9.

- Line 109 of Algorithm 9: As only  $c$  knows the private key  $\text{tlskey}(d_c)$  (Lemma 3), it follows that  $c$  created a term  $dpopProof$  such that  $\text{checksig}(dpopProof, \text{pub}(\text{tlskey}(d_c))) \equiv \top$  (Line 103 of Algorithm 9). However, as  $c$  is an honest client, it never creates such a term, as all signatures are created using  $\text{signkey}(c)$  (see Line 19 of Algorithm 3, Line 37 of Algorithm 3, Line 26 of Algorithm 4, Line 37 of Algorithm 3, Line 26 of Algorithm 4, Line 24 of Algorithm 6, and Definition 16).
- Line 120 of Algorithm 9: The structured access token contains the value  $\text{accessTokenContent}[cnf].1 \equiv x5t\#S256$ . Thus, the resource server executed Line 27 of Algorithm 12 (in the processing step  $P$ ). This means that  $e_{in}^P$  contains a value  $mTLSNonce$  in the body of the request such that  $\langle mTLSNonce, \text{pub}(\text{tlskey}(d_c)) \rangle \in \langle S^P(rs).mTLSRequests$ .

If the client  $c$  is honest in  $S^m$ , then it is also honest in  $S^P$ , and we can apply Lemma 8 and conclude that only  $c$  and  $rs$  can derive  $m_{req}.body[\text{TLS\_binding}]$ . As resource servers do not send requests containing `TLS_binding` in the request body, it follows that the HTTP request  $m_{req}$  was created by  $c$ .

## Case 2: Line 22 of Algorithm 13

**First Postcondition** In Line 22 of Algorithm 13, the resource server is processing an HTTP response  $resp_{introsp}$  (with the reference `TOKENINTROSPECTION`, see Line 2 of Algorithm 13). An honest resource server sends HTTP requests only by calling `HTTPS_SIMPLE_SEND` in Line 56 of Algorithm 12 (again with the reference `TOKENINTROSPECTION`). Let  $req_{introsp}$  be the corresponding request to  $resp_{introsp}$ . The processing step in which the resource server emitted  $req_{introsp}$  is  $P$  (as in the postcondition of the lemma): The input event of  $P$  contains an HTTP request  $m_{req}$  (again as in the first postcondition) with an access token  $t \equiv m_{req}.headers[\text{Authorization}].2$  (Line 19 of Algorithm 12). The resource  $r$  that the resource server sends out in Line 22 of Algorithm 13 (in the processing step  $Q$ ) was stored by the resource server in  $S^{P'}.pendingResponses$  in Line 50 of Algorithm 12, and the resource was generated in Line 46 of Algorithm 12 (in the processing step  $P$ ).

**Second Postcondition** The request  $req_{introsp}$  was sent by  $rs$  to a domain of  $as$ :  $responsibleAS$  in Line 15 of Algorithm 12 is a domain of  $as$ , as shown in the proof of the first case. Thus, it follows that  $S^P(rs).asInfo[responsibleAS][as\_introspect\_ep]$  is  $\langle \text{URL}, S, dom_{as}, /introspect, \langle \rangle, \perp \rangle$ , with  $dom_{as} \in \text{dom}(as)$  (see Definition 18).

Furthermore,  $req_{introsp}$  contains the value  $m_{req}.headers[\text{Authorization}].2$ , see Line 21 and Line 54 of Algorithm 12.

The authorization server  $as$  processes this request in the introspection endpoint in Line 143 of Algorithm 9. As the resource server did not stop in Line 11 of Algorithm 13, we conclude that the access token sent by the resource server in  $P$  is active, i.e., the authorization server executed Line 152 of Algorithm 9. Thus, there is a value  $record$  in the `records` state subterm of the authorization server's state with the access token (Line 148 of Algorithm 9), and in this record, there is a `cnf` and a `subject` entry (Line 152 of Algorithm 9) The `cnf` and `subject` values are added to the body of the introspection response, and the resource server checks that the `subject` value is contained in the list of `ids` that the resource server stores in  $S^Q(rs).ids$  (Line 16 of Algorithm 13).

An honest authorization server adds `cnf` values to an entry of its `records` state entry only in the token endpoint in Line 131 of Algorithm 9. Thus, this value is either  $[jkt: \text{hash}(k)]$  (see Line 109 of Algorithm 9), or  $[x5t\#S256: \text{hash}(k)]$  (see Line 120 of Algorithm 9), for some value  $k$ .

**Third Postcondition** The resource server checks in Line 13 of Algorithm 13 that the `cnf` value that the authorization server put into the response  $resp_{introsp}$  is equal to the  $cnfValue$  that the resource server stored in Line 50 of Algorithm 12 in the

processing step  $P$ . The resource server does the same checks in  $P$  as in the first case (i.e., when sending out the response in Line 69 of Algorithm 12). Thus, it holds true that the request processed in  $P$  either contains a DPoP proof that only  $c$  and  $rs$  can derive, or an mTLS nonce that only  $c$  and  $rs$  can derive. The proof is analogous to the proof of the first case, i.e., only  $c$  could have created the request  $e_{in}^P$ . ■

**Lemma 14 (Redirect URI Properties).** For any run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of a FAPI web system  $\mathcal{F}API$  with a network attacker, every configuration  $(S, E, N)$  in  $\rho$ , every authorization server  $as \in AS$  that is honest in  $S$ , every client  $c \in C$  that is honest in  $S$  with client identifier  $clientId \neq \langle \rangle$  at  $as$ , and every  $requestUri$ , all redirect URIs for  $c$  stored at  $as$  are HTTPS URIs and belong to  $c$ . Or, more formally: Let  $rec = S(as).authorizationRequests[requestUri]$ , then  $rec[client\_id] \equiv clientId$  implies both  $rec[redirect\_uri].protocol \equiv S$ , and  $rec[redirect\_uri].host \in \text{dom}(c)$

PROOF. The only places in which an honest authorization server writes to its `authorizationRequests` state subterm are:

- Line 28 of Algorithm 9: Here, the authorization server does not change or create values under the `client_id` or `redirect_uri` keys.
- Line 78 of Algorithm 9: See below.

In the latter case, the authorization server is processing a pushed authorization request, i.e., an HTTPS request  $req$  to the `/par` endpoint.

In order to get to Line 78 of Algorithm 9,  $req$  must contain valid client authentication data (see Lines 56 and 60), in particular,  $req.body$  must contain a client id (under key `client_id`) and either a value under key `TLS_AuthN` or `client_assertion`. In the latter case, Line 4 of Algorithm 10 together with Line 12 of Algorithm 10 and Line 60 of Algorithm 9 ensure that  $\text{extractmsg}(req.body[client\_assertion])[iss] \equiv req.body[client\_id]$ . We note that reaching Line 78 of Algorithm 9 implies that the current processing step will output an event (there are no **stops** between Line 78 and Line 80 of Algorithm 9). Hence, we can apply Lemma 10.

When reaching Line 78 of Algorithm 9,  $req$  also must contain a `redirectUri` value in  $req.body[redirect\_uri]$  (Line 63 of Algorithm 9). Furthermore, this `redirectUri` must be an HTTPS URI (Line 65 of Algorithm 9) and this is the value stored in the authorization server's `authorizationRequests` state subterm (in a record under the key `redirect_uri`), together with  $req.body[client\_id]$  (under key `client_id`).

Line 54 of Algorithm 9 ensures that  $req.body$  contains a field `code_challenge_method` with value S256.

From Lemma 10, we know that  $c$  must have created  $req$ . Since  $c$  is honest and the only place in which an honest client produces an HTTPS request with a `code_challenge_method` with value S256 is in Line 35 of Algorithm 6, we can conclude that the value of  $req.body[redirect\_uri]$  is the one selected in Lines 2f. of Algorithm 6. This implies  $req.body[redirect\_uri].host \in \text{dom}(c)$  and hence  $rec[redirect\_uri].host \in \text{dom}(c)$ . ■

**Lemma 15 (Integrity of Client's Session Storage).** For any run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of a FAPI web system  $\mathcal{F}API$  with a network attacker, every configuration  $(S, E, N)$  in  $\rho$ , every client  $c \in C$  that is honest in  $S$ , and every login session id  $lsid$ , we have that if  $lsid \in S(c).sessions$ , then all of the following hold true:

- 1)  $selected\_AS \in S(c).sessions[lsid]$
- 2) for all configurations  $(S', E', N')$  after  $(S, E, N)$  in  $\rho$  we have  $S'(c).sessions[lsid][selected\_AS] \equiv S(c).sessions[lsid][selected\_AS]$

PROOF. Since we have  $S^0(c).sessions$  (Definition 16), we know that if  $lsid \in S(c).sessions$ , such an entry must have been stored there by  $c$ . Clients only ever store/add such an entry in Line 10 of Algorithm 1, where a key `selected_AS` is part of the stored entry. The key used to refer to the entry inside `sessions` is a fresh nonce (i.e.,  $lsid$  is a fresh nonce there). Hence, whenever a client first stores an entry in `sessions` under key  $lsid$ , this entry contains a key `selected_AS`.

It is easy to see that Line 10 of Algorithm 1 is indeed the only place in which a client stores any value under key `selected_AS` in the `sessions` state subterm. Similarly, it is easy to check that this line is also the only place in which a client (over)writes a whole entry in the `sessions` state subterm. Hence, we can conclude: A key `selected_AS` is present whenever a client adds an entry to the `sessions` state subterm and neither the value stored under that key, nor the `sessions` entry itself are overwritten or removed anywhere, implying 1). In addition, if the client ever executes Line 10 of Algorithm 1 again, it will never overwrite an existing entry, because it will use a fresh login session id, thus we have 2). ■

**Lemma 16 (Integrity of Client's oauthConfigCache).** For any run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of a FAPI web system  $\mathcal{F}API$  with a network attacker, every configuration  $(S, E, N)$  in  $\rho$ , every authorization server  $as \in AS$  that is honest in  $S$ , every client  $c \in C$  that is honest in  $S$ , and every domain  $d \in \text{dom}(as)$ , it holds true that if  $d \in S(c).oauthConfigCache$ , we have all of the following:

- 1)  $S(c).oauthConfigCache[d][issuer] \equiv d$
- 2)  $S(c).oauthConfigCache[d][auth_ep] \equiv \langle \text{URL}, S, d, /auth, \langle \rangle, \perp \rangle$

- 3)  $S(c).oauthConfigCache[d][token\_ep] \equiv \langle URL, S, d, /token, \langle \rangle, \perp \rangle$
- 4)  $S(c).oauthConfigCache[d][par\_ep] \equiv \langle URL, S, d, /par, \langle \rangle, \perp \rangle$
- 5)  $S(c).oauthConfigCache[d][introspec\_ep] \equiv \langle URL, S, d, /introspection, \langle \rangle, \perp \rangle$
- 6)  $S(c).oauthConfigCache[d][jwks\_uri] \equiv \langle URL, S, d, /jwks, \langle \rangle, \perp \rangle$

We note that this implies that all these entries in  $S(c).oauthConfigCache[d]$  are never changed once they have been stored and that all entries are created in the same processing step.

PROOF. We start by noting that  $S^0(c).oauthConfigCache \equiv \langle \rangle$  (Definition 16), i.e., the `oauthConfigCache` state subterm is initially empty. An honest client only ever writes to its `oauthConfigCache` state subterm in Line 16 of Algorithm 2 when processing an HTTPS response. Hence,  $d \in S(c).oauthConfigCache$  implies that there must have been a processing step  $Q = (S^Q, E^Q, N^Q) \rightarrow (S^{Q'}, E^{Q'}, N^{Q'})$  in  $\rho$  such that  $d \notin S^Q(c).oauthConfigCache$  and  $d \in S^{Q'}(c).oauthConfigCache$ . In  $Q$ , `PROCESS_HTTPS_RESPONSE` must have been called with a *reference* as second argument, such that  $reference[responseTo] \equiv \text{CONFIG}$ . In addition,  $reference[session]$  must contain a value  $sessionId$  such that  $S^Q(c).sessions[sessionId][selected\_AS] \equiv m.body[issuer]$  (Line 14 of Algorithm 2). From Line 16 of Algorithm 2, we also know that  $S^Q(c).sessions[sessionId][selected\_AS] \equiv d$  (cf. Lemma 15). Hence, we already have that  $d \in S(c).oauthConfigCache$  implies 1).

With Lemma 4, we have that there must be a processing step  $P = (S^P, E^P, N^P) \rightarrow (S^{P'}, E^{P'}, N^{P'})$  prior to  $Q$  in  $\rho$  in which  $c$  called `HTTPS_SIMPLE_SEND` with *reference* as first argument. Such a *reference* (one with `responseTo` set to `CONFIG`) is only created in Line 9 of Algorithm 6. The accompanying message's host value there is  $S^P(c).sessions[sessionId][selected\_AS]$ , i.e., by Lemma 15,  $d$ . That same message's path value is either `/.well_known/openid-configuration` or `/.well_known/oauth-authorization-server`. From Lemma 27, Algorithm 26, and Lines 10ff. of Algorithm 31 (and because  $as$  does not leak  $tlskey(d)$ ), we know that the request given to `HTTPS_SIMPLE_SEND` in  $P$  can only be answered by  $as$  (and  $c$ , but clients do not reply to requests with the aforementioned path values).

Such a request, i.e., one with the path values mentioned above, is processed by  $as$  in Lines 2ff. of Algorithm 9. From looking at those Lines, it is obvious that the response sent in Line 10 of Algorithm 9 contains a body with a dictionary fulfilling 2)–6). Using Lemma 27 once more, we can conclude that  $c$  processes such a response in  $Q$  and thus we have 2)–6). ■

*Lemma 17 (Authorization code secrecy).* For any run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of a FAPI web system  $\mathcal{FAPI}$  with a network attacker, every configuration  $(S, E, N)$  in  $\rho$ , every authorization server  $as \in \text{AS}$  that is honest in  $S$ , every client  $c \in \text{C}$  that is honest in  $S$  with client identifier  $clientId$  at  $as$ , every identity  $id \in \text{ID}^{as}$  with  $b = \text{ownerOfID}(id)$  being an honest browser in  $S$ , every authorization code  $code \neq \perp$  for which there is a record  $rec \in {}^\langle \rangle S(as).records$  with  $rec[code] \equiv code$ ,  $rec[client\_id] \equiv clientId$ , and  $rec[subject] \equiv id$  and every process  $p \notin \{as, c, b\}$ , it holds true that  $code \notin d_\emptyset(S(p))$ .

PROOF.

- 1) For  $code$  to end up in  $(S(as).records.x)[code]$  (with  $x \in \mathbb{N}$ ), the  $as$  has to execute Line 44 of Algorithm 9, since the only other places where an honest authorization server writes to the – initially empty, see Definition 17 – `records` state subterm are:
  - Line 123 of Algorithm 9: This line overwrites the stored authorization code with  $\perp$ , i.e., codes written by this line are not relevant to this lemma.
  - Line 130 of Algorithm 9 and Line 131 of Algorithm 9: In these two places, the authorization server does not modify the code entry. Note that  $ptr$  in these places cannot point “into” one of the records (see condition in Line 92 of Algorithm 9).
- 2) A  $code$  stored in Line 44 of Algorithm 9 is a fresh nonce (Line 43 of Algorithm 9). Hence, a  $code$  generated by  $as$  in that line in some processing step  $s_i \rightarrow s_{i+1}$  is not known to any process up to and including  $s_i$ . Let  $e_{in}$  be the event processed by  $as$  in  $s_i \rightarrow s_{i+1}$ . In order to reach Line 44 of Algorithm 9,  $e_{in}$  must contain an HTTPS request  $req$  to the `/auth2` endpoint. The only place in which an honest  $as$  sends out the  $code$  value is the HTTPS response to  $req$  – i.e., if the sender of  $req$  is honest, this response is only readable by the sender of  $req$ .
- 3) In addition,  $req$  must contain a valid *identity–password* combination – because  $as$  stores  $code$  along with *identity* and *clientId* only if  $password \equiv \text{secretOfID}(identity)$ . Since  $as$  does not send requests to itself and  $\text{secretOfID}(identity)$  is only known to  $as$  and  $\text{ownerOfID}(identity)$ ,  $req$  must have been created by  $\text{ownerOfID}(identity)$  if the sender of  $req$  is honest. W.l.o.g., let  $identity \equiv id$ , i.e.,  $req$  was created by  $b$ .
- 4) Since the origin header of  $req$  must be a domain of  $as$  and  $req$  must use the POST method, we know that  $req$  was initiated by a script of  $as$ . In particular,  $req$  must have been initiated by  $script\_as\_form$  (as this is the only script ever sent by  $as$ ). This script does not leak  $code$  after it is returned from  $as$ , since it uses a form post to transmit the credentials to  $as$ , and the window is subsequently navigated away. Instead,  $as$  provides an empty script in its response to  $req$  (Line 52 of Algorithm 9). This response contains a location redirect header. It is now crucial to check that this redirect does not leak  $code$  to any process except for  $c$ . The value of the location header is taken from  $S(as).authorizationRequests[requestUri][redirect\_uri]$

where  $S(as).authorizationRequests[requestUri][client\_id] \equiv clientId$ . With Lemma 14, we have that this URI is an HTTPS URI and belongs to  $c$ . We therefore know that  $b$  will send an HTTPS request containing  $code$  to  $c$ . We now have to check whether  $c$  or a script delivered by  $c$  to  $b$  will leak  $code$ . Algorithm 1 processes all HTTPS requests delivered to  $c$ . As  $as$  redirected  $b$  using the 303 status code, the request must be a GET request. Hence,  $c$  does not process this request in Lines 5ff. of Algorithm 1. If the request is processed in Lines 2ff. of Algorithm 1,  $c$  only responds with a script and does not use  $code$  at all. This leaves us with Lines 12ff. of Algorithm 1; here, the  $code$  value is (a) stored in the `sessions` state subterm and (b) given the `SEND_TOKEN_REQUEST` function. The value from (a) is not accessed anywhere, hence, it cannot leak. As for (b), we have to look at Algorithm 3. There, the  $code$  is included in the body of an HTTPS request under the key `code` (Line 6 of Algorithm 3).

- 5) The HTTPS request (“token request”) prepared in Lines 6ff. of Algorithm 3 is sent to the token endpoint of  $as$  (which was selected in  $b$ ’s initial request and is bound to the authorization response via the  $\langle \_Host, sessionId \rangle$  cookie – see Line 13 of Algorithm 1 and Line 30 of Algorithm 2). Since an honest client does not change the contents of an element of `oauthConfigCache` once it is initialized with the selected authorization server’s metadata (see Line 9 of Algorithm 6, Line 16 of Algorithm 2, and Lemma 16), the token endpoint to which the  $code$  is sent is the one provided by  $as$  at its metadata endpoint. As  $as$  is honest, the token endpoint returned by its metadata endpoint uses a domain which belongs to  $as$  and protocol  $S$ . With Lemma 27 we can conclude that the token request as such does not leak  $code$ .
- 6) As the token request is a HTTPS request sent to a domain of  $as$  and  $as$  is honest, only  $as$  can decrypt the request and extract  $code$ . Requests to the token endpoint are processed in Lines 81ff. of Algorithm 9, It is easy to see that the  $code$  is not stored or send out there, hence, it cannot leak. ■

**Lemma 18 (Unique Code Verifier for Each Login Session ID at Client).** For any run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of a FAPI web system  $\mathcal{F}API$  with a network attacker, every configuration  $(S^i, E^i, N^i)$  in  $\rho$ , every client  $c \in \mathcal{C}$  that is honest in  $S$  with client identifier  $clientId$  at  $as$ , every login session id  $lsid$ , and every term  $codeVerifier$ , we have that  $S^i(c).sessions[lsid][code\_verifier] \equiv codeVerifier$  implies:

- 1)  $S^j(c).sessions[lsid][code\_verifier] \equiv codeVerifier$  for all  $j \geq i$ , and
- 2)  $S^i(c).sessions[lsid'][code\_verifier] \not\equiv codeVerifier$  for all  $lsid' \neq lsid$ .

**PROOF.** We start by noting that an honest client only ever stores something in an entry in `sessions` under key `code_verifier` in Line 34 of Algorithm 6. The value stored there is always a fresh nonce (see Line 26 of Algorithm 6). Hence, we can conclude 2).

To get 1), we need to prove that a stored code verifier is never overwritten. For this, we show that a client executes Line 34 of Algorithm 6 at most once with the same login session id (i.e.,  $sessionId$  in the context of said line). For this, we look at the places where Algorithm 6 ( `PREPARE_AND_SEND_PAR`) is called. Note that the first argument to `PREPARE_AND_SEND_PAR` is the aforementioned  $sessionId$ :

**Line 11 of Algorithm 1** Here, the first argument is a fresh nonce (see Line 9 of Algorithm 1), i.e., this line will never lead to `PREPARE_AND_SEND_PAR` being called a second time with a given  $sessionId$ .

**Line 17 of Algorithm 2** This line is only executed when the client processes an HTTPS response such that Algorithm 2 ( `PROCESS_HTTPS_RESPONSE`) was called with a  $reference$  containing a key `responseTo` with value `CONFIG`. The  $sessionId$  value used when calling `PREPARE_AND_SEND_PAR` is also taken from the  $reference$  (see Line 10 of Algorithm 2). I.e., we have to check where this  $reference$  came from.

$reference$  is one of the arguments to `PROCESS_HTTPS_RESPONSE` which is only called in Line 26 of Algorithm 31, where the value for  $reference$  is taken from the client’s `pendingRequests` state subterm. The `pendingRequests` state subterm is initially empty (Definition 16) and the only place where elements are added to this state subterm is Line 15 of Algorithm 31. There, in turn, the value for  $reference$  is taken (unchanged) from an entry in the `pendingDNS` state subterm. Once again, this state subterm is initially empty and there is only one place in which entries are added to it: In Line 2 of Algorithm 26, i.e., in `HTTPS_SIMPLE_SEND`, where  $reference$  is one of the arguments. Hence, we have to look at places where `HTTPS_SIMPLE_SEND` is called with a  $reference$  where  $reference[responseTo] \equiv CONFIG$ .

The only place where such a  $reference$  is passed to `HTTPS_SIMPLE_SEND` is Line 9 of Algorithm 6. However, this call always ends in a **stop** and the call happens *before* the client executes Line 34 of Algorithm 6 – hence, if an execution of `PREPARE_AND_SEND_PAR` leads to execution of Line 11 of Algorithm 1 and thus a subsequent call of `PREPARE_AND_SEND_PAR`, both calls use the same  $sessionId$ , but Line 34 of Algorithm 6 (i.e., storing a code verifier) is executed at most once.

**Line 20 of Algorithm 2** This case is very similar to the previous one, except for the following changes: The `responseTo` value in question is `JWKS` instead of `CONFIG`, and the place in which `HTTPS_SIMPLE_SEND` is called with a suitable  $reference$  is Line 14 of Algorithm 6. ■

**Lemma 19 (Request URIs do not Leak).** For any run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of a FAPI web system  $\mathcal{F}API$  with a network attacker, every processing step  $Q = (S, E, N) \rightarrow (S', E', N')$  in  $\rho$ , every client identifier  $clientId$ , every authorization



server  $as \in AS$  that is honest in  $S$ , every client  $c \in C$  with client identifier  $clientId$  at  $as$  that is honest in  $S$ , every browser  $b \in B$  that is honest in  $S$ , every domain  $d_c \in \text{dom}(c)$ , every login session id  $lsid$ , every nonce  $codeVerifier$  with

- (a)  $\langle \langle \_Host, sessionId \rangle, \langle lsid, T, T, T \rangle \rangle \in S(b).cookies[d_c]$ , and
- (b)  $S(c).sessions[lsid][code\_verifier] \equiv codeVerifier$ , and
- (c)  $S(c).sessions[lsid][selected\_AS] \in \text{dom}(as)$ , and
- (d)  $c$  does not leak the authorization request for  $lsid$  (see [Definition 23](#)),

then all of the following hold true

- 1) there is exactly one nonce  $requestUri$ , such that  $S(as).authorizationRequests[requestUri][code\_challenge] \equiv \text{hash}(codeVerifier)$ , and
- 2) only  $b$ ,  $c$ , and  $as$  know  $requestUri$ , i.e., for all processes  $p \notin \{b, c, as\}$ , we have  $requestUri \notin d_0(S(p))$ .

**PROOF.** We start by noting that all requests and responses at an authorization server’s pushed authorization request (PAR) endpoint must be HTTPS requests, i.e., as long as the sender of the request and the authorization server in question are honest, the contents of request and response are not leaked by these messages as such (they may still leak by other means).

- (I) **hash( $codeVerifier$ ) does not leak.** We start off by showing that  $\text{hash}(codeVerifier)$  does not leak to any process other than  $c$  and  $as$ . For this, we look at how  $codeVerifier$  (from (b)) is generated and stored by  $c$ . The only place in which an honest client – such as  $c$  – stores a value under key  $code\_verifier$  in its session storage is in `PREPARE_AND_SEND_PAR` in Line 34 of Algorithm 6. That value is generated in the same function in Line 26 by using a fresh nonce. Hence, at this point,  $\text{hash}(codeVerifier)$  is only derivable by  $c$ . `PREPARE_AND_SEND_PAR` ends with the client sending a PAR request which contains  $\text{hash}(codeVerifier)$  under the key  $code\_challenge$ . So we have to check who receives that request. The PAR request is sent to the pushed authorization request endpoint of the authorization server stored under key  $selected\_AS$  under  $lsid$  in the client’s session storage. As an honest client never changes this value once it is set ([Lemma 16](#)), we know from (c) that the PAR request is sent to  $as$ . An honest authorization server – such as  $as$  – only reads a value stored under the key  $code\_challenge$  in an incoming message when processing a request to its `/par` endpoint (Lines 53ff. of Algorithm 9). There, the value stored under  $code\_challenge$  – i.e.,  $\text{hash}(codeVerifier)$  – is stored in an authorization request record in the authorization server’s authorization requests storage (see Lines 67, 77, and 78). Since  $as$  is honest, it never sends out the  $code\_challenge$  value (neither from the authorization requests storage, nor from the records storage to which the  $code\_challenge$  is copied in Line 44 of Algorithm 9). Hence, the value  $\text{hash}(codeVerifier)$  sent in the PAR request is not leaked “directly”.

However, this value would be derivable if  $codeVerifier$  leaks, i.e., we also have to prove that  $codeVerifier$  does not leak. As noted above, this value is a fresh nonce stored in  $c$ ’s session storage under the key  $code\_verifier$ . The only place in which a client accesses such a value is in function `SEND_TOKEN_REQUEST`, where the value is included in the body of an HTTPS request under the key  $code\_verifier$  (Lines 6f. of Algorithm 3) which is sent to the token endpoint of the authorization server stored under key  $selected\_AS$  under  $lsid$  in the client’s session storage – i.e.,  $as$  by (c) and [Lemma 16](#). Hence, this request in itself does not leak  $codeVerifier$ .

The only place in which an honest authorization server reads a value stored under the key  $code\_verifier$  from an incoming message is when processing a token request in Line 89 of Algorithm 9. This value is not stored by the authorization server, neither is it sent anywhere. Hence,  $codeVerifier$  does not leak.

- (II)  **$as$  stores  $\text{hash}(codeVerifier)$ .** Because the cookie from (a) includes the `__Host` prefix and  $b$  is honest, that cookie must have been set by  $c$ . Clients only ever set such cookies when processing PAR responses in Lines 21ff. of Algorithm 2. With (b) (note that a client will never change the value stored under  $code\_verifier$ , see [Lemma 18](#)), this implies that  $c$  sent a PAR request containing  $\text{hash}(codeVerifier)$  to  $as$  (see (I)) and got a response. Hence,  $as$  must have processed that PAR request as described above. Part of that processing is to store the value of  $code\_challenge$  from the request – i.e.,  $\text{hash}(codeVerifier)$  here – in the authorization request storage. Thus, we can conclude that there must be some  $requestUri'$  such that  $S(as).authorizationRequests[requestUri'][code\_challenge] \equiv \text{hash}(codeVerifier)$ .
- (III) **Proof for 1).** From (I), we have that only  $c$  and  $as$  know the value  $\text{hash}(codeVerifier)$  and do not use it in any request except for a single PAR request from  $c$  to  $as$ . From (II), we have that  $as$  stores  $\text{hash}(codeVerifier)$  as part of processing that PAR request. As  $as$  will use a fresh nonce as request URI for every processed PAR request (see Line 70 of Algorithm 9), and never changes the stored values (except for  $code$ ), we can conclude that there is exactly one  $requestUri$  such that  $S(as).authorizationRequests[requestUri][code\_challenge] \equiv \text{hash}(codeVerifier)$ .
- (IV) **Proof for 2).** As shown above,  $requestUri$  is a fresh nonce chosen and stored by  $as$  when processing a PAR request send by  $c$ .  $requestUri$  is not sent out by authorization servers anywhere, except in the response to the PAR request (under the key `request_uri`) that lead to the “creation” of  $requestUri$ .

Since we already established that the receiver of that PAR response is  $c$  (see above), we now have to check how  $c$  uses  $requestUri$ .  $c$  only reads a value stored under the key `request_uri` from an incoming message when processing the response to a PAR request. While  $c$  does store that value in its session storage, it never accesses that stored value. However, after

processing the PAR response,  $c$  constructs an authorization request containing  $requestUri$  as part of the query parameters (under key `request_uri`). That authorization request is a redirect which “points” to the authorization endpoint of the authorization server stored under key `selected_AS` under  $lsid$  in  $c$ ’s session storage (i.e.,  $as$  by (c)). By (d), we also know that  $c$  does not execute Line 36 of Algorithm 2, i.e., does not leak the authorization request for  $lsid$ .

Before looking at the receiver of the aforementioned redirect, we note that  $as$  only ever reads the value of a request parameter `request_uri` in Line 19 of Algorithm 9 – that value is neither stored, nor sent out by  $as$ .

The redirect sent out by  $c$  when processing the PAR response is an HTTPS response which – among other things – contains a Set-Cookie header with a cookie of the form  $\langle\langle\_Host, sessionId\rangle, \langle lsid, \top, \top, \top\rangle\rangle$ . We note that this is the only place in which  $c$  sets such a cookie.

Since we know from (a) that  $b$  knows such a cookie, and (II) implies that  $c$  must have set this cookie, we know that the HTTPS response containing the redirect with  $requestUri$ , sent by  $c$ , was sent to  $b$ .

We now only have to show that  $b$  does not leak  $requestUri$ . The aforementioned redirect contains a Location header (Line 29 of Algorithm 2) and status code 303, hence  $b$  will enter the location header handling in Line 11 of Algorithm 21 when processing that redirect (note that the redirect is sent by  $c$  with an empty script, i.e., no leakage through a script is possible). This handling will either end in a **stop** without any changes to  $b$ ’s state and no output event – which means that  $b$  does neither store, nor send out  $requestUri$  – or with a call of `HTTP_SEND` in Line 27 of Algorithm 21. While `HTTP_SEND` does store the message to be send (containing  $requestUri$ ), that stored value is only ever accessed when processing a DNS response and is then encrypted and sent out. We already established above that the redirection target is one of  $as$ ’s authorization endpoints and that  $as$  does not leak any  $requestUri$  values received there. Hence, we have that only  $b$ ,  $c$ , and  $as$  know  $requestUri$ , i.e., for all processes  $p \notin \{b, c, as\}$ , we have  $requestUri \notin d_0(S(p))$ . ■

## B. Authorization Property

In this section, we show that the authorization property from Definition 2 holds.

*Lemma 20 (Authorization).* For

- every run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  of  $\mathcal{F}\mathcal{A}\mathcal{P}\mathcal{I}$  with a network attacker,
- every resource server  $rs \in \text{RS}$  that is honest in  $S^n$ ,
- every identity  $id \in {}^\diamond s_0^{rs}.ids$  with  $b = \text{ownerOfID}(id)$  being an honest browser in  $S^n$ ,
- every processing step in  $\rho$

$$Q = (S^Q, E^Q, N^Q) \xrightarrow[rs \rightarrow E_{\text{out}}^Q]{e_{\text{in}}^Q \rightarrow rs} (S^{Q'}, E^{Q'}, N^{Q'})$$

- every  $resourceID \in \mathbb{S}$  with  $as = \text{authorizationServerOfResource}^{rs}(resourceID)$  being honest in  $S^Q$ ,

it holds true that:

If  $\exists r, x, y, k, m_{\text{resp}}. \langle x, y, \text{enc}_s(m_{\text{resp}}, k) \rangle \in {}^\diamond E_{\text{out}}^Q$  such that  $m_{\text{resp}}$  is an HTTP response,  $r := m_{\text{resp}}.body[\text{resource}]$ , and  $r \in {}^\diamond S^{Q'}(rs).resourceNonce[id][resourceID]$ , then

- 1) There exists a processing step

$$P = (S^P, E^P, N^P) \xrightarrow[rs \rightarrow E_{\text{out}}^P]{e_{\text{in}}^P \rightarrow rs} (S^{P'}, E^{P'}, N^{P'})$$

such that

- a) either  $P = Q$  or  $P$  prior to  $Q$  in  $\rho$ , and
  - b)  $e_{\text{in}}^P$  is an event  $\langle x, y, \text{enc}_a(\langle m_{\text{req}}, k_1 \rangle, k_2) \rangle$  for some  $x, y, k_1$ , and  $k_2$  where  $m_{\text{req}} \in \mathcal{T}_{\mathcal{N}}$  is an HTTP request which contains a term (access token)  $t$  in its `Authorization` header, i.e.,  $t \equiv m_{\text{req}}.headers[\text{Authorization}].2$ , and
  - c)  $r$  is a fresh nonce generated in  $P$  at the resource endpoint of  $rs$  in Line 46 of Algorithm 12.
- 2)  $t$  is bound to a key  $k \in \mathcal{T}_{\mathcal{N}}$ ,  $as$ , and  $id$  in  $S^Q$  (see Definition 1).
  - 3) If there exists a client  $c \in \mathbb{C}$  such that  $k \equiv \text{pub}(\text{signkey}(c))$  or  $k \equiv \text{pub}(\text{tlskey}(d_c))$  and  $d_c \in \text{dom}(c)$ , and if  $c$  is honest in  $S^n$ , then  $r$  is not derivable from the attackers knowledge in  $S^n$  (i.e.,  $r \notin d_0(S^n(\text{attacker}))$ ).

**PROOF. Resource server sends resource to correct client.** The first and the second postcondition are shown in Lemma 13, where we also showed that the message contained in the event  $e_{\text{in}}^P$  was created by  $c$  (as intuitively, the access token is bound to  $c$  via mTLS or DPoP, and no other process can prove possession of the secret key to which the token is bound). The resource  $r$  is sent back as a response to  $e_{\text{in}}^P$ : If the resource server sends out the resource in Line 69 of Algorithm 12, then it encrypts the HTTPS response (symmetrically) with the key contained in  $e_{\text{in}}^P$ . Otherwise, the resource server sends out the response in Line 22 of Algorithm 13, encrypted (symmetrically) with the key contained in  $e_{\text{in}}^P$  (the resource server stored the key in its state).

Thus, the resource server sends out the resource  $r$  back to  $c$ , encrypted with a symmetric key that only  $c$  and  $rs$  can derive. This response can only be decrypted by  $c$ : A resource server can decrypt symmetrically only in Line 19 of Algorithm 31 (i.e., in the generic server model), where the decryption key is taken from the `pendingRequests` state subterm. The application-layer model of a resource server does not access this state subterm, and the generic HTTPS server model stores only fresh nonces as keys (see Line 15 of Algorithm 31).

**Client never sends resource  $r$  to attacker.** In the following, we show that  $c$  does not send the resource nonce  $r$  to the attacker by contradiction, i.e., we assume that the client does send  $r$  to the attacker.

**Redirection request was created by attacker.** The client processes the response of the resource server (containing the resource  $r$ ) in Line 49 of Algorithm 2 (as a client sends out requests with an `Authorization` header only by calling `HTTPS_SIMPLE_SEND` in Line 30 of Algorithm 4 with the reference `RESOURCE_USAGE`) in some processing step  $R = (S^R, E^R, N^R) \xrightarrow[c \rightarrow E_{out}^R]{e_{in}^R \rightarrow c} (S^{R'}, E^{R'}, N^{R'})$  ( $R$  happens after  $Q$ ). The client stores the resource into its sessions in Line 51 of Algorithm 2, but never access it again in any other location. The client sends the resource as a response to a request  $req_{redir}$  stored in  $S^R.sessions[sessionId[redirectEpRequest]]$ , for some value  $sessionId$  (and in particular, encrypts the response with the key contained in  $req_{redir}$ ), see Line 11, Line 53, and Line 54 of Algorithm 2.

An honest client sets `redirectEpRequest` values only in the redirection endpoint in Line 12 of Algorithm 1, i.e.,  $req_{redir}$  is a request that was previously received by the client. This request contains a value (an authorization code) in  $req_{redir}.parameters[code]$  (which the client puts into the token request in Algorithm 3).

As we assume that the client sends  $r$  to the attacker, it follows that  $req_{redir}$  was created by the attacker, in particular, the attacker can derive the symmetric key and all other values in the request.

**Access token was sent by correct authorization server.** Before sending the resource request, the client ensures that it sent the token request to the correct authorization server, i.e., the authorization server managing the resource: The client sends resource requests only in Algorithm 4. In Line 7 of Algorithm 4, the client checks whether the input argument `tokenEPDomain` is a domain of the authorization server managing the resource that the client wants to request at the resource server. Algorithm 4 is called only in Line 44 of Algorithm 2, and the value `tokenEPDomain` is the domain of the token request, i.e., the client received the access token from `authorizationServerOfResourcers(resourceID)` (see Definition 16). This authorization server is honest, as required by the precondition of the lemma.

**Attacker can derive authorization code issued for honest client and id.** The access token that the client received in the token request is bound to its signing key or TLS key, the authorization server  $as$ , and the identity  $id$  (as shown in the second postcondition, it is bound to some key,  $as$ , and  $id$ , and the third postcondition has the requirement that the key is a key of the client). The authorization server created the access token in the token endpoint in Line 81 of Algorithm 9 in some processing step  $T = (S^T, E^T, N^T) \xrightarrow[as \rightarrow E_{out}^T]{e_{in}^T \rightarrow as} (S^{T'}, E^{T'}, N^{T'})$ . The token request contains a code  $code$  such that there is a record  $rec \in \langle \rangle S^T(as).records$  with  $rec[code] \equiv code$  and  $code \not\equiv \perp$  (Line 92 of Algorithm 9). Furthermore, the record has the following values:

- $rec[clientId]$  is a client identifier of  $c$  at  $as$ , as the token request was sent by  $c$
- $rec[subject] \equiv id$  (as the access token is bound to this identity)

As shown in Lemma 5, the code that the client uses is the same code that it received in the request to the redirection endpoint, i.e.,  $req_{redir}$ .

However, this is a contradiction to Lemma 17, i.e., such an authorization code cannot leak to the attacker. ■

### C. Authentication Property

In this section, we show that the authentication property from Definition 4 holds. This will be a proof by contradiction, i.e., we assume that there is a FAPI web system  $\mathcal{FAPI}$  in which the authentication property is violated and deduce a contradiction.

*Assumption 1.* There exists a FAPI web system with a network attacker  $\mathcal{FAPI}$  such that there exists a run  $\rho$  of  $\mathcal{FAPI}$  with a configuration  $(S, E, N)$  in  $\rho$ , some  $c \in \mathcal{C}$  that is honest in  $S$ , some identity  $id \in \text{ID}$  with  $as = \text{governor}(id)$  being an honest AS and  $b = \text{secretOfID}(id)$  being browser honest in  $S$ , some service session identified by some nonce  $n$  for  $id$  at  $c$ , and  $n$  is derivable from the attacker's knowledge in  $S$  (i.e.,  $n \in d_\emptyset(S(\text{attacker}))$ ).

*Lemma 21 (Authentication Property Holds).* Assumption 1 is a contradiction.

**PROOF.** By Assumption 1, there is a service session identified by  $n$  for  $id$  at  $c$ , and hence, by Definition 3, we have that there is a session id  $x$  and a domain  $d \in \text{dom}(\text{governor}(id))$  with  $S(c).sessions[x][loggedInAs] \equiv \langle d, id \rangle$  and  $S(c).sessions[x][serviceSessionId] \equiv n$ . Assumption 1 says that  $n$  is derivable from the attacker's knowledge. Since we have  $S(c).sessions[x][serviceSessionId] \equiv n$ , we can check where such an entry in  $c$ 's state can be created.

The only place in which an honest client stores a service session id is in the function `CHECK_ID_TOKEN`, specifically in Line 15 of Algorithm 5. There, the client chooses a fresh nonce as the value for the service session id, in this case  $n$ . In the line before, it sets the value for  $S(c).sessions[x][loggedInAs]$ , in this case  $\langle d, id \rangle$ .

`CHECK_ID_TOKEN`, in turn, is only called in a single place: When processing an HTTPS response to a token request, in Line 48 of Algorithm 2. From the check in Line 47 of Algorithm 2, we know that this response came from (one of)  $as$ 's token endpoints (cf. Lemma 16). Since  $as$  is an honest authorization server, it will only reply to a token request if that request contains a valid authorization code, i.e., the token request must contain a  $code$  such that there is a record  $rec \in {}^\diamond S(as).records$  with  $rec[code] \equiv code$ ,  $rec[client\_id] \equiv clientId$ , and  $rec[subject] \equiv id$  where  $clientId$  must be one of  $c$ 's identifiers at  $as$  (otherwise, client authentication would fail and  $as$  would not output an id token, see Lemma 10).

By tracking backwards from Line 15 of Algorithm 5, it is easy to see that the same party that finally receives the service session id  $n$  in an HTTPS response sent in Line 19 of Algorithm 5 must have sent an HTTPS request  $req$  to  $c$  containing the aforementioned  $code$ .

We now have to differentiate between two cases: Either (a) the sender of  $req$  is one of  $b, c, as$ ; or (b) the sender of  $req$  is any other process (except for  $b, c$ , and  $as$ ).

In case (a), we know that the only party sending an HTTPS request with an authorization code (i.e., with a body dictionary containing a key `code`) is  $b$ . If  $b$  sent  $req$ ,  $b$  receives the service session id  $n$  in a set-cookie header with the `httpOnly` and `secure` flags set (see Line 17 of Algorithm 5). Hence,  $b$  will only ever send  $n$  to  $c$  in a cookie header as part of HTTPS requests, which does not leak  $n$ . Neither does  $c$  leak received service session id cookie values – in fact,  $c$  never even accesses a cookie named `serviceSessionId`. Furthermore, neither  $b$ , nor  $c$  leak  $n$  in any other way (the value is not even accessed), resulting in a contradiction to Assumption 1.

In case (b), that other process which sent  $req$  would need to know  $code$  in order to be able to include it in  $req$ . This contradicts Lemma 17. ■

#### D. Session Integrity for Authentication Property

In this section, we show that the authentication property from Definition 24 holds.

*Assumption 2.* There exists a FAPI web system with a network attacker  $\mathcal{FAPI}$  such that there exists a run  $\rho$  of  $\mathcal{FAPI}$  with a processing step  $Q = (S, E, N) \rightarrow (S', E', N')$  in  $\rho$ , a browser  $b$  honest in  $S$ , an authorization server  $as \in AS$ , an identity  $id$ , a client  $c \in C$  honest in  $S$ , and a nonce  $lsid$  with  $\text{loggedIn}_\rho^Q(b, c, id, as, lsid)$  and  $c$  did not leak the authorization request for  $lsid$ , such that

- (1) there is no processing step  $Q'$  prior to  $Q$  in  $\rho$  such that  $\text{started}_\rho^{Q'}(b, c, lsid)$ , or
- (2)  $as$  is honest in  $S$ , and there is no processing step  $Q''$  prior to  $Q$  in  $\rho$  such that  $\text{authenticated}_\rho^{Q''}(b, c, id, as, lsid)$ .

*Lemma 22 (Session Integrity for Authentication Property Holds).* Assumption 2 is a contradiction.

**PROOF. (1).** We have that  $\text{loggedIn}_\rho^Q(b, c, id, as, lsid)$ . With Definition 19, we know that  $c$  sent out a service session id associated with  $lsid$  to  $b$ . This can only happen when the client's function `CHECK_ID_TOKEN` was called with  $lsid$  as the first argument – which, in turn, can only happen in Line 48 of Algorithm 2 when  $c$  processes a response to a token request. Such a response is only accepted by  $c$  if  $c$  sent a corresponding token request before. Clients only send token requests in Line 21 of Algorithm 1, when processing an HTTPS request  $req_{\text{redir}}$ .

$req_{\text{redir}}$  must contain a cookie  $[\langle \_Host, sessionId \rangle: lsid]$ : The response (containing the service session id) sent by  $c$  to  $b$  in Line 19 of Algorithm 5 is sent to and encrypted for the sender of  $req_{\text{redir}}$ , because  $c$  looks these values up in the login session record stored in  $S(c).sessions[lsid]$  under the key `redirectEpRequest`. Such an entry in  $c$ 's session storage is only ever created in Line 20 of Algorithm 1 when processing an HTTPS request containing a login session cookie as described above.

We can now track how that cookie was stored in  $b$ : Since the cookie is stored under a domain of  $c$  (otherwise,  $b$  would not include it in requests to  $c$ ) and the cookie is set with the `\_Host` prefix, the cookie must have been set by  $c$ . A cookie with the properties shown above is only set in Line 30 of Algorithm 2. Similar to the `redirectEpRequest` session entry above,  $c$  sends this cookie as a response to a stored request, in this case, using the key `startRequest` to determine receiver and encryption key. A session entry with key `startRequest` is only ever created in Line 10 of Algorithm 1. Hence, for  $b$  to receive the cookie, there must have been a request from  $b$  to  $c$  to the `/startLogin` endpoint, using the POST method, and with an origin header for an origin of  $c$  (see Line 6 of Algorithm 1).

Due to the origin check, this request must have been sent by a script under one of  $c$ 's origins. There is only one script which could potentially send such a request: `script_client_index`. Hence, there must be a processing step  $Q'$  (prior to  $Q$ ) in  $\rho$  in which  $b$  executed `script_client_index` and in that script, executed Line 8 of Algorithm 8.

In addition, we already established above that  $c$  replied to this request (stored under the key `startRequest`) with a response containing a header of the form  $\langle \text{Set-Cookie}, [\langle \_Host, sessionId \rangle: \langle lsid, \top, \top, \top \rangle] \rangle$ .

Hence, we have that  $\text{started}_\rho^{Q'}(b, c, lsid)$ .

(2). Again, we have  $\text{loggedIn}_\rho^Q(b, c, id, as, lsid)$  and we know that  $c$  sent out a service session id associated with  $lsid$  to  $b$ . This can only happen in the client's function `CHECK_ID_TOKEN`, which only produces an output if  $c$  received an id token  $t$  (via a token response). From  $S(c).\text{sessions}[lsid][\text{loggedInAs}] \equiv \langle d, id \rangle$ , we know that for  $t_c := \text{extractMsg}(t)$ , we have  $t_c[\text{iss}] \equiv d$ ,  $t_c[\text{sub}] \equiv id$ , and  $t_c[\text{aud}] \equiv \text{clientId}$  (for some  $\text{clientId}$ ). Due to the check in Line 47 of Algorithm 2, this id token must have been sent by  $as$  (because  $d \in \text{dom}(as)$ ).  $as$  will only output such a term  $t$  if there is a record  $rec$  in  $as$ 's records state subterm with  $rec[\text{subject}] \equiv id$ ,  $rec[\text{client\_id}] \equiv \text{clientId}$ ,  $rec[\text{code\_challenge}] \equiv \text{codeChallenge}$  (for some value of  $\text{codeChallenge}$ ). We note that  $as$  issuing an id token with  $t_c[\text{aud}] \equiv \text{clientId}$  implies that  $c$  has client identifier  $\text{clientId}$  at  $as$  (see Definition 13 and Definition 17).

By construction of  $c$  and tracking of  $\text{sessions}[lsid]$  in  $c$ 's state, it is easy to see that once  $c$  reaches `CHECK_ID_TOKEN`, the session storage  $S(c).\text{sessions}[lsid]$  must contain a key `code_verifier` under which a nonce  $\text{codeVerifier}$  is stored. We note that  $S(b).\text{cookies}[d_c]$  must contain a cookie  $\langle \langle \_Host, \text{sessionId} \rangle, \langle lsid, \top, \top, \top \rangle \rangle$  for  $d_c \in \text{dom}(c)$ , because  $b$  sends a cookie  $\langle \_Host, \text{sessionId} \rangle : lsid$  as explained above,  $b$  is honest (and will thus not accept `\_Host` headers for  $d_c$  from parties other than  $c$ ), and if  $c$  sets a cookie  $\langle \_Host, \text{sessionId} \rangle$ , it will do with the attributes set as shown here.

Hence, we can apply Lemma 19 (note that  $S(c).\text{sessions}[lsid][\text{loggedInAs}] \equiv \langle d, id \rangle$  with  $d \in \text{dom}(as)$  implies  $S(c).\text{sessions}[lsid][\text{selected\_AS}] \equiv d \in \text{dom}(as)$ ). I.e., we now have that there is exactly one  $\text{requestUri}$  such that  $S(as).\text{authorizationRequests}[\text{requestUri}][\text{code\_challenge}] \equiv \text{hash}(\text{codeVerifier})$ , and only  $b$ ,  $c$ , and  $as$  know  $\text{requestUri}$ .

We know from Line 95 of Algorithm 9 that the token request which leads to  $as$  issuing  $t$  must contain a code verifier such that  $\text{hash}(\text{codeVerifier}) \equiv \text{rec}[\text{code\_challenge}]$  (with  $rec$  from above). Since we know that  $c$  must have sent the token request (otherwise,  $c$  would not have received  $t$ ), we can track where and how  $c$  creates such a request. This is only the case in function `SEND_TOKEN_REQUEST`. There,  $c$  selects the value for the code verifier and based on the session id which  $c$  received from  $b$  via the `sessionId` cookie. At the same time,  $c$  includes the `code` from  $b$ 's request's parameters.

Going back to  $as$ , we can track where a  $rec$  as described above can be stored into  $as$ 's state: This is only the case at  $as$ 's `/auth2` endpoint (Lines 31ff. of Algorithm 9). There,  $as$  will only store a record  $rec$ , if there is an  $\text{authZrec}$ , stored under the key  $\text{reqUri}$  in the `authorizationRequests` state subterm such that there is an  $\text{auth2Reference}$  with  $\text{authZrec}[\text{auth2\_reference}] \equiv \text{auth2Reference}$  and that  $\text{auth2Reference}$  is contained in the request to  $as$ 's `/auth2` endpoint. Such an  $\text{auth2Reference}$ , in turn, is only created at  $as$ 's `/auth` endpoint. For a request to this endpoint to lead to storing  $\text{auth2Reference}$ , the request must contain  $\text{reqUri}$  under the key `request\_uri`.

Note that by Lemma 19, we established that there is exactly one  $\text{requestUri}$  in  $as$ 's state such that  $S(as).\text{authorizationRequests}[\text{requestUri}][\text{code\_challenge}] \equiv \text{hash}(\text{codeVerifier})$ . Therefore,  $\text{reqUri} \equiv \text{requestUri}$ . In addition, it is easy to see that  $c$  and  $as$  do not send any requests to  $as$ 's `/auth` endpoint. Hence,  $b$  must have sent a request with  $\text{reqUri}$  to `/auth`.

Since  $\text{auth2Reference}$  from above is only sent to whoever sent the first request to `/auth` (and – if  $b$  receives it –  $b$  does not leak that value) we know that  $b$  must have sent the request to `/auth2` as well. As  $b$  is honest, this can only happen through a script – together with the origin header check in Line 31 of Algorithm 9, and  $\text{script\_as\_form}$  being the only script ever sent by  $as$ , we can conclude that there must have been a processing step  $Q''$  prior to  $Q'$  in  $\rho$  in which  $b$  was triggered, selected a document under one of  $as$ 's origins with script  $\text{script\_as\_form}$ , executed that script, selected  $id$  from its identities (because we know from above that  $rec[\text{subject}] \equiv id$  and such a  $rec$  is only stored at `/auth2` endpoint, if the identity in the request is equivalent to  $id$ ) and sent a request to  $as$ 's `/auth2` endpoint containing  $\text{auth2Reference}$  – hence, the  $\text{scriptstate}$  contained a key `auth2\_reference` with value  $\text{auth2Reference}$ .

Hence, we have  $\text{authenticated}_\rho^{Q''}(b, c, id, as, lsid)$  which – together with (1). from above – contradicts Assumption 2, therefore proving the lemma. ■

### E. Session Integrity for Authorization Property

In this section, we show that the session integrity property from Definition 25 holds.

*Lemma 23 (Session Integrity for Authorization Property Holds).* For every run  $\rho$  of  $\mathcal{F}\mathcal{A}\mathcal{P}\mathcal{I}$ , every processing step  $Q = (S, E, N) \rightarrow (S', E', N')$  in  $\rho$ , every browser  $b$  that is honest in  $S$ , every  $as \in \text{AS}$ , every identity  $u$ , every client  $c \in \mathcal{C}$  that is honest in  $S$ , every  $rs \in \text{RS}$  that is honest in  $S$ , every nonce  $r$ , every nonce  $lsid$ , we have that if  $\text{accessesResource}_\rho^Q(b, r, u, c, rs, as, lsid)$  and  $c$  did not leak the authorization request for  $lsid$  (see Definition 23), then (1) there exists a processing step  $Q'$  in  $\rho$  (before  $Q$ ) such that  $\text{started}_\rho^{Q'}(b, c, lsid)$ , and (2) if  $as$  is honest in  $S$ , then there exists a processing step  $Q''$  in  $\rho$  (before  $Q$ ) such that  $\text{authenticated}_\rho^{Q''}(b, c, u, as, lsid)$ .

**PROOF. (1).** Due to  $\text{accessesResource}_\rho^Q(b, r, u, c, rs, as, lsid)$ , it holds true that the browser  $b$  has a `sessionId` cookie with the session identifier  $lsid$  for the domain of the client  $c$ . This cookie is set with the `\_Host` prefix, i.e., it follows that the cookie was set by  $c$ , which responds with a `Set-Cookie` (with `sessionId`) only in Line 30 of Algorithm 2. The remaining proof is analogous

to the proof of the first postcondition of Lemma 22.

(2).

**Client received resource from  $rs$ .** As the client executes Line 55 of Algorithm 2 (precondition of the lemma), and as

$S'(c).sessions[lsid][resourceServer] \in \text{dom}(rs)$  is only set in Line 52 of Algorithm 2, it follows that  $c$  received the resource  $r$  in a response from  $rs$ , i.e., it sent the corresponding resource request to  $rs$ .

**Resource request contains access token associated with  $u$  at  $as$ .** An honest resource server sends out an HTTP response  $resp_{resource}$  with  $resource \in \langle \rangle resp_{resource}.body$  either in Line 69 of Algorithm 12 or Line 22 of Algorithm 13. As shown in the proof of Lemma 20, the resource server received a resource request  $req_{resource}$  containing an access token  $t$  (either in the same processing step when storing the resource or in a previous processing step). Furthermore, as the resource server stores the resource in  $S(rs).resourceNonces[u][resourceId]$  (with  $resourceId \in \mathcal{T}_{\mathcal{R}}(c)$ ), it follows that  $req_{resource}$  has the path  $resourceId$ . Thus, it follows that the value  $responsibleAS$  chosen by the resource server in Line 15 of Algorithm 12 is a domain of  $as$  (as the resource server never changes the  $resourceASMapping$  subterm of its state, see also Definition 18).

If  $rs$  returns the resource  $r$  in Line 69 of Algorithm 12, then the access token is a structured JWT signed by  $as$  (Line 60 of Algorithm 12) and containing the sub value  $u$  (Line 62 of Algorithm 12). Otherwise, if  $r$  is returned in Line 22 of Algorithm 13, then the resource server received a response from  $as$  containing the sub value  $u$  and asserting that the access token contained in  $req_{resource}$  is valid. In both cases (structured access token or opaque token) it follows that the authorization server  $as$  has a sequence  $rec$  in the  $records$  subterm of its state with  $rec[access\_token] \equiv t$  and  $rec[subject] \equiv u$ .

**Token request was sent to  $as$ .** An honest client sends resource requests only in Algorithm 4, which is called only in Line 44 of Algorithm 2, i.e., after receiving the token response. The check in Line 7 of Algorithm 4 ensures that the token request  $req_{token}$  was sent to  $as$  (as the client calls Algorithm 4 with the domain of the token request, see also Definition 16). From this, it follows that  $S(c).sessions[lsid][selected\_AS]$  is a domain of  $as$ , as the client sends the token request to this domain, see Line 4, Line 10, and Line 11 of Algorithm 3.

**PAR request was sent to  $as$ .** The token request  $req_{token}$  sent from  $c$  to  $as$  contains an authorization code  $code$  and a PKCE code verifier  $pkce\_cv$  (see Line 6 of Algorithm 3). As the authorization server responds with an access token, it follows that the checks at the token endpoint in Line 81 of Algorithm 9 passed successfully. In particular, this implies that the token request contains the correct PKCE verifier for the code, i.e., the authorization code and the PKCE challenge corresponding to the PKCE verifier were stored in the same record entry in the  $records$  state subterm (see Line 89 and Line 95 of Algorithm 9). An authorization server adds these records to its  $records$  state subterm only in Line 44 of Algorithm 9, and the sequence that is added is taken from the  $authorizationRequests$  state subterm, see Line 40 of Algorithm 9. In this processing step, the authorization server also creates the authorization code (Line 43 of Algorithm 9) and associates the identity with the code (Line 41 of Algorithm 9).

Thus, as the authorization server  $as$  exchanged the authorization code  $code$  at the token endpoint and the issued access token is associated with the identity  $u$ , it follows that identity  $u$  logged in at the  $/auth2$  endpoint of  $as$ , and the request to  $/auth2$  contained a value  $auth2\_reference$  in its body equal to  $S''(as).authorizationRequests[requestUri][auth2\_reference]$  (with  $S''$  being the state of a configuration prior to  $Q$ ; see also Line 39 of Algorithm 9). The authorization server received the  $requestUri$  value at the auth endpoint, i.e., the process that can derive the request URI value can also derive the  $auth2$  reference.

As  $S(c).sessions[lsid][selected\_AS]$  is a domain of  $as$ , it follows that the client sent the pushed authorization request to  $as$  in Line 36 of Algorithm 6 in a previous processing step of the trace. In this processing step, the client chose the PKCE verifier  $pkce\_cv$  in Line 26 of Algorithm 6 and stored this value into the  $lsid$  session in Line 34 of Algorithm 6.

Now, we can apply Lemma 19 and conclude that the request URI can only be derived by  $b$ ,  $c$ , and  $as$ . As  $as$  does not send requests to itself and  $c$  does not send any request to an  $auth$  endpoint, it follows that the request to the  $auth$  endpoint of  $as$  was sent by  $b$ . The remaining argumentation is the same as for the proof of Lemma 22. ■

## F. Proof of Theorem

Theorem 1 follows from Lemma 20, Lemma 21, Lemma 22, and Lemma 23.

$$\text{dec}_a(\text{enc}_a(x, \text{pub}(y)), y) = x \quad (6)$$

$$\text{dec}_s(\text{enc}_s(x, y), y) = x \quad (7)$$

$$\text{checksig}(\text{sig}(x, y), \text{pub}(y)) = \top \quad (8)$$

$$\text{extractmsg}(\text{sig}(x, y)) = x \quad (9)$$

$$\text{checkmac}(\text{mac}(x, y), y) = \top \quad (10)$$

$$\text{extractmsg}(\text{mac}(x, y)) = x \quad (11)$$

$$\pi_i(\langle x_1, \dots, x_n \rangle) = x_i \text{ if } 1 \leq i \leq n \quad (12)$$

$$\pi_j(\langle x_1, \dots, x_n \rangle) = \diamond \text{ if } j \notin \{1, \dots, n\} \quad (13)$$

$$\pi_j(t) = \diamond \text{ if } t \text{ is not a sequence} \quad (14)$$

**Figure 8.** Equational theory for  $\Sigma$ .

## APPENDIX G TECHNICAL DEFINITIONS

Here, we provide technical definitions of the WIM. These follow the descriptions in [28, 33–38].

### A. Terms and Notations

**Definition 26 (Signature  $\Sigma$ ).** We define the signature  $\Sigma$ , over which we will define formal terms, as the union of the following pairwise disjoint sets:

**Constants**  $\mathcal{C} = \mathbb{S} \cup \text{IPs} \cup \{\perp, \top, \diamond\}$  with the three sets pairwise disjoint.  $\mathbb{S}$  is the set of all (ASCII) strings, including the empty string  $\varepsilon$ .  $\text{IPs}$  is the set of IP addresses.

**Function Symbols** to represent public keys, asymmetric encryption and decryption, symmetric encryption and decryption, signatures, signature verification, MACs, MAC verification, message extraction from signatures and MACs, and hashing, respectively:  $\text{pub}(\cdot)$ ,  $\text{enc}_a(\cdot, \cdot)$ ,  $\text{dec}_a(\cdot, \cdot)$ ,  $\text{enc}_s(\cdot, \cdot)$ ,  $\text{dec}_s(\cdot, \cdot)$ ,  $\text{sig}(\cdot, \cdot)$ ,  $\text{checksig}(\cdot, \cdot)$ ,  $\text{mac}(\cdot, \cdot)$ ,  $\text{checkmac}(\cdot, \cdot)$ ,  $\text{extractmsg}(\cdot)$ ,  $\text{hash}(\cdot)$ .

**Sequences** of any length  $\langle \cdot \rangle$ ,  $\langle \cdot, \cdot \rangle$ ,  $\langle \cdot, \cdot, \cdot \rangle$ , etc. Note that formally, these sequence “constructors” are also function symbols.

**Projection Symbols** to access sequence elements:  $\pi_i(\cdot)$  for all  $i \in \mathbb{N}_0$ . Note that formally, projection symbols are also function symbols.

**Definition 27 (Nonces and Terms).** By  $X = \{x_0, x_1, \dots\}$  we denote a set of variables and by  $\mathcal{N}$  we denote an infinite set of constants (*nonces*) such that  $\Sigma$ ,  $X$ , and  $\mathcal{N}$  are pairwise disjoint. For  $N \subseteq \mathcal{N}$ , we define the set  $\mathcal{T}_N(X)$  of *terms* over  $\Sigma \cup N \cup X$  inductively as usual: (1) If  $t \in N \cup X \cup \mathcal{C}$ , then  $t$  is a term. (2) If  $f \in \Sigma$  is an  $n$ -ary function symbol for some  $n \geq 0$  and  $t_1, \dots, t_n$  are terms, then  $f(t_1, \dots, t_n)$  is a term.

By  $\equiv$  we denote the congruence relation on  $\mathcal{T}_{\mathcal{N}}(X)$  induced by the theory associated with  $\Sigma$  (see Figure 8). For example, we have that  $\pi_1(\text{dec}_a(\text{enc}_a(\langle a, b \rangle, \text{pub}(k)), k)) \equiv a$ .

**Definition 28 (Ground Terms, Messages, Placeholders, Protomessages).** By  $\mathcal{T}_N = \mathcal{T}_N(\emptyset)$ , we denote the set of all terms over  $\Sigma \cup N$  without variables, called *ground terms*. The set  $\mathcal{M}$  of messages (over  $\mathcal{N}$ ) is defined to be the set of ground terms  $\mathcal{T}_{\mathcal{N}}$ .

We define the set  $V_{\text{process}} = \{\nu_1, \nu_2, \dots\}$  of variables (called placeholders). The set  $\mathcal{M}^\nu := \mathcal{T}_{\mathcal{N}}(V_{\text{process}})$  is called the set of *protomessages*, i.e., messages that can contain placeholders.

**Example 1.** For example,  $k \in \mathcal{N}$  and  $\text{pub}(k)$  are messages, where  $k$  typically models a private key and  $\text{pub}(k)$  the corresponding public key. For constants  $a, b, c$  and the nonce  $k \in \mathcal{N}$ , the message  $\text{enc}_a(\langle a, b, c \rangle, \text{pub}(k))$  is interpreted to be the message  $\langle a, b, c \rangle$  (the sequence of constants  $a, b, c$ ) encrypted by the public key  $\text{pub}(k)$ .

**Definition 29 (Events and Protoevents).** An *event* (over IPs and  $\mathcal{M}$ ) is a term of the form  $\langle a, f, m \rangle$ , for  $a, f \in \text{IPs}$  and  $m \in \mathcal{M}$ , where  $a$  is interpreted to be the receiver address and  $f$  is the sender address. We denote by  $\mathcal{E}$  the set of all events. Events over IPs and  $\mathcal{M}^\nu$  are called *protoevents* and are denoted  $\mathcal{E}^\nu$ . By  $2^{\mathcal{E}^\nu}$  (or  $2^{\mathcal{E}^\nu \downarrow}$ , respectively) we denote the set of all sequences of (proto)events, including the empty sequence (e.g.,  $\langle \rangle$ ,  $\langle \langle a, f, m \rangle, \langle a', f', m' \rangle, \dots \rangle$ , etc.).

**Definition 30 (Normal Form).** Let  $t$  be a term. The *normal form* of  $t$  is acquired by reducing the function symbols from left to right as far as possible using the equational theory shown in Figure 8. For a term  $t$ , we denote its normal form as  $t \downarrow$ .

**Definition 31 (Pattern Matching).** Let *pattern*  $\in \mathcal{T}_{\mathcal{N}}(\{*\})$  be a term containing the wildcard (variable  $*$ ). We say that a term  $t$  *matches pattern* iff  $t$  can be acquired from *pattern* by replacing each occurrence of the wildcard with an arbitrary term (which may

be different for each instance of the wildcard). We write  $t \sim pattern$ . For a sequence of patterns  $patterns$  we write  $t \dot{\sim} patterns$  to denote that  $t$  matches at least one pattern in  $patterns$ .

For a term  $t'$  we write  $t'|pattern$  to denote the term that is acquired from  $t'$  by removing all immediate subterms of  $t'$  that do not match  $pattern$ .

*Example 2.* For example, for a pattern  $p = \langle \top, * \rangle$  we have that  $\langle \top, 42 \rangle \sim p$ ,  $\langle \perp, 42 \rangle \not\sim p$ , and

$$\langle \langle \perp, \top \rangle, \langle \top, 23 \rangle, \langle a, b \rangle, \langle \top, \perp \rangle \rangle | p = \langle \langle \top, 23 \rangle, \langle \top, \perp \rangle \rangle .$$

*Definition 32 (Variable Replacement).* Let  $N \subseteq \mathcal{N}$ ,  $\tau \in \mathcal{T}_N(\{x_1, \dots, x_n\})$ , and  $t_1, \dots, t_n \in \mathcal{T}_N$ .

By  $\tau[t_1/x_1, \dots, t_n/x_n]$  we denote the (ground) term obtained from  $\tau$  by replacing all occurrences of  $x_i$  in  $\tau$  by  $t_i$ , for all  $i \in \{1, \dots, n\}$ .

*Definition 33 (Sequence Notations).* Let  $t = \langle t_1, \dots, t_n \rangle$  and  $r = \langle r_1, \dots, r_m \rangle$  be sequences,  $s$  a set, and  $x, y$  terms. We define the following operations:

- $t \subset^{\langle \rangle} s \iff t_1, \dots, t_n \in s$
- $x \in^{\langle \rangle} t \iff \exists i: t_i = x$
- $t +^{\langle \rangle} y := \langle t_1, \dots, t_n, y \rangle$
- $t \cup r := \langle t_1, \dots, t_n, r_1, \dots, r_m \rangle$
- $t -^{\langle \rangle} y := \begin{cases} \langle t_1, \dots, t_{i-1}, t_{i+1}, \dots, t_n \rangle & \text{if } \exists i: t_i = x \text{ (i.e., } y \in^{\langle \rangle} t) \\ t & \text{otherwise (i.e., } y \notin^{\langle \rangle} t) \end{cases}$

If  $y$  occurs more than once in  $t$ ,  $t -^{\langle \rangle} y$  non-deterministically removes one of the occurrences.

- $t -^{\langle \rangle * } y$  is  $t$  with all occurrences of  $y$  removed.
- $|t| := n$ . If  $t'$  is not a sequence, we set  $|t'| := \diamond$ .
- For a finite set  $M$  with  $M = \{m_1, \dots, m_n\}$  we use  $\langle M \rangle$  to denote the term of the form  $\langle m_1, \dots, m_n \rangle$ . The order of the elements does not matter; one is chosen arbitrarily.

*Definition 34 (Dictionaries).* A dictionary over  $X$  and  $Y$  is a term of the form

$$\langle \langle k_1, v_1 \rangle, \dots, \langle k_n, v_n \rangle \rangle$$

where  $k_1, \dots, k_n \in X$ ,  $v_1, \dots, v_n \in Y$ . We call every term  $\langle k_i, v_i \rangle$ ,  $i \in \{1, \dots, n\}$ , an *element* of the dictionary with key  $k_i$  and value  $v_i$ . We often write  $[k_1: v_1, \dots, k_n: v_n]$  instead of  $\langle \langle k_1, v_1 \rangle, \dots, \langle k_n, v_n \rangle \rangle$ . We denote the set of all dictionaries over  $X$  and  $Y$  by  $[X \times Y]$ . Note that the empty dictionary is equivalent to the empty sequence, i.e.,  $\square = \langle \rangle$ ; and dictionaries as such may contain duplicate keys (however, all dictionary operations are only defined on dictionaries with unique keys).

*Definition 35 (Operations on Dictionaries).* Let  $z = [k_1: v_1, k_2: v_2, \dots, k_n: v_n]$  be a dictionary with unique keys, i.e.,  $\forall i, j: k_i \neq k_j$ . In addition, let  $t$  and  $v$  be terms. We define the following operations:

- $t \in z \iff \exists i \in \{1, \dots, n\}: k_i = t$
- $z[t] := \begin{cases} v_i & \text{if } \exists k_i \in z: t = k_i \\ \langle \rangle & \text{otherwise (i.e., if } t \notin z) \end{cases}$
- $z - t := \begin{cases} [k_1: v_1, \dots, k_{i-1}: v_{i-1}, k_{i+1}: v_{i+1}, \dots, k_n: v_n] & \text{if } \exists k_i \in z: t = k_i \\ z & \text{otherwise (i.e., if } t \notin z) \end{cases}$
- In our algorithm descriptions, we often write **let**  $z[t] := v$ . If  $t \notin z$  prior to this operation, an element  $\langle t, v \rangle$  is appended to  $z$ . Otherwise, i.e., if there already is an element  $\langle t, x \rangle \in^{\langle \rangle} z$ , this element is updated to  $\langle t, v \rangle$ .

We emphasize that these operations are only defined on dictionaries with unique keys.

Given a term  $t = \langle t_1, \dots, t_n \rangle$ , we can refer to any subterm using a sequence of integers. The subterm is determined by repeated application of the projection  $\pi_i$  for the integers  $i$  in the sequence. We call such a sequence a *pointer*:

*Definition 36 (Pointers).* A pointer is a sequence of non-negative integers. We write  $\tau.\bar{p}$  for the application of the pointer  $\bar{p}$  to the term  $\tau$ . This operator is applied from left to right. For pointers consisting of a single integer, we may omit the sequence braces for brevity.

*Example 3.* For the term  $\tau = \langle a, b, \langle c, d, \langle e, f \rangle \rangle \rangle$  and the pointer  $\bar{p} = \langle 3, 1 \rangle$ , the subterm of  $\tau$  at the position  $\bar{p}$  is  $c = \pi_1(\pi_3(\tau))$ . Also,  $\tau.3.\langle 3, 1 \rangle = \tau.3.\bar{p} = \tau.3.3.1 = e$ .

To improve readability, we try to avoid writing, e.g.,  $o.2$  or  $\pi_2(o)$  in this document. Instead, we will use the names of the components of a sequence that is of a defined form as pointers that point to the corresponding subterms. E.g., if an *Origin* term is defined as  $\langle host, protocol \rangle$  and  $o$  is an Origin term, then we can write  $o.protocol$  instead of  $\pi_2(o)$  or  $o.2$ . See also Example 4.



*Definition 37 (Concatenation of Sequences).* For a sequence  $a = \langle a_1, \dots, a_i \rangle$  and a sequence  $b = \langle b_1, b_2, \dots \rangle$ , we define the concatenation as  $a \cdot b := \langle a_1, \dots, a_i, b_1, b_2, \dots \rangle$ .

*Definition 38 (Subtracting from Sequences).* For a sequence  $X$  and a set or sequence  $Y$  we define  $X \setminus Y$  to be the sequence  $X$  where for each element in  $Y$ , a non-deterministically chosen occurrence of that element in  $X$  is removed.

## B. Message and Data Formats

We now provide some more details about data and message formats that are needed for the formal treatment of the web model presented in the following.

### 1) URLs:

*Definition 39.* A URL is a term of the form

$$\langle \text{URL}, \text{protocol}, \text{host}, \text{path}, \text{parameters}, \text{fragment} \rangle$$

with  $\text{protocol} \in \{\text{P}, \text{S}\}$  (for plain (HTTP) and secure (HTTPS)), a domain  $\text{host} \in \text{Doms}$ ,  $\text{path} \in \mathbb{S}$ ,  $\text{parameters} \in [\mathbb{S} \times \mathcal{T}_{\mathcal{N}}]$ , and  $\text{fragment} \in \mathcal{T}_{\mathcal{N}}$ . The set of all valid URLs is URLs.

The *fragment* part of a URL can be omitted when writing the URL. Its value is then defined to be  $\perp$ . We sometimes also write  $\text{URL}_{\text{path}}^{\text{host}}$  to denote the URL  $\langle \text{URL}, \text{S}, \text{host}, \text{path}, \langle \rangle, \perp \rangle$ .

As mentioned above, for specific terms, such as URLs, we typically use the names of its components as pointers (see Definition 36):

*Example 4.* For the URL  $u = \langle \text{URL}, a, b, c, d \rangle$ ,  $u.\text{protocol} = a$ . If, in the algorithms described later, we say  $u.\text{path} := e$  then  $u = \langle \text{URL}, a, b, c, e \rangle$  afterwards.

### 2) Origins:

*Definition 40.* An origin is a term of the form  $\langle \text{host}, \text{protocol} \rangle$  with  $\text{host} \in \text{Doms}$  and  $\text{protocol} \in \{\text{P}, \text{S}\}$ . We write Origins for the set of all origins.

*Example 5.* For example,  $\langle \text{F00}, \text{S} \rangle$  is the HTTPS origin for the domain F00, while  $\langle \text{BAR}, \text{P} \rangle$  is the HTTP origin for the domain BAR.

### 3) Cookies:

*Definition 41.* A cookie is a term of the form  $\langle \text{name}, \text{content} \rangle$  where  $\text{name} \in \mathcal{T}_{\mathcal{N}}$ , and  $\text{content}$  is a term of the form  $\langle \text{value}, \text{secure}, \text{session}, \text{httpOnly} \rangle$  where  $\text{value} \in \mathcal{T}_{\mathcal{N}}$ ,  $\text{secure}, \text{session}, \text{httpOnly} \in \{\top, \perp\}$ . As  $\text{name}$  is a term, it may also be a sequence consisting of two parts. If the name consists of two parts, we call the first part of the sequence (i.e.,  $\text{name}.1$ ) the *prefix* of the name. We write Cookies for the set of all cookies and Cookies<sup>v</sup> for the set of all cookies where names and values are defined over  $\mathcal{T}_{\mathcal{N}}(V)$ .

If the *secure* attribute of a cookie is set, the browser will not transfer this cookie over unencrypted HTTP connections.<sup>8</sup> If the *session* flag is set, this cookie will be deleted as soon as the browser is closed. The *httpOnly* attribute controls whether scripts have access to this cookie.

When the `__Host` prefix (see [15]) of a cookie is set (i.e.,  $\text{name}$  consists of two parts and  $\text{name}.1 \equiv \text{__Host}$ ), the browser accepts the cookie only if the *secure* attribute is set. As such cookies are only transferred over secure channels (i.e., with TLS), the cookie cannot be set by a network attacker. Note that the WIM does not model the domain attribute of the Set-Cookie header, so cookies in the WIM are always sent to the originating domain and not some subdomain. Therefore, the WIM models only the `__Host` prefix, but not the `__Secure` prefix.

Also note that cookies of the form described here are only contained in HTTP(S) responses. In HTTP(S) requests, only the components  $\text{name}$  and  $\text{value}$  are transferred as a pairing of the form  $\langle \text{name}, \text{value} \rangle$ .

### 4) HTTP Messages:

*Definition 42.* An HTTP request is a term of the form shown in (15). An HTTP response is a term of the form shown in (16).

$$\langle \text{HTTPReq}, \text{nonce}, \text{method}, \text{host}, \text{path}, \text{parameters}, \text{headers}, \text{body} \rangle \quad (15)$$

$$\langle \text{HTTPResp}, \text{nonce}, \text{status}, \text{headers}, \text{body} \rangle \quad (16)$$

The components are defined as follows:

- $\text{nonce} \in \mathcal{N}$  serves to map each response to the corresponding request.
- $\text{method} \in \text{Methods}$  is one of the HTTP methods.

<sup>8</sup>Note that *secure* cookies can be set over unencrypted connections (c.f. RFC 6265).

- $host \in \text{Doms}$  is the host name in the HOST header of HTTP/1.1.
- $path \in \mathbb{S}$  indicates the resource path at the server side.
- $status \in \mathbb{S}$  is the HTTP status code (i.e., a number between 100 and 505, as defined by the HTTP standard).
- $parameters \in [\mathbb{S} \times \mathcal{T}_{\mathcal{N}}]$  contains URL parameters.
- $headers \in [\mathbb{S} \times \mathcal{T}_{\mathcal{N}}]$  contains request/response headers. The dictionary elements are terms of one of the following forms:
  - $\langle \text{Origin}, o \rangle$  where  $o$  is an origin,
  - $\langle \text{Set-Cookie}, c \rangle$  where  $c$  is a sequence of cookies,
  - $\langle \text{Cookie}, c \rangle$  where  $c \in [\mathbb{S} \times \mathcal{T}_{\mathcal{N}}]$  (note that in this header, only names and values of cookies are transferred),
  - $\langle \text{Location}, l \rangle$  where  $l \in \text{URLs}$ ,
  - $\langle \text{Referer}, r \rangle$  where  $r \in \text{URLs}$ ,
  - $\langle \text{Strict-Transport-Security}, \top \rangle$ ,
  - $\langle \text{Authorization}, \langle username, password \rangle \rangle$  where  $username, password \in \mathbb{S}$  (this header models the ‘Basic’ HTTP Authentication Scheme, see [67]),
  - $\langle \text{ReferrerPolicy}, p \rangle$  where  $p \in \{\text{noreferrer}, \text{origin}\}$ .
- $body \in \mathcal{T}_{\mathcal{N}}$  in requests and responses.

We write HTTPRequests/HTTPResponses for the set of all HTTP requests or responses, respectively.

*Example 6 (HTTP Request and Response).*

$$r := \langle \text{HTTPReq}, n_1, \text{POST}, \text{example.com}, /show, \langle \langle \text{index}, 1 \rangle \rangle, [\text{Origin} : \langle \text{example.com}, \mathbb{S} \rangle], \langle \text{foo}, \text{bar} \rangle \rangle \quad (17)$$

$$s := \langle \text{HTTPResp}, n_1, 200, \langle \langle \text{Set-Cookie}, \langle \langle \text{SID}, \langle n_2, \perp, \perp, \top \rangle \rangle \rangle \rangle, \langle \text{somescript}, x \rangle \rangle \quad (18)$$

An HTTP POST request for the URL <http://example.com/show?index=1> is shown in (17), with an Origin header and a body that contains  $\langle \text{foo}, \text{bar} \rangle$ . A possible response is shown in (18), which contains an httpOnly cookie with name SID and value  $n_2$  as well as a string `somescript` representing a script that can later be executed in the browser (see Section G-K) and the scripts initial state  $x$ .

*a) Encrypted HTTP Messages:* For HTTPS, requests are encrypted using the public key of the server. Such a request contains an (ephemeral) symmetric key chosen by the client that issued the request. The server is supposed to encrypt the response using the symmetric key.

*Definition 43.* An *encrypted HTTP request* is of the form  $\text{enc}_a(\langle m, k' \rangle, k)$ , where  $k \in \text{terms}$ ,  $k' \in \mathcal{N}$ , and  $m \in \text{HTTPRequests}$ . The corresponding *encrypted HTTP response* would be of the form  $\text{enc}_s(m', k')$ , where  $m' \in \text{HTTPResponses}$ . We call the sets of all encrypted HTTP requests and responses HTTPSRequests or HTTPSResponses, respectively.

We say that an HTTP(S) response matches or corresponds to an HTTP(S) request if both terms contain the same nonce.

*Example 7.*

$$\text{enc}_a(\langle r, k' \rangle, \text{pub}(k_{\text{example.com}})) \quad (19)$$

$$\text{enc}_s(s, k') \quad (20)$$

The term (19) shows an encrypted request (with  $r$  as in (17)). It is encrypted using the public key  $\text{pub}(k_{\text{example.com}})$ . The term (20) is a response (with  $s$  as in (18)). It is encrypted symmetrically using the (symmetric) key  $k'$  that was sent in the request (19).

*5) DNS Messages:*

*Definition 44.* A *DNS request* is a term of the form  $\langle \text{DNSResolve}, domain, nonce \rangle$  where  $domain \in \text{Doms}$ ,  $nonce \in \mathcal{N}$ . We call the set of all DNS requests DNSRequests.

*Definition 45.* A *DNS response* is a term of the form  $\langle \text{DNSResolved}, domain, result, nonce \rangle$  with  $domain \in \text{Doms}$ ,  $result \in \text{IPs}$ ,  $nonce \in \mathcal{N}$ . We call the set of all DNS responses DNSResponses.

DNS servers are supposed to include the nonce they received in a DNS request in the DNS response that they send back so that the party which issued the request can match it with the request.

### C. Atomic Processes, Systems and Runs

Entities that take part in a network are modeled as atomic processes. An atomic process takes a term that describes its current state and an event as input, and then (non-deterministically) outputs a new state and a sequence of events.

*Definition 46 (Generic Atomic Processes and Systems).* A (generic) atomic process is a tuple

$$p = (I^p, Z^p, R^p, s_0^p)$$

where  $I^p \subseteq \text{IPs}$ ,  $Z^p \subseteq \mathcal{T}_{\mathcal{N}}$  is a set of states,  $R^p \subseteq (\mathcal{E} \times Z^p) \times (2^{\mathcal{E}^\nu} \times \mathcal{T}_{\mathcal{N}}(V_{\text{process}}))$  (input event and old state map to sequence of output events and new state), and  $s_0^p \in Z^p$  is the initial state of  $p$ . For any new state  $s$  and any sequence of nonces  $(\eta_1, \eta_2, \dots)$  we demand that  $s[\eta_1/\nu_1, \eta_2/\nu_2, \dots] \in Z^p$ . A system  $\mathcal{P}$  is a (possibly infinite) set of atomic processes.

*Definition 47 (Configurations).* A configuration of a system  $\mathcal{P}$  is a tuple  $(S, E, N)$  where the state of the system  $S$  maps every atomic process  $p \in \mathcal{P}$  to its current state  $S(p) \in Z^p$ , the sequence of waiting events  $E$  is an infinite sequence<sup>9</sup>  $(e_1, e_2, \dots)$  of events waiting to be delivered, and  $N$  is an infinite sequence of nonces  $(n_1, n_2, \dots)$ .

*Definition 48 (Processing Steps).* A processing step of the system  $\mathcal{P}$  is of the form

$$(S, E, N) \xrightarrow[p \rightarrow E_{\text{out}}]{e_{\text{in}} \rightarrow p} (S', E', N')$$

where

- 1)  $(S, E, N)$  and  $(S', E', N')$  are configurations of  $\mathcal{P}$ ,
- 2)  $e_{\text{in}} = \langle a, f, m \rangle \in E$  is an event,
- 3)  $p \in \mathcal{P}$  is a process,
- 4)  $E_{\text{out}}$  is a sequence (term) of events

such that there exists

- 1) a sequence (term)  $E_{\text{out}}^\nu \subseteq 2^{\mathcal{E}^\nu}$  of protoevents,
- 2) a term  $s^\nu \in \mathcal{T}_{\mathcal{N}}(V_{\text{process}})$ ,
- 3) a sequence  $(v_1, v_2, \dots, v_i)$  of all placeholders appearing in  $E_{\text{out}}^\nu$  (ordered lexicographically),
- 4) a sequence  $N^\nu = (\eta_1, \eta_2, \dots, \eta_i)$  of the first  $i$  elements in  $N$

with

- 1)  $((e_{\text{in}}, S(p)), (E_{\text{out}}^\nu, s^\nu)) \in R^p$  and  $a \in I^p$ ,
- 2)  $E_{\text{out}} = E_{\text{out}}^\nu[\eta_1/v_1, \dots, \eta_i/v_i]$ ,
- 3)  $S'(p) = s^\nu[\eta_1/v_1, \dots, \eta_i/v_i]$  and  $S'(p') = S(p')$  for all  $p' \neq p$ ,
- 4)  $E' = E_{\text{out}} \cdot (E \setminus \{e_{\text{in}}\})$ ,
- 5)  $N' = N \setminus N^\nu$ .

We may omit the superscript and/or subscript of the arrow.

Intuitively, for a processing step, we select one of the processes in  $\mathcal{P}$ , and call it with one of the events in the list of waiting events  $E$ . In its output (new state and output events), we replace any occurrences of placeholders  $\nu_x$  by “fresh” nonces from  $N$  (which we then remove from  $N$ ). The output events are then prepended to the list of waiting events, and the state of the process is reflected in the new configuration.

*Definition 49 (Runs).* Let  $\mathcal{P}$  be a system,  $E^0$  be sequence of events, and  $N^0$  be a sequence of nonces. A run  $\rho$  of a system  $\mathcal{P}$  initiated by  $E^0$  with nonces  $N^0$  is a finite sequence of configurations  $((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  or an infinite sequence of configurations  $((S^0, E^0, N^0), \dots)$  such that  $S^0(p) = s_0^p$  for all  $p \in \mathcal{P}$  and  $(S^i, E^i, N^i) \rightarrow (S^{i+1}, E^{i+1}, N^{i+1})$  for all  $0 \leq i < n$  (finite run) or for all  $i \geq 0$  (infinite run).

We denote the state  $S^n(p)$  of a process  $p$  at the end of a finite run  $\rho$  by  $\rho(p)$ .

Usually, we will initiate runs with a set  $E^0$  containing infinite trigger events of the form  $\langle a, a, \text{TRIGGER} \rangle$  for each  $a \in \text{IPs}$ , interleaved by address.

#### D. Atomic Dolev-Yao Processes

We next define atomic Dolev-Yao processes, for which we require that the messages and states that they output can be computed (more formally, derived) from the current input event and state. For this purpose, we first define what it means to derive a message from given messages.

*Definition 50 (Deriving Terms).* Let  $M$  be a set of ground terms. We say that a term  $m$  can be derived from  $M$  with placeholders  $V$  if there exist  $n \geq 0$ ,  $m_1, \dots, m_n \in M$ , and  $\tau \in \mathcal{T}_\emptyset(\{x_1, \dots, x_n\} \cup V)$  such that  $m \equiv \tau[m_1/x_1, \dots, m_n/x_n]$ . We denote by  $d_V(M)$  the set of all messages that can be derived from  $M$  with variables  $V$ .

For example, the term  $a$  can be derived from the set of terms  $\{\text{enc}_a(\langle a, b, c \rangle, \text{pub}(k)), k\}$ , i.e.,  $a \in d_\emptyset(\{\text{enc}_a(\langle a, b, c \rangle, \text{pub}(k)), k\})$ .

<sup>9</sup>Here: Not in the sense of terms as defined earlier.

A (Dolev-Yao) process consists of a set of addresses the process listens to, a set of states (terms), an initial state, and a relation that takes an event and a state as input and (non-deterministically) returns a new state and a sequence of events. The relation models a computation step of the process. It is required that the output can be derived from the input event and the state.

*Definition 51 (Atomic Dolev-Yao Process).* An atomic Dolev-Yao process (or simply, a DY process) is a tuple  $p = (I^p, Z^p, R^p, s_0^p)$  such that  $p$  is an atomic process and for all events  $e \in \mathcal{E}$ , sequences of protoevents  $E, s \in \mathcal{T}_{\mathcal{N}}$ ,  $s' \in \mathcal{T}_{\mathcal{N}}(V_{\text{process}})$ , with  $((e, s), (E, s')) \in R^p$  it holds true that  $E, s' \in d_{V_{\text{process}}}(\{e, s\})$ .

### E. Attackers

The so-called *attacker process* is a Dolev-Yao process which records all messages it receives and outputs any finite sequence of events it can possibly derive from its recorded messages. Hence, an attacker process carries out all attacks any Dolev-Yao process could possibly perform. Attackers can corrupt other parties (using corrupt messages).

*Definition 52 (Atomic Attacker Process).* An (atomic) attacker process for a set of sender addresses  $A \subseteq \text{IPs}$  is an atomic DY process  $p = (I, Z, R, s_0)$  such that for all events  $e$ , and  $s \in \mathcal{T}_{\mathcal{N}}$  we have that  $((e, s), (E, s')) \in R$  iff  $s' = \langle e, E, s \rangle$  and  $E = \langle \langle a_1, f_1, m_1 \rangle, \dots, \langle a_n, f_n, m_n \rangle \rangle$  with  $n \in \mathbb{N}$ ,  $a_1, \dots, a_n \in \text{IPs}$ ,  $f_1, \dots, f_n \in A$ ,  $m_1, \dots, m_n \in d_{V_{\text{process}}}(\{e, s\})$ .

Note that in a web system, we distinguish between two kinds of attacker processes: web attackers and network attackers. Both kinds match the definition above, but differ in the set of assigned addresses in the context of a web system. While for web attackers, the set of addresses  $I^p$  is disjoint from other web attackers and honest processes, i.e., web attackers participate in the network as any other party, the set of addresses  $I^p$  of a network attacker is not restricted. Hence, a network attacker can intercept events addressed to any party as well as spoof all addresses. Note that one network attacker subsumes any number of web attackers as well as any number of network attackers.

### F. Notations for Functions and Algorithms

When describing algorithms, we use the following notations:

1) *Non-deterministic choosing and iteration:* The notation **let**  $n \leftarrow N$  is used to describe that  $n$  is chosen non-deterministically from the set  $N$ . If  $N$  is empty, the corresponding processing step in which this selection happens does not finish. We write **for**  $s \in M$  **do** to denote that the following commands are repeated for every element in  $M$ , where the variable  $s$  is the current element. The order in which the elements are processed is chosen non-deterministically. We write, for example,

**let**  $x, y$  **such that**  $\langle \text{Constant}, x, y \rangle \equiv t$  **if possible; otherwise** doSomethingElse

for some variables  $x, y$ , a string **Constant**, and some term  $t$  to express that  $x := \pi_2(t)$ , and  $y := \pi_3(t)$  if  $\text{Constant} \equiv \pi_1(t)$  and if  $|\langle \text{Constant}, x, y \rangle| = |t|$ , and that otherwise  $x$  and  $y$  are not set and doSomethingElse is executed.

2) *Function calls:* When calling functions that do not return anything, we write

**call** FUNCTION\_NAME( $x, y$ )

to describe that a function FUNCTION\_NAME is called with two variables  $x$  and  $y$  as parameters. If that function executes the command **stop**  $E, s'$ , the processing step terminates, where  $E$  is the sequence of events output by the associated process and  $s'$  is its new state. If that function does not terminate with a **stop**, the control flow returns to the calling function at the next line after the call.

When calling a function that has a return value, we omit the **call** and directly write

**let**  $z := \text{FUNCTION\_NAME}(x, y)$

to assign the return value to a variable  $z$  after the function returns. Note that the semantics for execution of **stop** within such functions is the same as for functions without a return value.

3) *Stop without output:* We write **stop** (without further parameters) to denote that there is no output and no change in the state.

4) *Placeholders:* In several places throughout the algorithms we use placeholders to generate “fresh” nonces as described in our communication model (see Definition 27). Table I shows a list of some of the placeholders, generally denoted by  $\nu$  with some subscript to distinguish between multiple fresh values.

5) *Abbreviations for URLs and Origins:* We sometimes use an abbreviation for URLs. We write  $\text{URL}_{path}^d$  to describe the following URL term:  $\langle \text{URL}, S, d, path, \langle \rangle \rangle$ . If the domain  $d$  belongs to some distinguished process  $P$  and it is the only domain associated to this process, we may also write  $\text{URL}_{path}^P$ . For a (secure) origin  $\langle d, S \rangle$  of some domain  $d$ , we also write  $\text{origin}_d$ . Again, if the domain  $d$  belongs to some distinguished process  $P$  and  $d$  is the only domain associated to this process, we may write  $\text{origin}_P$ .

### G. Browsers

Here, we present the formal model of browsers.

Placeholder	Usage
$\nu_1$	Algorithm 22, new window nonces
$\nu_2$	Algorithm 22, new HTTP request nonce
$\nu_3$	Algorithm 22, lookup key for pending HTTP requests entry
$\nu_4$	Algorithm 20, new HTTP request nonce (multiple lines)
$\nu_5$	Algorithm 20, new subwindow nonce
$\nu_6$	Algorithm 21, new HTTP request nonce
$\nu_7$	Algorithm 21, new document nonce
$\nu_8$	Algorithm 17, lookup key for pending DNS entry
$\nu_9$	Algorithm 14, new window nonce
$\nu_{10}, \dots$	Algorithm 20, replacement for placeholders in script output

Table I: List of placeholders used in browser algorithms.

1) *Scripts*: Recall that a *script* models JavaScript running in a browser. Scripts are defined similarly to Dolev-Yao processes. When triggered by a browser, a script is provided with state information. The script then outputs a term representing a new internal state and a command to be interpreted by the browser (see also the specification of browsers below).

*Definition 53 (Placeholders for Scripts)*. By  $V_{\text{script}} = \{\lambda_1, \dots\}$  we denote an infinite set of variables used in scripts.

*Definition 54 (Scripts)*. A *script* is a relation  $R \subseteq \mathcal{T}_{\mathcal{N}} \times \mathcal{T}_{\mathcal{N}}(V_{\text{script}})$  such that for all  $s \in \mathcal{T}_{\mathcal{N}}$ ,  $s' \in \mathcal{T}_{\mathcal{N}}(V_{\text{script}})$  with  $(s, s') \in R$  it follows that  $s' \in d_{V_{\text{script}}}(s)$ .

A script is called by the browser which provides it with state information (such as the script's last scriptstate and limited information about the browser's state)  $s$ . The script then outputs a term  $s'$ , which represents the new scriptstate and some command which is interpreted by the browser. The term  $s'$  may contain variables  $\lambda_1, \dots$  which the browser will replace by (otherwise unused) placeholders  $\nu_1, \dots$  which will be replaced by nonces once the browser DY process finishes (effectively providing the script with a way to get "fresh" nonces).

Similarly to an attacker process, the so-called *attacker script* outputs everything that is derivable from the input.

*Definition 55 (Attacker Script)*. The attacker script  $R^{\text{att}}$  outputs everything that is derivable from the input, i.e.,  $R^{\text{att}} = \{(s, s') \mid s \in \mathcal{T}_{\mathcal{N}}, s' \in d_{V_{\text{script}}}(s)\}$ .

2) *Web Browser State*: Before we can define the state of a web browser, we first have to define windows and documents.

*Definition 56*. A *window* is a term of the form  $w = \langle \text{nonce}, \text{documents}, \text{opener} \rangle$  with  $\text{nonce} \in \mathcal{N}$ ,  $\text{documents} \subset^{\langle \rangle} \text{Documents}$  (defined below),  $\text{opener} \in \mathcal{N} \cup \{\perp\}$  where  $d.\text{active} = \top$  for exactly one  $d \in^{\langle \rangle} \text{documents}$  if  $\text{documents}$  is not empty (we then call  $d$  the *active document of*  $w$ ). We write  $\text{Windows}$  for the set of all windows. We write  $w.\text{activedocument}$  to denote the active document inside window  $w$  if it exists and  $\langle \rangle$  else.

We will refer to the window nonce as (*window*) *reference*.

The documents contained in a window term to the left of the active document are the previously viewed documents (available to the user via the "back" button) and the documents in the window term to the right of the currently active document are documents available via the "forward" button.

A window  $a$  may have opened a top-level window  $b$  (i.e., a window term which is not a subterm of a document term). In this case, the *opener* part of the term  $b$  is the nonce of  $a$ , i.e.,  $b.\text{opener} = a.\text{nonce}$ .

*Definition 57*. A *document*  $d$  is a term of the form

$$\langle \text{nonce}, \text{location}, \text{headers}, \text{referrer}, \text{script}, \text{scriptstate}, \text{scriptinputs}, \text{subwindows}, \text{active} \rangle$$

where  $\text{nonce} \in \mathcal{N}$ ,  $\text{location} \in \text{URLs}$ ,  $\text{headers} \in [\mathbb{S} \times \mathcal{T}_{\mathcal{N}}]$ ,  $\text{referrer} \in \text{URLs} \cup \{\perp\}$ ,  $\text{script} \in \mathcal{T}_{\mathcal{N}}$ ,  $\text{scriptstate} \in \mathcal{T}_{\mathcal{N}}$ ,  $\text{scriptinputs} \in \mathcal{T}_{\mathcal{N}}$ ,  $\text{subwindows} \subset^{\langle \rangle} \text{Windows}$ ,  $\text{active} \in \{\top, \perp\}$ . A *limited document* is a term of the form  $\langle \text{nonce}, \text{subwindows} \rangle$  with  $\text{nonce}, \text{subwindows}$  as above. A window  $w \in^{\langle \rangle} \text{subwindows}$  is called a *subwindow* (of  $d$ ). We write  $\text{Documents}$  for the set of all documents. For a document term  $d$  we write  $d.\text{origin}$  to denote the origin of the document, i.e., the term  $\langle d.\text{location}.\text{host}, d.\text{location}.\text{protocol} \rangle \in \text{Origins}$ .

We will refer to the document nonce as (*document*) *reference*.

*Definition 58*. For two window terms  $w$  and  $w'$  we write

$$w \xrightarrow{\text{childof}} w'$$

if  $w \in^{\langle \rangle} w'.\text{activedocument}.\text{subwindows}$ . We write  $\xrightarrow{\text{childof}^+}$  for the transitive closure and we write  $\xrightarrow{\text{childof}^*}$  for the reflexive transitive closure.

In the web browser state, HTTP(S) messages are tracked using *references*, where we distinguish between references for XMLHttpRequests and references for normal HTTP(S) requests.

*Definition 59.* A reference for a normal HTTP(S) request is a sequence of the form  $\langle \text{REQ}, \textit{nonce} \rangle$ , where *nonce* is a window reference. A reference for a XMLHttpRequest is a sequence of the form  $\langle \text{XHR}, \textit{nonce}, \textit{xhrreference} \rangle$ , where *nonce* is a document reference and *xhrreference* is some nonce that was chosen by the script that initiated the request.

We can now define the set of states of web browsers. Note that we use the dictionary notation that we introduced in Definition 34.

*Definition 60.* The set of states  $Z_{\text{webbrowser}}$  of a web browser atomic Dolev-Yao process consists of the terms of the form

$$\langle \textit{windows}, \textit{ids}, \textit{secrets}, \textit{cookies}, \textit{localStorage}, \textit{sessionStorage}, \textit{keyMapping}, \textit{sts}, \textit{DNSaddress}, \textit{pendingDNS}, \textit{pendingRequests}, \textit{isCorrupted}, \textit{cibaBindingMessages}, \textit{tlskeys} \rangle$$

with the subterms as follows:

- $\textit{windows} \subset \langle \rangle$  Windows contains a list of window terms (modeling top-level windows, or browser tabs) which contain documents, which in turn can contain further window terms (iframes).
- $\textit{ids} \subset \langle \rangle \mathcal{T}_{\mathcal{N}}$  is a list of identities that are owned by this browser (i.e., belong to the user of the browser).
- $\textit{secrets} \in [\text{Origins} \times \mathcal{T}_{\mathcal{N}}]$  contains a list of secrets that are associated with certain origins (i.e., passwords of the user of the browser at certain websites). Note that this structure allows to have a single secret under an origin or a list of secrets under an origin.
- $\textit{cookies}$  is a dictionary over Doms and sequences of Cookies modeling cookies that are stored for specific domains.
- $\textit{localStorage} \in [\text{Origins} \times \mathcal{T}_{\mathcal{N}}]$  stores the data saved by scripts using the localStorage API (separated by origins).
- $\textit{sessionStorage} \in [OR \times \mathcal{T}_{\mathcal{N}}]$  for  $OR := \{\langle o, r \rangle \mid o \in \text{Origins}, r \in \mathcal{N}\}$  similar to localStorage, but the data in sessionStorage is additionally separated by top-level windows.
- $\textit{keyMapping} \in [\text{Doms} \times \mathcal{T}_{\mathcal{N}}]$  maps domains to TLS encryption keys.
- $\textit{sts} \subset \langle \rangle$  Doms stores the list of domains that the browser only accesses via TLS (strict transport security).
- $\textit{DNSaddress} \in \text{IPs}$  defines the IP address of the DNS server.
- $\textit{pendingDNS} \in [\mathcal{N} \times \mathcal{T}_{\mathcal{N}}]$  contains one pairing per unanswered DNS query of the form  $\langle \textit{reference}, \textit{request}, \textit{url} \rangle$ . In these pairings, *reference* is an HTTP(S) request reference (as above), *request* contains the HTTP(S) message that awaits DNS resolution, and *url* contains the URL of said HTTP request. The pairings in *pendingDNS* are indexed by the DNS request/response nonce.
- $\textit{pendingRequests} \in \mathcal{T}_{\mathcal{N}}$  contains pairings of the form  $\langle \textit{reference}, \textit{request}, \textit{url}, \textit{key}, \textit{f} \rangle$  with *reference*, *request*, and *url* as in *pendingDNS*, *key* is the symmetric encryption key if HTTPS is used or  $\perp$  otherwise, and *f* is the IP address of the server to which the request was sent.
- $\textit{isCorrupted} \in \{\perp, \text{FULLCORRUPT}, \text{CLOSECORRUPT}\}$  specifies the corruption level of the browser.
- $\textit{cibaBindingMessages} \in \mathcal{T}_{\mathcal{N}}$  contains pairings of the form  $\langle \textit{dom}, \textit{bindingMsg} \rangle$ , where *bindingMsg* is a CIBA binding message received from the (client) domain *dom*. The browser compares this binding message to the value received from an AS.
- $\textit{tlskeys} \in [\text{Doms} \times \mathcal{N}]$  is a mapping from domains to private keys.

In corrupted browsers, certain subterms are used in different ways (e.g., *pendingRequests* is used to store all observed messages).

3) *Web Browser Relation:* We will now define the relation  $R_{\text{webbrowser}}$  of a standard HTTP browser. We first introduce some notations and then describe the functions that are used for defining the browser main algorithm. We then define the browser relation.

a) *Helper Functions:* In the following description of the web browser relation  $R_{\text{webbrowser}}$  we use the helper functions Subwindows, Docs, Clean, CookieMerge, AddCookie, and NavigableWindows.

**Subwindows and Docs.** Given a browser state  $s$ , Subwindows( $s$ ) denotes the set of all pointers<sup>10</sup> to windows in the window list  $s.\textit{windows}$  and (recursively) the subwindows of their active documents. We exclude subwindows of inactive documents and their subwindows. With Docs( $s$ ) we denote the set of pointers to all active documents in the set of windows referenced by Subwindows( $s$ ).

*Definition 61.* For a browser state  $s$  we denote by Subwindows( $s$ ) the minimal set of pointers that satisfies the following conditions: (1) For all windows  $w \in \langle \rangle s.\textit{windows}$  there is a  $\bar{p} \in \text{Subwindows}(s)$  such that  $s.\bar{p} = w$ . (2) For all  $\bar{p} \in \text{Subwindows}(s)$ , the active document  $d$  of the window  $s.\bar{p}$  and every subwindow  $w$  of  $d$  there is a pointer  $\bar{p}' \in \text{Subwindows}(s)$  such that  $s.\bar{p}' = w$ .

Given a browser state  $s$ , the set Docs( $s$ ) of pointers to active documents is the minimal set such that for every  $\bar{p} \in \text{Subwindows}(s)$  with  $s.\bar{p}.\textit{activedocument} \neq \langle \rangle$ , there exists a pointer  $\bar{p}' \in \text{Docs}(s)$  with  $s.\bar{p}' = s.\bar{p}.\textit{activedocument}$ .

By Subwindows<sup>+</sup>( $s$ ) and Docs<sup>+</sup>( $s$ ) we denote the respective sets that also include the inactive documents and their subwindows.

<sup>10</sup>Recall the definition of a pointer in Definition 36.

**Clean.** The function Clean will be used to determine which information about windows and documents the script running in the document  $d$  has access to.

*Definition 62.* Let  $s$  be a browser state and  $d$  a document. By  $\text{Clean}(s, d)$  we denote the term that equals  $s.\text{windows}$  but with (1) all inactive documents removed (including their subwindows etc.), (2) all subterms that represent non-same-origin documents w.r.t.  $d$  replaced by a limited document  $d'$  with the same nonce and the same subwindow list, and (3) the values of the subterms `headers` for all documents set to  $\langle \rangle$ . (Note that non-same-origin documents on all levels are replaced by their corresponding limited document.)

**CookieMerge.** The function CookieMerge merges two sequences of cookies together: When used in the browser, *oldcookies* is the sequence of existing cookies for some origin, *newcookies* is a sequence of new cookies that was output by some script. The sequences are merged into a set of cookies using an algorithm that is based on the *Storage Mechanism* algorithm described in RFC6265.

*Definition 63.* For a sequence of cookies (with pairwise different names) *oldcookies*, a sequence of cookies *newcookies*, and a string  $\text{protocol} \in \{\text{P}, \text{S}\}$ , the set  $\text{CookieMerge}(\text{oldcookies}, \text{newcookies}, \text{protocol})$  is defined by the following algorithm: From *newcookies* remove all cookies  $c$  that have  $c.\text{content.httpOnly} \equiv \top$  or where  $(c.\text{name}.1 \equiv \_\text{Host}) \wedge ((\text{protocol} \equiv \text{P}) \vee (c.\text{secure} \equiv \perp))$ . For any  $c, c' \in \langle \rangle \text{newcookies}$ ,  $c.\text{name} \equiv c'.\text{name}$ , remove the cookie that appears left of the other in *newcookies*. Let  $m$  be the set of cookies that have a name that either appears in *oldcookies* or in *newcookies*, but not in both. For all pairs of cookies  $(c_{\text{old}}, c_{\text{new}})$  with  $c_{\text{old}} \in \langle \rangle \text{oldcookies}$ ,  $c_{\text{new}} \in \langle \rangle \text{newcookies}$ ,  $c_{\text{old}}.\text{name} \equiv c_{\text{new}}.\text{name}$ , add  $c_{\text{new}}$  to  $m$  if  $c_{\text{old}}.\text{content.httpOnly} \equiv \perp$  and add  $c_{\text{old}}$  to  $m$  otherwise. The result of  $\text{CookieMerge}(\text{oldcookies}, \text{newcookies}, \text{protocol})$  is  $m$ .

**AddCookie.** The function AddCookie adds a cookie  $c$  received in an HTTP response to the sequence of cookies contained in the sequence *oldcookies*. It is again based on the algorithm described in RFC6265 but simplified for the use in the browser model.

*Definition 64.* For a sequence of cookies (with pairwise different names) *oldcookies*, a cookie  $c$ , and a string  $\text{protocol} \in \{\text{P}, \text{S}\}$  (denoting whether the HTTP response was received from an insecure or a secure origin), the sequence  $\text{AddCookie}(\text{oldcookies}, c, \text{protocol})$  is defined by the following algorithm: Let  $m := \text{oldcookies}$ . If  $(c.\text{name}.1 \equiv \_\text{Host}) \wedge \neg((\text{protocol} \equiv \text{S}) \wedge (c.\text{secure} \equiv \top))$ , then return  $m$ , else: Remove any  $c'$  from  $m$  that has  $c.\text{name} \equiv c'.\text{name}$ . Append  $c$  to  $m$  and return  $m$ .

**NavigableWindows.** The function NavigableWindows returns a set of windows that a document is allowed to navigate. We closely follow [7], Section 5.1.4 for this definition.

*Definition 65.* The set  $\text{NavigableWindows}(\bar{w}, s')$  is the set  $\bar{W} \subseteq \text{Subwindows}(s')$  of pointers to windows that the active document in  $\bar{w}$  is allowed to navigate. The set  $\bar{W}$  is defined to be the minimal set such that for every  $\bar{w}' \in \text{Subwindows}(s')$  the following is true:

- If  $s'.\bar{w}'.\text{activedocument.origin} \equiv s'.\bar{w}.\text{activedocument.origin}$  (i.e., the active documents in  $\bar{w}$  and  $\bar{w}'$  are same-origin), then  $\bar{w}' \in \bar{W}$ , and
- If  $s'.\bar{w} \xrightarrow{\text{childof}^*} s'.\bar{w}' \wedge \nexists \bar{w}'' \in \text{Subwindows}(s') \text{ with } s'.\bar{w}' \xrightarrow{\text{childof}^*} s'.\bar{w}''$  ( $\bar{w}'$  is a top-level window and  $\bar{w}$  is an ancestor window of  $\bar{w}'$ ), then  $\bar{w}' \in \bar{W}$ , and
- If  $\exists \bar{p} \in \text{Subwindows}(s')$  such that  $s'.\bar{w}' \xrightarrow{\text{childof}^+} s'.\bar{p}$   
 $\wedge s'.\bar{p}.\text{activedocument.origin} = s'.\bar{w}.\text{activedocument.origin}$  ( $\bar{w}'$  is not a top-level window but there is an ancestor window  $\bar{p}$  of  $\bar{w}'$  with an active document that has the same origin as the active document in  $\bar{w}$ ), then  $\bar{w}' \in \bar{W}$ , and
- If  $\exists \bar{p} \in \text{Subwindows}(s')$  such that  $s'.\bar{w}'.\text{opener} = s'.\bar{p}.\text{nonce} \wedge \bar{p} \in \bar{W}$  ( $\bar{w}'$  is a top-level window—it has an opener—and  $\bar{w}$  is allowed to navigate the opener window of  $\bar{w}'$ ,  $\bar{p}$ ), then  $\bar{w}' \in \bar{W}$ .

b) *Functions:*

- The function GETNAVIGABLEWINDOW (Algorithm 14) is called by the browser to determine the window that is *actually* navigated when a script in the window  $s'.\bar{w}$  provides a window reference for navigation (e.g., for opening a link). When it is given a window reference (nonce) *window*, this function returns a pointer to a selected window term in  $s'$ :
  - If *window* is the string `_BLANK`, a new window is created and a pointer to that window is returned.
  - If *window* is a nonce (reference) and there is a window term with a reference of that value in the windows in  $s'$ , a pointer  $\bar{w}'$  to that window term is returned, as long as the window is navigable by the current window's document (as defined by NavigableWindows above).

In all other cases,  $\bar{w}$  is returned instead (the script navigates its own window).

- The function GETWINDOW (Algorithm 15) takes a window reference as input and returns a pointer to a window as above, but it checks only that the active documents in both windows are same-origin. It creates no new windows.

---

**Algorithm 14** Web Browser Model: Determine window for navigation.

```
1: function GETNAVIGABLEWINDOW( $\bar{w}$ ,  $window$ ,  $noreferrer$ ,  $s'$ )
2:   if  $window \equiv \_BLANK$  then  $\rightarrow$  Open a new window when  $\_BLANK$  is used
3:   if  $noreferrer \equiv \top$  then
4:     let  $w' := \langle \nu_9, \langle \rangle, \perp \rangle$ 
5:   else
6:     let  $w' := \langle \nu_9, \langle \rangle, s'.\bar{w}.nonce \rangle$ 
7:     let  $s'.windows := s'.windows + \langle \rangle w'$ 
       $\hookrightarrow$  and let  $\bar{w}'$  be a pointer to this new element in  $s'$ 
8:     return  $w'$ 
9:   let  $w' \leftarrow \text{NavigableWindows}(\bar{w}, s')$  such that  $s'.\bar{w}'.nonce \equiv window$ 
       $\hookrightarrow$  if possible; otherwise return  $\bar{w}$ 
10:  return  $w'$ 
```

---

---

**Algorithm 15** Web Browser Model: Determine same-origin window.

```
1: function GETWINDOW( $\bar{w}$ ,  $window$ ,  $s'$ )
2:   let  $w' \leftarrow \text{Subwindows}(s')$  such that  $s'.\bar{w}'.nonce \equiv window$ 
       $\hookrightarrow$  if possible; otherwise return  $\bar{w}$ 
3:   if  $s'.\bar{w}'.activatedocument.origin \equiv s'.\bar{w}.activatedocument.origin$  then
4:     return  $w'$ 
5:   return  $\bar{w}$ 
```

---

---

**Algorithm 16** Web Browser Model: Cancel pending requests for given window.

```
1: function CANCELNAV( $reference$ ,  $s'$ )
2:   remove all  $\langle reference, req, url, key, f \rangle$  from  $s'.pendingRequests$  for any  $req, url, key, f$ 
3:   remove all  $\langle x, \langle reference, message, url \rangle \rangle$  from  $s'.pendingDNS$ 
       $\hookrightarrow$  for any  $x, message, url$ 
4:   return  $s'$ 
```

---

---

**Algorithm 17** Web Browser Model: Prepare headers, do DNS resolution, save message.

```
1: function HTTP_SEND( $reference$ ,  $message$ ,  $url$ ,  $origin$ ,  $referrer$ ,  $referrerPolicy$ ,  $a$ ,  $s'$ )
2:   if  $message.host \in \langle \rangle s'.sts$  then
3:     let  $url.protocol := S$ 
4:     let  $cookies := \{ \{ \langle c.name, c.content.value \rangle \mid c \in \langle \rangle s'.cookies[message.host] \}$ 
       $\hookrightarrow \wedge (c.content.secure \equiv \top \implies (url.protocol \equiv S)) \}$ 
5:     let  $message.headers[Cookie] := cookies$ 
6:     if  $origin \neq \perp$  then
7:       let  $message.headers[Origin] := origin$ 
8:     if  $referrerPolicy \equiv \text{no-referrer}$  then
9:       let  $referrer := \perp$ 
10:    if  $referrer \neq \perp$  then
11:      if  $referrerPolicy \equiv \text{origin}$  then
12:        let  $referrer := \langle \text{URL}, referrer.protocol, referrer.host, /, \langle \rangle, \perp \rangle$ 
           $\rightarrow$  Referrer stripped down to origin.
13:        let  $referrer.fragment := \perp$ 
           $\rightarrow$  Browsers do not send fragment identifiers in the Referer header.
14:        let  $message.headers[Referer] := referrer$ 
15:    let  $s'.pendingDNS[\nu_8] := \langle reference, message, url \rangle$ 
16:    stop  $\langle \langle s'.DNSaddress, a, \langle \text{DNSResolve}, message.host, \nu_8 \rangle \rangle, s' \rangle$ 
```

---

---

**Algorithm 18** Web Browser Model: Navigate a window backward.

```
1: function NAVBACK( $\bar{w}'$ ,  $s'$ )
2:   if  $\exists \bar{j} \in \mathbb{N}, \bar{j} > 1$  such that  $s'.\bar{w}'.documents.\bar{j}.active \equiv \top$  then
3:     let  $s'.\bar{w}'.documents.\bar{j}.active := \perp$ 
4:     let  $s'.\bar{w}'.documents.(\bar{j} - 1).active := \top$ 
5:     let  $s' := \text{CANCELNAV}(s'.\bar{w}'.nonce, s')$ 
6:   stop  $\langle \rangle, s'$ 
```

---



---

**Algorithm 19** Web Browser Model: Navigate a window forward.

---

```
1: function NAVFORWARD( $\bar{w}'$ ,  $s'$ )
2:   if  $\exists \bar{j} \in \mathbb{N}$  such that  $s'.\bar{w}'.documents.\bar{j}.active \equiv \top$ 
    $\hookrightarrow \wedge s'.\bar{w}'.documents.(\bar{j} + 1) \in \text{Documents}$  then
3:     let  $s'.\bar{w}'.documents.\bar{j}.active := \perp$ 
4:     let  $s'.\bar{w}'.documents.(\bar{j} + 1).active := \top$ 
5:     let  $s' := \text{CANCELNAV}(s'.\bar{w}'.nonce, s')$ 
6:   stop  $\langle \rangle$ ,  $s'$ 
```

---

---

**Algorithm 20** Web Browser Model: Execute a script.

---

```
1: function RUNSCRIPT( $\bar{w}$ ,  $\bar{d}$ ,  $a$ ,  $s'$ )
2:   let  $tree := \text{Clean}(s', s'.\bar{d})$ 
3:   let  $cookies := \langle \{ \langle c.name, c.content.value \rangle \mid c \in \langle \rangle s'.cookies [s'.\bar{d}.origin.host] \}$ 
    $\hookrightarrow \wedge c.content.httpOnly \equiv \perp$ 
    $\hookrightarrow \wedge (c.content.secure \equiv \top \implies (s'.\bar{d}.origin.protocol \equiv S)) \rangle \rangle$ 
4:   let  $tlw \leftarrow s'.windows$  such that  $tlw$  is the top-level window containing  $\bar{d}$ 
5:   let  $sessionStorage := s'.sessionStorage [s'.\bar{d}.origin, tlw.nonce]$ 
6:   let  $localStorage := s'.localStorage [s'.\bar{d}.origin]$ 
7:   let  $secrets := s'.secrets [s'.\bar{d}.origin]$ 
8:   let  $R := \text{script}^{-1}(s'.\bar{d}.script)$  if possible; otherwise stop
9:   let  $in := \langle tree, s'.\bar{d}.nonce, s'.\bar{d}.scriptstate, s'.\bar{d}.scriptinputs, cookies,$ 
    $\hookrightarrow localStorage, sessionStorage, s'.ids, secrets \rangle$ 
10:  let  $state' \leftarrow \mathcal{T}_{\mathcal{N}}(V_{\text{process}})$ ,  $cookies' \leftarrow \text{Cookies}^\nu$ ,  $localStorage' \leftarrow \mathcal{T}_{\mathcal{N}}(V_{\text{process}})$ ,
    $\hookrightarrow sessionStorage' \leftarrow \mathcal{T}_{\mathcal{N}}(V_{\text{process}})$ ,  $command \leftarrow \mathcal{T}_{\mathcal{N}}(V_{\text{process}})$ ,
    $\hookrightarrow out := \langle state', cookies', localStorage', sessionStorage', command \rangle$ 
    $\hookrightarrow$  such that  $out := out^\lambda[\nu_{10}/\lambda_1, \nu_{11}/\lambda_2, \dots]$  with  $(in, out^\lambda) \in R$ 
11:  let  $s'.cookies [s'.\bar{d}.origin.host] :=$ 
    $\hookrightarrow \langle \text{CookieMerge}(s'.cookies [s'.\bar{d}.origin.host], cookies', s'.\bar{d}.origin.protocol) \rangle$ 
12:  let  $s'.localStorage [s'.\bar{d}.origin] := localStorage'$ 
13:  let  $s'.sessionStorage [s'.\bar{d}.origin, tlw.nonce] := sessionStorage'$ 
14:  let  $s'.\bar{d}.scriptstate := state'$ 
15:  let  $referrer := s'.\bar{d}.location$ 
16:  let  $referrerPolicy := s'.\bar{d}.headers[ReferrerPolicy]$ 
17:  let  $docorigin := s'.\bar{d}.origin$ 
18:  switch  $command$  do
19:    case  $\langle \text{HREF}, url, hrefwindow, noreferrer \rangle$ 
20:      let  $\bar{w}' := \text{GETNAVIGABLEWINDOW}(\bar{w}, hrefwindow, noreferrer, s')$ 
21:      let  $reference := \langle \text{REQ}, s'.\bar{w}'.nonce \rangle$ 
22:      let  $req := \langle \text{HTTPReq}, \nu_4, \text{GET}, url.host, url.path, url.parameters, \langle \rangle, \langle \rangle \rangle$ 
23:      if  $noreferrer \equiv \top$  then
24:        let  $referrerPolicy := noreferrer$ 
25:        let  $s' := \text{CANCELNAV}(reference, s')$ 
26:        call  $\text{HTTP\_SEND}(reference, req, url, \perp, referrer, referrerPolicy, a, s')$ 
27:      case  $\langle \text{IFRAME}, url, window \rangle$ 
28:        if  $window \equiv \_SELF$  then
29:          let  $\bar{w}' := \bar{w}$ 
30:        else
31:          let  $\bar{w}' := \text{GETWINDOW}(\bar{w}, window, s')$ 
32:          let  $req := \langle \text{HTTPReq}, \nu_4, \text{GET}, url.host, url.path, url.parameters, \langle \rangle, \langle \rangle \rangle$ 
33:          let  $w' := \langle \nu_5, \langle \rangle, \perp \rangle$ 
34:          let  $s'.\bar{w}'.activedocument.subwindows := s'.\bar{w}'.activedocument.subwindows + \langle \rangle w'$ 
35:          call  $\text{HTTP\_SEND}(\langle \text{REQ}, \nu_5 \rangle, req, url, \perp, referrer, referrerPolicy, a, s')$ 
```

$\rightarrow$  Algorithm continues on next page.

---

---

```

36:  case ⟨FORM, url, method, data, hrefwindow⟩
37:  if method ∉ {GET, POST} then
38:    stop
39:    let  $\overline{w'}$  := GETNAVIGABLEWINDOW( $\overline{w}$ , hrefwindow,  $\perp$ ,  $s'$ )
40:    let reference := ⟨REQ,  $s'.\overline{w'}$ .nonce⟩
41:    if method = GET then
42:      let body := ⟨⟩
43:      let parameters := data
44:      let origin :=  $\perp$ 
45:    else
46:      let body := data
47:      let parameters := url.parameters
48:      let origin := docorigin
49:    let req := ⟨HTTPReq,  $\nu_4$ , method, url.host, url.path, parameters, ⟨⟩, body⟩
50:    let  $s'$  := CANCELNAV(reference,  $s'$ )
51:    call HTTP_SEND(reference, req, url, origin, referrer, referrerPolicy, a,  $s'$ )
52:  case ⟨CIBAFORM, url, method, data, hrefwindow, clientDomain, cibaBindingMessage⟩
    → Custom CIBA FORM command: When starting a CIBA flow, the client returns a binding message. When authenticating at the
    AS, the end-user has to make sure that they receive the same value. For modeling this behavior, we extend the browser state
    by the cibaBindingMessages subterm and define this command which first checks if the cibaBindingMessage is stored
    by the browser and then continues as the FORM command. Note that this command is a modeling artifact.
53:  if ⟨clientDomain, cibaBindingMessage⟩ ∉∅  $s'.cibaBindingMessages$  then
54:    stop
55:  if method ∉ {GET, POST} then
56:    stop
57:    let  $\overline{w'}$  := GETNAVIGABLEWINDOW( $\overline{w}$ , hrefwindow,  $\perp$ ,  $s'$ )
58:    let reference := ⟨REQ,  $s'.\overline{w'}$ .nonce⟩
59:    if method = GET then
60:      let body := ⟨⟩
61:      let parameters := data
62:      let origin :=  $\perp$ 
63:    else
64:      let body := data
65:      let parameters := url.parameters
66:      let origin := docorigin
67:    let req := ⟨HTTPReq,  $\nu_4$ , method, url.host, url.path, parameters, ⟨⟩, body⟩
68:    let  $s'$  := CANCELNAV(reference,  $s'$ )
69:    call HTTP_SEND(reference, req, url, origin, referrer, referrerPolicy, a,  $s'$ )
70:  case ⟨SETSCRIPT, window, script⟩
71:    let  $\overline{w'}$  := GETWINDOW( $\overline{w}$ , window,  $s'$ )
72:    let  $s'.\overline{w'}$ .activedocument.script := script
73:    stop ⟨⟩,  $s'$ 
74:  case ⟨SETSCRIPTSTATE, window, scriptstate⟩
75:    let  $\overline{w'}$  := GETWINDOW( $\overline{w}$ , window,  $s'$ )
76:    let  $s'.\overline{w'}$ .activedocument.scriptstate := scriptstate
77:    stop ⟨⟩,  $s'$ 
78:  case ⟨XMLHTTPREQUEST, url, method, data, xhrreference⟩
79:  if method ∈ {CONNECT, TRACE, TRACK} ∨ xhrreference ∉  $V_{\text{process}} \cup \{\perp\}$  then
80:    stop
81:  if url.host ≠ docorigin.host ∨ url.protocol ≠ docorigin.protocol then
82:    stop
83:  if method ∈ {GET, HEAD} then
84:    let data := ⟨⟩
85:    let origin :=  $\perp$ 
86:  else
87:    let origin := docorigin
88:  let req := ⟨HTTPReq,  $\nu_4$ , method, url.host, url.path, url.parameters, ⟨⟩, data⟩
89:  let reference := ⟨XHR,  $s'.\overline{d}$ .nonce, xhrreference⟩
90:  call HTTP_SEND(reference, req, url, origin, referrer, referrerPolicy, a,  $s'$ )
91:  case ⟨BACK, window⟩
92:  let  $\overline{w'}$  := GETNAVIGABLEWINDOW( $\overline{w}$ , window,  $\perp$ ,  $s'$ )
93:  call NAVBACK( $\overline{w'}$ ,  $s'$ )

```

→ Algorithm continues on next page.

---

---

```

94:   case ⟨FORWARD, window⟩
95:     let  $\bar{w}' := \text{GETNAVIGABLEWINDOW}(\bar{w}, \text{window}, \perp, s')$ 
96:     call NAVFORWARD( $\bar{w}'$ ,  $s'$ )
97:   case ⟨CLOSE, window⟩
98:     let  $\bar{w}' := \text{GETNAVIGABLEWINDOW}(\bar{w}, \text{window}, \perp, s')$ 
99:     remove  $s'.\bar{w}'$  from the sequence containing it
100:    stop  $\langle \rangle$ ,  $s'$ 
101:   case ⟨POSTMESSAGE, window, message, origin⟩
102:     let  $\bar{w}' \leftarrow \text{Subwindows}(s')$  such that  $s'.\bar{w}'.\text{nonce} \equiv \text{window}$ 
103:     if  $\exists \bar{j} \in \mathbb{N}$  such that  $s'.\bar{w}'.\text{documents}.\bar{j}.\text{active} \equiv \top$ 
104:        $\hookrightarrow \wedge (\text{origin} \neq \perp \implies s'.\bar{w}'.\text{documents}.\bar{j}.\text{origin} \equiv \text{origin})$  then
105:         let  $s'.\bar{w}'.\text{documents}.\bar{j}.\text{scriptinputs} := s'.\bar{w}'.\text{documents}.\bar{j}.\text{scriptinputs}$ 
106:          $\hookrightarrow + \langle \rangle$  ⟨POSTMESSAGE,  $s'.\bar{w}'.\text{nonce}$ ,  $\text{docorigin}$ ,  $\text{message}$ ⟩
107:     stop  $\langle \rangle$ ,  $s'$ 
108:   case else
109:     stop

```

---

- The function CANCELNAV (Algorithm 16) is used to stop any pending requests for a specific window. From the pending requests and pending DNS requests it removes any requests with the given window reference.
- The function HTTP\_SEND (Algorithm 17) takes an HTTP request *message* as input, adds cookie and origin headers to the message, creates a DNS request for the hostname given in the request and stores the request in  $s'.\text{pendingDNS}$  until the DNS resolution finishes. *reference* is a reference as defined in Definition 59. *wrl* contains the full URL of the request (this is mainly used to retrieve the protocol that should be used for this message, and to store the fragment identifier for use after the document was loaded). *origin* is the origin header value that is to be added to the HTTP request.
- The functions NAVBACK (Algorithm 18) and NAVFORWARD (Algorithm 19), navigate a window backward or forward. More precisely, they deactivate one document and activate that document's preceding document or succeeding document, respectively. If no such predecessor/successor exists, the functions do not change the state.
- The function RUNSCRIPT (Algorithm 20) performs a script execution step of the script in the document  $s'.\bar{d}$  (which is part of the window  $s'.\bar{w}$ ). A new script and document state is chosen according to the relation defined by the script and the new script and document state is saved. Afterwards, the *command* that the script issued is interpreted.
- The function PROCESSRESPONSE (Algorithm 21) is responsible for processing an HTTP response (*response*) that was received as the response to a request (*request*) that was sent earlier. *reference* is a reference as defined in Definition 59. *requestUrl* contains the URL used when retrieving the document.

The function first saves any cookies that were contained in the response to the browser state, then checks whether a redirection is requested (Location header). If that is not the case, the function creates a new document (for normal requests) or delivers the contents of the response to the respective receiver (for XHR responses).

c) *Browser Relation*: We can now define the relation  $R_{\text{webbrowser}}$  of a web browser atomic process as follows:

**Definition 66.** The pair  $((\langle a, f, m \rangle, s), (M, s'))$  belongs to  $R_{\text{webbrowser}}$  iff the non-deterministic Algorithm 22 (or any of the functions called therein), when given  $(\langle a, f, m \rangle, s)$  as input, terminates with **stop**  $M, s'$ , i.e., with output  $M$  and  $s'$ .

Recall that  $\langle a, f, m \rangle$  is an (input) event and  $s$  is a (browser) state,  $M$  is a sequence of (output) protoevents, and  $s'$  is a new (browser) state (potentially with placeholders for nonces).

#### H. Definition of Web Browsers

Finally, we define web browser atomic Dolev-Yao processes as follows:

**Definition 67 (Web Browser atomic Dolev-Yao Process).** A web browser atomic Dolev-Yao process is an atomic Dolev-Yao process of the form  $p = (I^p, Z_{\text{webbrowser}}, R_{\text{webbrowser}}, s_0^p)$  for a set  $I^p$  of addresses,  $Z_{\text{webbrowser}}$  and  $R_{\text{webbrowser}}$  as defined above, and an initial state  $s_0^p \in Z_{\text{webbrowser}}$ .

**Definition 68 (Web Browser Initial State).** An initial state  $s_0^p \in Z_{\text{webbrowser}}$  for a browser process  $p$  is a web browser state (Definition 60) with the following properties:

- $s_0^p.\text{windows} \equiv \langle \rangle$
- $s_0^p.\text{ids} \subset \langle \rangle \mathcal{T}_{\mathcal{N}}$  (intended to be constrained by instantiations of the Web Infrastructure Model)
- $s_0^p.\text{secrets} \in [\text{Origins} \times \mathcal{T}_{\mathcal{N}}]$  (intended to be constrained by instantiations of the Web Infrastructure Model)
- $s_0^p.\text{cookies} \equiv \langle \rangle$
- $s_0^p.\text{localStorage} \equiv \langle \rangle$
- $s_0^p.\text{sessionStorage} \equiv \langle \rangle$

---

**Algorithm 21** Web Browser Model: Process an HTTP response.

---

```
1: function PROCESSRESPONSE(response, reference, request, requestUrl, a, f, s')
2:   if Set-Cookie  $\in$  response.headers then
3:     for each  $c \in \langle \rangle$  response.headers [Set-Cookie],  $c \in$  Cookies do
4:       let s'.cookies [request.host]
          $\hookrightarrow :=$  AddCookie(s'.cookies [request.host], c, requestUrl.protocol)
5:   if Strict-Transport-Security  $\in$  response.headers  $\wedge$  requestUrl.protocol  $\equiv$  S then
6:     let s'.sts := s'.sts +  $\langle \rangle$  request.host
7:   if Referer  $\in$  request.headers then
8:     let referrer := request.headers[Referer]
9:   else
10:    let referrer :=  $\perp$ 
11:   if Location  $\in$  response.headers  $\wedge$  response.status  $\in$  {303, 307} then
12:     let url := response.headers [Location]
13:     if url.fragment  $\equiv$   $\perp$  then
14:       let url.fragment := requestUrl.fragment
15:     let method' := request.method
16:     let body' := request.body
17:     if Origin  $\in$  request.headers
          $\hookrightarrow \wedge$  request.headers[Origin]  $\neq$   $\diamond$ 
          $\hookrightarrow \wedge$  ( $\langle$ url.host, url.protocol $\rangle \equiv \langle$ request.host, requestUrl.protocol $\rangle$ )
          $\hookrightarrow \vee$  ( $\langle$ request.host, requestUrl.protocol $\rangle \equiv$  request.headers[Origin]) then
18:       let origin := request.headers[Origin]
19:     else
20:       let origin :=  $\diamond$ 
21:     if response.status  $\equiv$  303  $\wedge$  request.method  $\notin$  {GET, HEAD} then
22:       let method' := GET
23:       let body' :=  $\langle \rangle$ 
24:     if  $\exists \bar{w} \in$  Subwindows(s') such that s'. $\bar{w}$ .nonce  $\equiv$   $\pi_2$ (reference) then  $\rightarrow$  Do not redirect XHRs.
25:       let req :=  $\langle$ HTTPReq, v6, method', url.host, url.path, url.parameters,  $\langle \rangle$ , body' $\rangle$ 
26:       let referrerPolicy := response.headers[RefererPolicy]
27:       call HTTP_SEND(reference, req, url, origin, referrer, referrerPolicy, a, s')
28:     else
29:       stop  $\langle \rangle$ , s'
30:   switch  $\pi_1$ (reference) do
31:     case REQ
32:       let  $\bar{w} \leftarrow$  Subwindows(s') such that s'. $\bar{w}$ .nonce  $\equiv$   $\pi_2$ (reference) if possible;
          $\hookrightarrow$  otherwise stop  $\rightarrow$  normal response
33:       if response.body  $\not\sim$   $\langle$ *, * $\rangle$  then
34:         stop  $\langle \rangle$ , s'
35:       let script :=  $\pi_1$ (response.body)
36:       let scriptstate :=  $\pi_2$ (response.body)
37:       let d :=  $\langle$ v7, requestUrl, response.headers, referrer, script, scriptstate,  $\langle \rangle$ ,  $\langle \rangle$ ,  $\top$  $\rangle$ 
38:       if s'. $\bar{w}$ .documents  $\equiv$   $\langle \rangle$  then
39:         let s'. $\bar{w}$ .documents :=  $\langle$ d $\rangle$ 
40:       else
41:         let  $\bar{i} \leftarrow$   $\mathbb{N}$  such that s'. $\bar{w}$ .documents. $\bar{i}$ .active  $\equiv$   $\top$ 
42:         let s'. $\bar{w}$ .documents. $\bar{i}$ .active :=  $\perp$ 
43:         remove s'. $\bar{w}$ .documents. $(\bar{i} + 1)$  and all following documents
          $\hookrightarrow$  from s'. $\bar{w}$ .documents
44:         let s'. $\bar{w}$ .documents := s'. $\bar{w}$ .documents +  $\langle \rangle$  d
45:       stop  $\langle \rangle$ , s'
46:     case XHR
47:       let  $\bar{w} \leftarrow$  Subwindows(s'),  $\bar{d}$  such that s'. $\bar{d}$ .nonce  $\equiv$   $\pi_2$ (reference)
          $\hookrightarrow \wedge$  s'. $\bar{d}$  = s'. $\bar{w}$ .activedocument if possible; otherwise stop
          $\rightarrow$  process XHR response
48:       let headers := response.headers – Set-Cookie
49:       let s'. $\bar{d}$ .scriptinputs := s'. $\bar{d}$ .scriptinputs +  $\langle \rangle$ 
          $\langle$ XMLHTTPREQUEST, headers, response.body,  $\pi_3$ (reference) $\rangle$ 
50:       stop  $\langle \rangle$ , s'
```

---

- $s_0^p.\text{keyMapping} \in [\text{Doms} \times \mathcal{T}_{\mathcal{N}}]$  (intended to be constrained by instantiations of the Web Infrastructure Model)
- $s_0^p.\text{sts} \equiv \langle \rangle$
- $s_0^p.\text{DNSaddress} \in \text{IPs}$  (note that this includes the possibility of using an attacker-controlled address)
- $s_0^p.\text{pendingDNS} \equiv \langle \rangle$
- $s_0^p.\text{pendingRequests} \equiv \langle \rangle$
- $s_0^p.\text{isCorrupted} \equiv \perp$
- $s_0^p.\text{cibaBindingMessages} \equiv \langle \rangle$
- $s_0^p.\text{tlskeys} \equiv \text{tlskeys}^p$  (see Appendix D-C)

Note that instantiations of the Web Infrastructure Model may define different conditions for a web browser's initial state.

## I. Helper Functions

In order to simplify the description of scripts, we use several helper functions.

a) *CHOOSEINPUT* (Algorithm 23): The state of a document contains a term, say *scriptinputs*, which records the input this document has obtained so far (via XHRs and postMessages). If the script of the document is activated, it will typically need to pick one input message from *scriptinputs* and record which input it has already processed. For this purpose, the function  $\text{CHOOSEINPUT}(s', \text{scriptinputs})$  is used, where  $s'$  denotes the script's current state. It saves the indexes of already handled messages in the scriptstate  $s'$  and chooses a yet unhandled input message from *scriptinputs*. The index of this message is then saved in the scriptstate (which is returned to the script).

b) *CHOOSEFIRSTINPUTPAT* (Algorithm 24): Similar to the function *CHOOSEINPUT* above, we define the function *CHOOSEFIRSTINPUTPAT*. This function takes the term *scriptinputs*, which as above records the input this document has obtained so far (via XHRs and postMessages, append-only), and a pattern. If called, this function chooses the first message in *scriptinputs* that matches *pattern* and returns it. This function is typically used in places, where a script only processes the first message that matches the pattern. Hence, we omit recording the usage of an input.

c) *PARENTWINDOW*: To determine the nonce referencing the parent window in the browser, the function  $\text{PARENTWINDOW}(\text{tree}, \text{docnonce})$  is used. It takes the term *tree*, which is the (partly cleaned) tree of browser windows the script is able to see and the document nonce *docnonce*, which is the nonce referencing the current document the script is running in, as input. It outputs the nonce referencing the window which directly contains in its subwindows the window of the document referenced by *docnonce*. If there is no such window (which is the case if the script runs in a document of a top-level window), *PARENTWINDOW* returns  $\perp$ .

d) *PARENTDOCNONCE*: The function  $\text{PARENTDOCNONCE}(\text{tree}, \text{docnonce})$  determines (similar to *PARENTWINDOW* above) the nonce referencing the active document in the parent window in the browser. It takes the term *tree*, which is the (partly cleaned) tree of browser windows the script is able to see and the document nonce *docnonce*, which is the nonce referencing the current document the script is running in, as input. It outputs the nonce referencing the active document in the window which directly contains in its subwindows the window of the document referenced by *docnonce*. If there is no such window (which is the case if the script runs in a document of a top-level window) or no active document, *PARENTDOCNONCE* returns *docnonce*.

e) *SUBWINDOWS*: This function takes a term *tree* and a document nonce *docnonce* as input just as the function above. If *docnonce* is not a reference to a document contained in *tree*, then  $\text{SUBWINDOWS}(\text{tree}, \text{docnonce})$  returns  $\langle \rangle$ . Otherwise, let  $\langle \text{docnonce}, \text{location}, \langle \rangle, \text{referrer}, \text{script}, \text{scriptstate}, \text{scriptinputs}, \text{subwindows}, \text{active} \rangle$  denote the subterm of *tree* corresponding to the document referred to by *docnonce*. Then,  $\text{SUBWINDOWS}(\text{tree}, \text{docnonce})$  returns *subwindows*.

f) *AUXWINDOW*: This function takes a term *tree* and a document nonce *docnonce* as input as above. From all window terms in *tree* that have the window containing the document identified by *docnonce* as their opener, it selects one non-deterministically and returns its nonce. If there is no such window, it returns the nonce of the window containing *docnonce*.

g) *AUXDOCNONCE*: Similar to *AUXWINDOW* above, the function *AUXDOCNONCE* takes a term *tree* and a document nonce *docnonce* as input. From all window terms in *tree* that have the window containing the document identified by *docnonce* as their opener, it selects one non-deterministically and returns its active document's nonce. If there is no such window or no active document, it returns *docnonce*.

h) *OPENERWINDOW*: This function takes a term *tree* and a document nonce *docnonce* as input as above. It returns the window nonce of the opener window of the window that contains the document identified by *docnonce*. Recall that the nonce identifying the opener of each window is stored inside the window term. If no document with nonce *docnonce* is found in the tree *tree* or the document with nonce *docnonce* is not directly contained in a top-level window,  $\diamond$  is returned.

i) *GETWINDOW*: This function takes a term *tree* and a document nonce *docnonce* as input as above. It returns the nonce of the window containing *docnonce*.

j) *GETORIGIN*: To extract the origin of a document, the function  $\text{GETORIGIN}(\text{tree}, \text{docnonce})$  is used. This function searches for the document with the identifier *docnonce* in the (cleaned) tree *tree* of the browser's windows and documents. It returns the origin *o* of the document. If no document with nonce *docnonce* is found in the tree *tree*,  $\diamond$  is returned.

---

**Algorithm 22** Web Browser Model: Main Algorithm.

---

**Input:**  $\langle a, f, m \rangle, s$

- 1: **let**  $s' := s$
- 2: **if**  $s.isCorrupted \neq \perp$  **then**
- 3:   **let**  $s'.pendingRequests := \langle m, s.pendingRequests \rangle$   $\rightarrow$  Collect incoming messages
- 4:   **let**  $m' \leftarrow dv(s')$
- 5:   **let**  $a' \leftarrow IPs$
- 6:   **stop**  $\langle \langle a', a, m' \rangle \rangle, s'$
- 7: **if**  $m \equiv TRIGGER$  **then**  $\rightarrow$  A special trigger message.
- 8:   **let**  $switch \leftarrow \{script, urlbar, reload, forward, back\}$
- 9:   **if**  $switch \equiv script$  **then**  $\rightarrow$  Run some script.
- 10:    **let**  $\bar{w} \leftarrow Subwindows(s')$  **such that**  $s'.\bar{w}.documents \neq \langle \rangle$   
     $\hookrightarrow$  **if possible; otherwise stop**  $\rightarrow$  Pointer to some window.
- 11:    **let**  $\bar{d} := \bar{w} + \langle \rangle$  **activedocument**
- 12:    **call**  $RUNSCRIPT(\bar{w}, \bar{d}, a, s')$
- 13:   **else if**  $switch \equiv urlbar$  **then**  $\rightarrow$  Create some new request.
- 14:    **let**  $newwindow \leftarrow \{\top, \perp\}$
- 15:    **if**  $newwindow \equiv \top$  **then**  $\rightarrow$  Create a new window.
- 16:      **let**  $windownonce := \nu_1$
- 17:      **let**  $w' := \langle windownonce, \langle \rangle, \perp \rangle$
- 18:      **let**  $s'.windows := s'.windows + \langle \rangle w'$
- 19:    **else**  $\rightarrow$  Use existing top-level window.
- 20:    **let**  $tlw \leftarrow N$  **such that**  $s'.tlw.documents \neq \langle \rangle$   
     $\hookrightarrow$  **if possible; otherwise stop**  $\rightarrow$  Pointer to some top-level window.
- 21:    **let**  $windownonce := s'.tlw.nonce$
- 22:    **let**  $protocol \leftarrow \{P, S\}$
- 23:    **let**  $host \leftarrow Doms$
- 24:    **let**  $path \leftarrow S$
- 25:    **let**  $fragment \leftarrow S$
- 26:    **let**  $parameters \leftarrow [S \times S]$
- 27:    **let**  $body := \langle \rangle$
- 28:    **let**  $startciba \leftarrow \{\top, \perp\}$
- 29:    **if**  $startciba \equiv \top$  **then**
- 30:      **let**  $body[authServ] \leftarrow Doms$
- 31:      **let**  $body[identity] \leftarrow s'.ids$
- 32:    **let**  $url := \langle URL, protocol, host, path, parameters, fragment \rangle$
- 33:    **let**  $req := \langle HTTPReq, \nu_2, GET, host, path, parameters, \langle \rangle, body \rangle$
- 34:    **call**  $HTTP\_SEND(\langle REQ, windownonce \rangle, req, url, \perp, \perp, \perp, a, s')$
- 35:   **else if**  $switch \equiv reload$  **then**  $\rightarrow$  Reload some document.
- 36:    **let**  $\bar{w} \leftarrow Subwindows(s')$  **such that**  $s'.\bar{w}.documents \neq \langle \rangle$   
     $\hookrightarrow$  **if possible; otherwise stop**  $\rightarrow$  Pointer to some window.
- 37:    **let**  $url := s'.\bar{w}.activedocument.location$
- 38:    **let**  $req := \langle HTTPReq, \nu_2, GET, url.host, url.path, url.parameters, \langle \rangle, \langle \rangle \rangle$
- 39:    **let**  $referrer := s'.\bar{w}.activedocument.referrer$
- 40:    **let**  $s' := CANCELNAV(s'.\bar{w}.nonce, s')$
- 41:    **call**  $HTTP\_SEND(\langle REQ, s'.\bar{w}.nonce \rangle, req, url, \perp, referrer, \perp, a, s')$
- 42:   **else if**  $switch \equiv forward$  **then**
- 43:    **let**  $\bar{w} \leftarrow Subwindows(s')$  **such that**  $s'.\bar{w}.documents \neq \langle \rangle$   
     $\hookrightarrow$  **if possible; otherwise stop**  $\rightarrow$  Pointer to some window.
- 44:    **call**  $NAVFORWARD(\bar{w}, s')$
- 45:   **else if**  $switch \equiv back$  **then**
- 46:    **let**  $\bar{w} \leftarrow Subwindows(s')$  **such that**  $s'.\bar{w}.documents \neq \langle \rangle$   
     $\hookrightarrow$  **if possible; otherwise stop**  $\rightarrow$  Pointer to some window.
- 47:    **call**  $NAVBACK(\bar{w}, s')$
- 48: **else if**  $m \equiv FULLCORRUPT$  **then**  $\rightarrow$  Request to corrupt browser
- 49:   **let**  $s'.isCorrupted := FULLCORRUPT$
- 50:   **stop**  $\langle \rangle, s'$
- 51: **else if**  $m \equiv CLOSECORRUPT$  **then**  $\rightarrow$  Close the browser
- 52:   **let**  $s'.secrets := \langle \rangle$
- 53:   **let**  $s'.windows := \langle \rangle$
- 54:   **let**  $s'.pendingDNS := \langle \rangle$
- 55:   **let**  $s'.pendingRequests := \langle \rangle$
- 56:   **let**  $s'.sessionStorage := \langle \rangle$
- 57:   **let**  $s'.cookies \subset \langle \rangle$  **Cookies such that**  
     $\hookrightarrow (c \in \langle \rangle s'.cookies) \iff (c \in \langle \rangle s.cookies \wedge c.content.session \equiv \perp)$
- 58:   **let**  $s'.isCorrupted := CLOSECORRUPT$
- 59:   **stop**  $\langle \rangle, s'$

---

---

```

60: else if  $\exists \langle \text{reference}, \text{request}, \text{url}, \text{key}, f \rangle \in {}^\diamond s'.\text{pendingRequests}$  such that
     $\hookrightarrow \pi_1(\text{dec}_s(m, \text{key})) \equiv \text{HTTPResp}$  then  $\rightarrow$  Encrypted HTTP response
61:   let  $m' := \text{dec}_s(m, \text{key})$ 
62:   if  $m'.\text{nonce} \neq \text{request.nonce}$  then
63:     stop
64:   remove  $\langle \text{reference}, \text{request}, \text{url}, \text{key}, f \rangle$  from  $s'.\text{pendingRequests}$ 
65:   if  $\text{binding\_message} \in {}^\diamond m'.\text{body}$  then
66:     let  $s'.\text{cibaBindingMessages} := s'.\text{cibaBindingMessages} + {}^\diamond \langle \text{request.host}, m'.\text{body}[\text{binding\_message}] \rangle$ 
67:     call  $\text{PROCESSRESPONSE}(m', \text{reference}, \text{request}, \text{url}, a, f, s')$ 
68: else if  $\pi_1(m) \equiv \text{HTTPResp} \wedge \exists \langle \text{reference}, \text{request}, \text{url}, \perp, f \rangle \in {}^\diamond s'.\text{pendingRequests}$  such that
     $\hookrightarrow m.\text{nonce} \equiv \text{request.nonce}$  then  $\rightarrow$  Plain HTTP Response
69:   remove  $\langle \text{reference}, \text{request}, \text{url}, \perp, f \rangle$  from  $s'.\text{pendingRequests}$ 
70:   call  $\text{PROCESSRESPONSE}(m, \text{reference}, \text{request}, \text{url}, a, f, s')$ 
71: else if  $m \in \text{DNSResponses}$  then  $\rightarrow$  Successful DNS response
72:   if  $m.\text{nonce} \notin s'.\text{pendingDNS} \vee m.\text{result} \notin \text{IPs}$ 
     $\hookrightarrow \vee m.\text{domain} \neq s'.\text{pendingDNS}[m.\text{nonce}].\text{request.host}$  then
73:     stop
74:   let  $\langle \text{reference}, \text{message}, \text{url} \rangle := s'.\text{pendingDNS}[m.\text{nonce}]$ 
75:   if  $\text{url.protocol} \equiv \text{S}$  then
76:     let  $s'.\text{pendingRequests} := s'.\text{pendingRequests}$ 
     $\hookrightarrow + {}^\diamond \langle \text{reference}, \text{message}, \text{url}, \nu_3, m.\text{result} \rangle$ 
77:     let  $\text{message} := \text{enc}_a(\langle \text{message}, \nu_3 \rangle, s'.\text{keyMapping}[\text{message.host}])$ 
78:   else
79:     let  $s'.\text{pendingRequests} := s'.\text{pendingRequests}$ 
     $\hookrightarrow + {}^\diamond \langle \text{reference}, \text{message}, \text{url}, \perp, m.\text{result} \rangle$ 
80:   let  $s'.\text{pendingDNS} := s'.\text{pendingDNS} - m.\text{nonce}$ 
81:   stop  $\langle \langle m.\text{result}, a, \text{message} \rangle \rangle, s'$ 
82: else if  $\exists m_{\text{dec}}, k, k', \text{inDomain}$  such that  $\langle m_{\text{dec}}, k \rangle \equiv \text{dec}_a(m, k') \wedge \langle \text{inDomain}, k' \rangle \in s'.\text{tlskeys}$  then
     $\rightarrow$  For modelling CIBA, we allow the browser to receive requests. By this, the AS can contact its users and ask to give their consent
    for a given CIBA flow
83:   let  $n, \text{method}, \text{path}, \text{parameters}, \text{headers}, \text{body}$  such that
     $\hookrightarrow \langle \text{HTTPReq}, n, \text{method}, \text{inDomain}, \text{path}, \text{parameters}, \text{headers}, \text{body} \rangle \equiv m_{\text{dec}}$ 
     $\hookrightarrow$  if possible; otherwise stop
84:   if  $\text{path} \neq /start - \text{ciba} - \text{authentication}$  then stop
85:   let  $\text{newwindow} \leftarrow \{\top, \perp\}$ 
86:   if  $\text{newwindow} \equiv \top$  then  $\rightarrow$  Create a new window.
87:     let  $\text{windownonce} := \nu_1$ 
88:     let  $w' := \langle \text{windownonce}, \langle \rangle, \perp \rangle$ 
89:     let  $s'.\text{windows} := s'.\text{windows} + {}^\diamond w'$ 
90:   else  $\rightarrow$  Use existing top-level window.
91:     let  $\text{tlw} \leftarrow \mathbb{N}$  such that  $s'.\text{tlw}.\text{documents} \neq \langle \rangle$ 
     $\hookrightarrow$  if possible; otherwise stop  $\rightarrow$  Pointer to some top-level window.
92:     let  $\text{windownonce} := s'.\text{tlw}.\text{nonce}$ 
93:     let  $\text{url} := \text{body}[\text{ciba\_url}]$ 
94:     let  $\text{req} := \langle \text{HTTPReq}, \nu_{\text{ciba\_req}}, \text{POST}, \text{url.host}, \varepsilon, \langle \rangle, \langle \rangle, \text{body} \rangle$ 
95:     call  $\text{HTTP\_SEND}(\langle \text{REQ}, \text{windownonce} \rangle, \text{req}, \text{url}, \perp, \perp, \perp, a, s')$ 
96: stop

```

---

**Algorithm 23** Function to retrieve an unhandled input message for a script.

---

```

1: function  $\text{CHOOSEINPUT}(s', \text{scriptinputs})$ 
2:   let  $iid$  such that  $iid \in \{1, \dots, |\text{scriptinputs}|\} \wedge iid \notin {}^\diamond s'.\text{handledInputs}$  if possible;
     $\hookrightarrow$  otherwise return  $(\perp, s')$ 
3:   let  $\text{input} := \pi_{iid}(\text{scriptinputs})$ 
4:   let  $s'.\text{handledInputs} := s'.\text{handledInputs} + {}^\diamond iid$ 
5:   return  $(\text{input}, s')$ 

```

---

**Algorithm 24** Function to extract the first script input message matching a specific pattern.

---

```

1: function  $\text{CHOOSEFIRSTINPUTPAT}(\text{scriptinputs}, \text{pattern})$ 
2:   let  $i$  such that  $i = \min\{j : \pi_j(\text{scriptinputs}) \sim \text{pattern}\}$  if possible; otherwise return  $\perp$ 
3:   return  $\pi_i(\text{scriptinputs})$ 

```

---

---

**Algorithm 25** Relation of a DNS server  $R^d$ .

---

**Input:**  $\langle a, f, m \rangle, s$ 

```
1: let  $domain, n$  such that  $\langle \text{DNSResolve}, domain, n \rangle \equiv m$  if possible; otherwise stop  $\langle \rangle, s$ 
2: if  $domain \in s$  then
3:   let  $addr := s[domain]$ 
4:   let  $m' := \langle \text{DNSResolved}, domain, addr, n \rangle$ 
5:   stop  $\langle \langle f, a, m' \rangle \rangle, s$ 
6: stop  $\langle \rangle, s$ 
```

---

k) *GETPARAMETERS*: Works exactly as *GETORIGIN*, but returns the document's parameters instead.

### J. DNS Servers

*Definition 69.* A *DNS server*  $d$  (in a flat DNS model) is modeled in a straightforward way as an atomic DY process  $(I^d, \{s_0^d\}, R^d, s_0^d)$ . It has a finite set of addresses  $I^d$  and its initial (and only) state  $s_0^d$  encodes a mapping from domain names to addresses of the form

$$s_0^d = \langle \langle \text{domain}_1, a_1 \rangle, \langle \text{domain}_2, a_2 \rangle, \dots \rangle .$$

DNS queries are answered according to this table (if the requested DNS name cannot be found in the table, the request is ignored).

The relation  $R^d \subseteq (\mathcal{E} \times \{s_0^d\}) \times (2^{\mathcal{E}} \times \{s_0^d\})$  of  $d$  above is defined by Algorithm 25.

### K. Web Systems

The web infrastructure and web applications are formalized by what is called a web system. A web system contains, among others, a (possibly infinite) set of DY processes, modeling web browsers, web servers, DNS servers, and attackers (which may corrupt other entities, such as browsers).

*Definition 70.* A *web system*  $\mathcal{WS} = (\mathcal{W}, \mathcal{S}, \text{script}, E^0)$  is a tuple with its components defined as follows:

The first component,  $\mathcal{W}$ , denotes a system (a set of DY processes) and is partitioned into the sets Hon, Web, and Net of honest, web attacker, and network attacker processes, respectively.

Every  $p \in \text{Web} \cup \text{Net}$  is an attacker process for some set of sender addresses  $A \subseteq \text{IPs}$ . For a web attacker  $p \in \text{Web}$ , we require its set of addresses  $I^p$  to be disjoint from the set of addresses of all other web attackers and honest processes, i.e.,  $I^p \cap I^{p'} = \emptyset$  for all  $p' \neq p, p' \in \text{Hon} \cup \text{Web}$ . Hence, a web attacker cannot listen to traffic intended for other processes. Also, we require that  $A = I^p$ , i.e., a web attacker can only use sender addresses it owns. Conversely, a network attacker may listen to all addresses (i.e., no restrictions on  $I^p$ ) and may spoof all addresses (i.e., the set  $A$  may be IPs).

Every  $p \in \text{Hon}$  is a DY process which models either a *web server*, a *web browser*, or a *DNS server*. Just as for web attackers, we require that  $p$  does not spoof sender addresses and that its set of addresses  $I^p$  is disjoint from those of other honest processes and the web attackers.

The second component,  $\mathcal{S}$ , is a finite set of scripts such that  $R^{\text{att}} \in \mathcal{S}$ . The third component, *script*, is an injective mapping from  $\mathcal{S}$  to  $\mathbb{S}$ , i.e., by *script* every  $s \in \mathcal{S}$  is assigned its string representation *script*( $s$ ).

Finally,  $E^0$  is an (infinite) sequence of events, containing an infinite number of events of the form  $\langle a, a, \text{TRIGGER} \rangle$  for every  $a \in \bigcup_{p \in \mathcal{W}} I^p$ .

A *run* of  $\mathcal{WS}$  is a run of  $\mathcal{W}$  initiated by  $E^0$ .

### L. Generic HTTPS Server Model

This base model can be used to ease modeling of HTTPS server atomic processes. It defines placeholder algorithms that can be superseded by more detailed algorithms to describe a concrete relation for an HTTPS server.

*Definition 71 (Base state for an HTTPS server).* The state of each HTTPS server that is an instantiation of this relation must contain at least the following subterms:  $pendingDNS \in [\mathcal{N} \times \mathcal{T}_{\mathcal{N}}]$ ,  $pendingRequests \in \mathcal{T}_{\mathcal{N}}$  (both containing arbitrary terms),  $DNSaddress \in \text{IPs}$  (containing the IP address of a DNS server),  $keyMapping \in [\text{Doms} \times \mathcal{T}_{\mathcal{N}}]$  (containing a mapping from domains to public keys),  $tlskeys \in [\text{Doms} \times \mathcal{N}]$  (containing a mapping from domains to private keys), and  $corrupt \in \mathcal{T}_{\mathcal{N}}$  (either  $\perp$  if the server is not corrupted, or an arbitrary term otherwise).

We note that in concrete instantiations of the generic HTTPS server model, there is no need to extract information from these subterms or alter these subterms.

Let  $\nu_{n_0}$  and  $\nu_{n_1}$  denote placeholders for nonces that are not used in the concrete instantiation of the server. We now define the default functions of the generic web server in Algorithms 26–30, and the main relation in Algorithm 31.



---

**Algorithm 26** Generic HTTPS Server Model: Sending a DNS message (in preparation for sending an HTTPS message).

---

```
1: function HTTPS_SIMPLE_SEND(reference, message, a, s')
2:   let s'.pendingDNS[ $\nu_{n0}$ ] := (reference, message)
3:   stop  $\langle\langle s'.\text{DNSaddress}, a, \langle \text{DNSResolve}, \text{message.host}, \nu_{n0} \rangle \rangle\rangle, s'$ 
```

---

---

**Algorithm 27** Generic HTTPS Server Model: Default HTTPS response handler.

---

```
1: function PROCESS_HTTPS_RESPONSE(m, reference, request, a, f, s')
2:   stop
```

---

---

**Algorithm 28** Generic HTTPS Server Model: Default trigger event handler.

---

```
1: function PROCESS_TRIGGER(a, s')
2:   stop
```

---

---

**Algorithm 29** Generic HTTPS Server Model: Default HTTPS request handler.

---

```
1: function PROCESS_HTTPS_REQUEST(m, k, a, f, s')
2:   stop
```

---

---

**Algorithm 30** Generic HTTPS Server Model: Default handler for other messages.

---

```
1: function PROCESS_OTHER(m, a, f, s')
2:   stop
```

---

---

**Algorithm 31** Generic HTTPS Server Model: Main relation of a generic HTTPS server

---

**Input:**  $\langle a, f, m \rangle, s$

- 1: **let**  $s' := s$
- 2: **if**  $s'.\text{corrupt} \neq \perp \vee m \equiv \text{CORRUPT}$  **then**
- 3:     **let**  $s'.\text{corrupt} := \langle \langle a, f, m \rangle, s'.\text{corrupt} \rangle$
- 4:     **let**  $m' \leftarrow d_V(s')$
- 5:     **let**  $a' \leftarrow \text{IPs}$
- 6:     **stop**  $\langle \langle a', a, m' \rangle \rangle, s'$
- 7: **if**  $\exists m_{\text{dec}}, k, k', \text{inDomain}$  **such that**  $\langle m_{\text{dec}}, k \rangle \equiv \text{dec}_a(m, k') \wedge \langle \text{inDomain}, k' \rangle \in s.\text{tlskeys}$  **then**
- 8:     **let**  $n, \text{method}, \text{path}, \text{parameters}, \text{headers}, \text{body}$  **such that**  
       $\hookrightarrow \langle \text{HTTPReq}, n, \text{method}, \text{inDomain}, \text{path}, \text{parameters}, \text{headers}, \text{body} \rangle \equiv m_{\text{dec}}$   
       $\hookrightarrow$  **if possible; otherwise stop**
- 9:     **call**  $\text{PROCESS\_HTTPS\_REQUEST}(m_{\text{dec}}, k, a, f, s')$
- 10: **else if**  $m \in \text{DNSResponses}$  **then**  $\rightarrow$  **Successful DNS response**
- 11:     **if**  $m.\text{nonce} \notin s.\text{pendingDNS} \vee m.\text{result} \notin \text{IPs}$   
       $\hookrightarrow \vee m.\text{domain} \neq s.\text{pendingDNS}[m.\text{nonce}].2.\text{host}$  **then**
- 12:         **stop**
- 13:     **let**  $\text{reference} := s.\text{pendingDNS}[m.\text{nonce}].1$
- 14:     **let**  $\text{request} := s.\text{pendingDNS}[m.\text{nonce}].2$
- 15:     **let**  $s'.\text{pendingRequests} := s'.\text{pendingRequests} + \langle \text{reference}, \text{request}, \nu_{n1}, m.\text{result} \rangle$
- 16:     **let**  $\text{message} := \text{enc}_a(\langle \text{request}, \nu_{n1} \rangle, s'.\text{keyMapping}[\text{request}.\text{host}])$
- 17:     **let**  $s'.\text{pendingDNS} := s'.\text{pendingDNS} - m.\text{nonce}$
- 18:     **stop**  $\langle \langle m.\text{result}, a, \text{message} \rangle \rangle, s'$
- 19: **else if**  $\exists \langle \text{reference}, \text{request}, \text{key}, f \rangle \in s'.\text{pendingRequests}$   
       $\hookrightarrow$  **such that**  $\pi_1(\text{dec}_s(m, \text{key})) \equiv \text{HTTPResp}$  **then**  $\rightarrow$  **Encrypted HTTP response**
- 20:     **let**  $m' := \text{dec}_s(m, \text{key})$
- 21:     **if**  $m'.\text{nonce} \neq \text{request}.\text{nonce}$  **then**
- 22:         **stop**
- 23:     **if**  $m' \notin \text{HTTPResponses}$  **then**
- 24:         **call**  $\text{PROCESS\_OTHER}(m, a, f, s')$
- 25:     **remove**  $\langle \text{reference}, \text{request}, \text{key}, f \rangle$  **from**  $s'.\text{pendingRequests}$
- 26:     **call**  $\text{PROCESS\_HTTPS\_RESPONSE}(m', \text{reference}, \text{request}, a, f, s')$
- 27: **else if**  $m \equiv \text{TRIGGER}$  **then**  $\rightarrow$  **Process was triggered**
- 28:     **call**  $\text{PROCESS\_TRIGGER}(a, s')$
- 29: **else**
- 30:     **call**  $\text{PROCESS\_OTHER}(m, a, f, s')$
- 31: **stop**

---

### M. General Security Properties of the WIM

We now repeat general application independent security properties of the WIM [34].

Let  $\mathcal{WS} = (\mathcal{W}, \mathcal{S}, \text{script}, E_0)$  be a web system. In the following, we write  $s_x = (S_x, E_x)$  for the states of a web system.

**Definition 72 (Emitting Events).** Given an atomic process  $p$ , an event  $e$ , and a finite run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  or an infinite run  $\rho = ((S^0, E^0, N^0), \dots)$  we say that  $p$  emits  $e$  iff there is a processing step in  $\rho$  of the form

$$(S^i, E^i, N^i) \xrightarrow[p \rightarrow E]{} (S^{i+1}, E^{i+1}, N^{i+1})$$

for some  $i \geq 0$  and a sequence of events  $E$  with  $e \in E$ . We also say that  $p$  emits  $m$  iff  $e = \langle x, y, m \rangle$  for some addresses  $x, y$ .

**Definition 73.** We say that a term  $t$  is derivably contained in (a term)  $t'$  for (a set of DY processes)  $P$  (in a processing step  $s_i \rightarrow s_{i+1}$  of a run  $\rho = (s_0, s_1, \dots)$ ) if  $t$  is derivable from  $t'$  with the knowledge available to  $P$ , i.e.,

$$t \in d_\emptyset(\{t'\} \cup \bigcup_{p \in P} S^{i+1}(p))$$

**Definition 74.** We say that a set of processes  $P$  leaks a term  $t$  (in a processing step  $s_i \rightarrow s_{i+1}$ ) to a set of processes  $P'$  if there exists a message  $m$  that is emitted (in  $s_i \rightarrow s_{i+1}$ ) by some  $p \in P$  and  $t$  is derivably contained in  $m$  for  $P'$  in the processing step  $s_i \rightarrow s_{i+1}$ . If we omit  $P'$ , we define  $P' := \mathcal{W} \setminus P$ . If  $P$  is a set with a single element, we omit the set notation.

**Definition 75.** We say that a DY process  $p$  created a message  $m$  in a processing step

$$(S^i, E^i, N^i) \xrightarrow[p \rightarrow E_{\text{out}}]{e_{\text{in}} \rightarrow p} (S^{i+1}, E^{i+1}, N^{i+1})$$

of a run  $\rho = ((S^0, E^0, N^0), \dots, (S^n, E^n, N^n))$  if all of the following hold true

- $m$  is a subterm of one of the events in  $E_{\text{out}}$
- $m$  is and was not derivable by any other set of processes

$$m \notin d_\emptyset\left(\bigcup_{\substack{p' \in \mathcal{W} \setminus \{p\} \\ 0 \leq j \leq i+1}} S^j(p')\right)$$

We note a process  $p$  creating a message does not imply that  $p$  can derive that message.

**Definition 76.** We say that a browser  $b$  accepted a message (as a response to some request) if the browser decrypted the message (if it was an HTTPS message) and called the function PROCESSRESPONSE, passing the message and the request (see Algorithm 21).

**Definition 77.** We say that an atomic DY process  $p$  knows a term  $t$  in some state  $s = (S, E, N)$  of a run if it can derive the term from its knowledge, i.e.,  $t \in d_\emptyset(S(p))$ .

**Definition 78.** Let  $N \subseteq \mathcal{N}$ ,  $t \in \mathcal{T}_N(X)$ , and  $k \in \mathcal{T}_N(X)$ . We say that  $k$  appears only as a public key in  $t$ , if

- 1) If  $t \in N \cup X$ , then  $t \neq k$
- 2) If  $t = f(t_1, \dots, t_n)$ , for  $f \in \Sigma$  and  $t_i \in \mathcal{T}_{\mathcal{N}}(X)$  ( $i \in \{1, \dots, n\}$ ), then  $f = \text{pub}$  or for all  $t_i$ ,  $k$  appears only as a public key in  $t_i$ .

**Definition 79.** We say that a script initiated a request  $r$  if a browser triggered the script (in Line 10 of Algorithm 20) and the first component of the *command* output of the script relation is either HREF, IFRAME, FORM, or XMLHTTPREQUEST such that the browser issues the request  $r$  in the same step as a result.

**Definition 80.** We say that an instance of the generic HTTPS server  $s$  accepted a message (as a response to some request) if the server decrypted the message (if it was an HTTPS message) and called the function PROCESS\_HTTPS\_RESPONSE, passing the message and the request (see Algorithm 31).

For a run  $\rho = s_0, s_1, \dots$  of any  $\mathcal{WS}$ , we state the following lemmas:

**Lemma 24.** If in the processing step  $s_i \rightarrow s_{i+1}$  of a run  $\rho$  of  $\mathcal{WS}$  an honest browser  $b$

- (I) emits an HTTPS request of the form

$$m = \text{enc}_a(\langle \text{req}, k \rangle, \text{pub}(k'))$$

(where  $\text{req}$  is an HTTP request,  $k$  is a nonce (symmetric key), and  $k'$  is the private key of some other DY process  $u$ ), and

- (II) in the initial state  $s_0$ , for all processes  $p \in \mathcal{W} \setminus \{u\}$ , the private key  $k'$  appears only as a public key in  $S^0(p)$ , and
- (III)  $u$  never leaks  $k'$ ,

then all of the following statements are true:

- (1) There is no state of  $\mathcal{MS}$  where any party except for  $u$  knows  $k'$ , thus no one except for  $u$  can decrypt  $m$  to obtain  $req$ .
- (2) If there is a processing step  $s_j \rightarrow s_{j+1}$  where the browser  $b$  leaks  $k$  to  $\mathcal{W} \setminus \{u, b\}$  there is a processing step  $s_h \rightarrow s_{h+1}$  with  $h < j$  where  $u$  leaks the symmetric key  $k$  to  $\mathcal{W} \setminus \{u, b\}$  or the browser is fully corrupted in  $s_j$ .
- (3) The value of the host header in  $req$  is the domain that is assigned the public key  $\text{pub}(k')$  in the browsers' keymapping  $s_0.\text{keyMapping}$  (in its initial state).
- (4) If  $b$  accepts a response (say,  $m'$ ) to  $m$  in a processing step  $s_j \rightarrow s_{j+1}$  and  $b$  is honest in  $s_j$  and  $u$  did not leak the symmetric key  $k$  to  $\mathcal{W} \setminus \{u, b\}$  prior to  $s_j$ , then  $u$  created the HTTPS response  $m'$  to the HTTPS request  $m$ , i.e., the nonce of the HTTP request  $req$  is not known to any atomic process  $p$ , except for the atomic processes  $b$  and  $u$ .

PROOF. (1) follows immediately from the preconditions.

The process  $u$  never leaks  $k'$ , and initially, the private key  $k'$  appears only as a public key in all other process states. As the equational theory does not allow the extraction of a private key  $x$  from a public key  $\text{pub}(x)$ , the other processes can never derive  $k'$ .

Thus, even with the knowledge of all nonces (except for those of  $u$ ),  $k'$  can never be derived from any network output of  $u$ , and  $k'$  cannot be known to any other party. Thus, nobody except for  $u$  can derive  $req$  from  $m$ .

(2) We assume that  $b$  leaks  $k$  to  $\mathcal{W} \setminus \{u, b\}$  in the processing step  $s_j \rightarrow s_{j+1}$  without  $u$  prior leaking the key  $k$  to anyone except for  $u$  and  $b$  and that the browser is not fully corrupted in  $s_j$ , and lead this to a contradiction.

The browser is honest in  $s_i$ . From the definition of the browser  $b$ , we see that the key  $k$  is always chosen as a fresh nonce (placeholder  $\nu_3$  in Lines 71ff. of Algorithm 22) that is not used anywhere else. Further, the key is stored in the browser's state in *pendingRequests*. The information from *pendingRequests* is not extracted or used anywhere else (in particular it is not accessible by scripts). If the browser becomes closecorrupted prior to  $s_j$  (and after  $s_i$ ), the key cannot be used anymore (compare Lines 51ff. of Algorithm 22). Hence,  $b$  does not leak  $k$  to any other party in  $s_j$  (except for  $u$  and  $b$ ). This proves (2).

(3) Per the definition of browsers (Algorithm 22), a host header is always contained in HTTP requests by browsers. From Line 77 of Algorithm 22 we can see that the encryption key for the request  $req$  was chosen using the host header of the message. It is chosen from the *keyMapping* in the browser's state, which is never changed during  $\rho$ . This proves (3).

(4) An HTTPS response  $m'$  that is accepted by  $b$  as a response to  $m$  has to be encrypted with  $k$ . The nonce  $k$  is stored by the browser in the *pendingRequests* state information. The browser only stores freshly chosen nonces there (i.e., the nonces are not used twice, or for other purposes than sending one specific request). The information cannot be altered afterwards (only deleted) and cannot be read except when the browser checks incoming messages. The nonce  $k$  is only known to  $u$  (which did not leak it to any other party prior to  $s_j$ ) and  $b$  (which did not leak it either, as  $u$  did not leak it and  $b$  is honest, see (2)). The browser  $b$  cannot send responses. This proves (4).

*Corollary 1.* In the situation of Lemma 24, as long as  $u$  does not leak the symmetric key  $k$  to  $\mathcal{W} \setminus \{u, b\}$  and the browser does not become fully corrupted,  $k$  is not known to any DY process  $p \notin \{u, b\}$  (i.e.,  $\nexists s' = (S', E') \in \rho: k \in d_{N^p}(S'(p))$ ).

*Lemma 25.* If for some  $s_i \in \rho$  an honest browser  $b$  has a document  $d$  in its state  $S_i(b).\text{windows}$  with the origin  $\langle \text{dom}, S \rangle$  where  $\text{dom} \in \text{Domain}$ , and  $S_i(b).\text{keyMapping}[\text{dom}] \equiv \text{pub}(k)$  with  $k \in \mathcal{K}$  being a private key, and there is only one DY process  $p$  that knows the private key  $k$  in all  $s_j, j \leq i$ , then  $b$  extracted (in Line 37 in Algorithm 21) the script in that document from an HTTPS response that was created by  $p$ .

PROOF. The origin of the document  $d$  is set only once: In Line 37 of Algorithm 21. The values (domain and protocol) used there stem from the information about the request (say,  $req$ ) that led to the loading of  $d$ . These values have been stored in *pendingRequests* between the request and the response actions. The contents of *pendingRequests* are indexed by freshly chosen nonces and can never be altered or overwritten (only deleted when the response to a request arrives). The information about the request  $req$  was added to *pendingRequests* in Line 76 (or Line 79 which we can exclude as we will see later) of Algorithm 22. In particular, the request was an HTTPS request iff a (symmetric) key was added to the information in *pendingRequests*. When receiving the response to  $req$ , it is checked against that information and accepted only if it is encrypted with the proper key and contains the same nonce as the request (say,  $n$ ). Only then the protocol part of the origin of the newly created document becomes  $S$ . The domain part of the origin (in our case  $\text{dom}$ ) is taken directly from the *pendingRequests* and is thus guaranteed to be unaltered.

From Line 77 of Algorithm 22 we can see that the encryption key for the request  $req$  was actually chosen using the host header of the message which will finally be the value of the origin of the document  $d$ . Since  $b$  therefore selects the public key  $S_i(b).\text{keyMapping}[\text{dom}] = S_0(b).\text{keyMapping}[\text{dom}] \equiv \text{pub}(k)$  for  $p$  (the key mapping cannot be altered during a run), we can see that  $req$  was encrypted using a public key that matches a private key which is only (if at all) known to  $p$ . With Lemma 24 we see that the symmetric encryption key for the response,  $k$ , is only known to  $b$  and the respective web server. The same holds for the nonce  $n$  that was chosen by the browser and included in the request. Thus, no other party than  $p$  can encrypt a response that is accepted by the browser  $b$  and which finally defines the script of the newly created document.

*Lemma 26.* If in a processing step  $s_i \rightarrow s_{i+1}$  of a run  $\rho$  of  $\mathcal{MS}$  an honest browser  $b$  issues an HTTP(S) request with the Origin header value  $\langle \text{dom}, S \rangle$  where  $S_i(b).\text{keyMapping}[\text{dom}] \equiv \text{pub}(k)$  with  $k \in \mathcal{K}$  being a private key, and there is only one DY process  $p$  that knows the private key  $k$  in all  $s_j$ ,  $j \leq i$ , then

- that request was initiated by a script that  $b$  extracted (in Line 37 in Algorithm 21) from an HTTPS response that was created by  $p$ , or
- that request is a redirect to a response of a request that was initiated by such a script.

**PROOF.** The browser algorithms create HTTP requests with an origin header by calling the `HTTP_SEND` function (Algorithm 17), with the origin being the fourth input parameter. This function adds the origin header only if this input parameter is not  $\perp$ .

The browser calls the `HTTP_SEND` function with an origin that is not  $\perp$  only in the following places:

- Line 51 of Algorithm 20
- Line 90 of Algorithm 20
- Line 27 of Algorithm 21

In the first two cases, the request was initiated by a script. The Origin header of the request is defined by the origin of the script's document. With Lemma 25 we see that the content of the document, in particular the script, was indeed provided by  $p$ .

In the last case (Location header redirect), as the origin is not  $\diamond$ , the condition of Line 17 of Algorithm 21 must have been true and the origin value is set to the value of the origin header of the request. In particular, this implies that an origin header does not change during redirects (unless set to  $\diamond$ ; in this case, the value stays the same in the subsequent redirects). Thus, the original request must have been created by the first two cases shown above.

The following lemma is similar to Lemma 24, but is applied to the generic HTTPS server (instead of the web browser).

*Lemma 27.* If in the processing step  $s_i \rightarrow s_{i+1}$  of a run  $\rho$  of  $\mathcal{MS}$  an honest instance  $s$  of the generic HTTPS server model (I) emits an HTTPS request of the form

$$m = \text{enc}_a(\langle \text{req}, k \rangle, \text{pub}(k'))$$

- (where  $\text{req}$  is an HTTP request,  $k$  is a nonce (symmetric key), and  $k'$  is the private key of some other DY process  $u$ ), and
- (II) in the initial state  $s_0$ , for all processes  $p \in \mathcal{W} \setminus \{u\}$ , the private key  $k'$  appears only as a public key in  $S^0(p)$ ,
- (III)  $u$  never leaks  $k'$ ,
- (IV) the instance model defined on top of the HTTPS server does not read or write the *pendingRequests* subterm of its state,
- (V) the instance model defined on top of the HTTPS server does not emit messages in `HTTPSRequests`,
- (VI) the instance model defined on top of the HTTPS server does not change the values of the *keyMapping* subterm of its state, and
- (VII) when receiving HTTPS requests of the form  $\text{enc}_a(\langle \text{req}', k_2 \rangle, \text{pub}(k'))$ ,  $u$  uses the nonce of the HTTP request  $\text{req}'$  only as nonce values of HTTPS responses encrypted with the symmetric key  $k_2$ ,
- (VIII) when receiving HTTPS requests of the form  $\text{enc}_a(\langle \text{req}', k_2 \rangle, \text{pub}(k'))$ ,  $u$  uses the symmetric key  $k_2$  only for symmetrically encrypting HTTP responses (and in particular,  $k_2$  is not part of a payload of any messages sent out by  $u$ ),

then all of the following statements are true:

- (1) There is no state of  $\mathcal{MS}$  where any party except for  $u$  knows  $k'$ , thus no one except for  $u$  can decrypt  $m$  to obtain  $\text{req}$ .
- (2) If there is a processing step  $s_j \rightarrow s_{j+1}$  where some process leaks  $k$  to  $\mathcal{W} \setminus \{u, s\}$ , there is a processing step  $s_h \rightarrow s_{h+1}$  with  $h < j$  where  $u$  leaks the symmetric key  $k$  to  $\mathcal{W} \setminus \{u, s\}$  or the process  $s$  is corrupted in  $s_j$ .
- (3) The value of the host header in  $\text{req}$  is the domain that is assigned the public key  $\text{pub}(k')$  in  $S^0(s).\text{keyMapping}$  (i.e., in the initial state of  $s$ ).
- (4) If  $s$  accepts a response (say,  $m'$ ) to  $m$  in a processing step  $s_j \rightarrow s_{j+1}$  and  $s$  is honest in  $s_j$  and  $u$  did not leak the symmetric key  $k$  to  $\mathcal{W} \setminus \{u, s\}$  prior to  $s_j$ , then  $u$  created the HTTPS response  $m'$  to the HTTPS request  $m$ , i.e., the nonce of the HTTP request  $\text{req}$  is not known to any atomic process  $p$ , except for the atomic processes  $s$  and  $u$ .

**PROOF.** (1) follows immediately from the preconditions. The proof is the same as for Lemma 24:

The process  $u$  never leaks  $k'$ , and initially, the private key  $k'$  appears only as a public key in all other process states. As the equational theory does not allow the extraction of a private key  $x$  from a public key  $\text{pub}(x)$ , the other processes can never derive  $k'$ .

Thus, even with the knowledge of all nonces (except for those of  $u$ ),  $k'$  can never be derived from any network output of  $u$ , and  $k'$  cannot be known to any other party. Thus, nobody except for  $u$  can derive  $\text{req}$  from  $m$ .

(2) We assume that some process leaks  $k$  to  $\mathcal{W} \setminus \{u, s\}$  in the processing step  $s_j \rightarrow s_{j+1}$  without  $u$  prior leaking the key  $k$  to anyone except for  $u$  and  $s$  and that the process  $s$  is not corrupted in  $s_j$ , and lead this to a contradiction.

The process  $s$  is honest in  $s_i$ .  $s$  emits HTTPS requests like  $m$  only in Line 18 of Algorithm 31:

- The message emitted in Line 3 of Algorithm 26 has a different message structure
- As  $s$  is honest, it does not send the message of Line 6 of Algorithm 31
- There is no other place in the generic HTTPS server model where messages are emitted and due to precondition (V), the application-specific model does not emit HTTPS requests. ■

The value  $k$ , which is the placeholder  $\nu_{n1}$  in Algorithm 31, is only stored in the *pendingRequests* subterm of the state of  $s$ , i.e., in  $S^{i+1}(s).pendingRequests$ . Other than that,  $s$  only accesses this value in Line 19 of Algorithm 31, where it is only used to decrypt the response in Line 20 (in particular, the key is not propagated to the application-specific model, and the key cannot be contained within the payload of an response due to (VIII)). We note that there is no other line in the model of the generic HTTPS server where this subterm is accessed and the application-specific model does not access this subterm due to precondition (IV). Hence,  $s$  does not leak  $k$  to any other party in  $s_j$  (except for  $u$  and  $s$ ). This proves (2).

(3) From Line 16 of Algorithm 31 we can see that the encryption key for the message  $m$  was chosen using the host header of the request. It is chosen from the *keyMapping* subterm of the state of  $s$ , which is never changed during  $\rho$  by the HTTPS server and never changed by the application-specific model due to precondition (VI). This proves (3).

(4)

**Response was encrypted with  $k$ .** An HTTPS response  $m'$  that is accepted by  $s$  as a response to  $m$  has to be encrypted with  $k$ : The decryption key is taken from the *pendingRequests* subterm of its state in Line 19 of Algorithm 31, where  $s$  only stores fresh nonces as keys that are added to requests as symmetric keys (see also Lines 15 and 16). The nonces (symmetric keys) are not used twice, or for other purposes than sending one specific request.

**Only  $s$  and  $u$  can create the response.** As shown previously, only  $s$  and  $u$  can derive the symmetric key (as  $s$  is honest in  $s_j$ ). Thus,  $m'$  must have been created by either  $s$  or  $u$ .

$s$  **cannot have created the response.** We assume that  $s$  emitted the message  $m'$  and lead this to a contradiction.

The generic server algorithms of  $s$  (when being honest) emit messages only in two places: In Line 3 of Algorithm 26, where a DNS request is sent, and in Line 18 of Algorithm 31, where a message with a different structure than  $m'$  is created (as  $m'$  is accepted by the server,  $m'$  must be a symmetrically encrypted ciphertext).

Thus, the instance model of  $s$  must have created the response  $m'$ .

Due to Precondition (IV), the instance model of  $s$  cannot read the *pendingRequests* subterm of its state. The symmetric key is generated freshly by the generic server algorithm in Lines 15 and 16 of Algorithm 31 and stored only in *pendingRequests*.

As the generic algorithms do not call any of the handlers with a symmetric key stored in *pendingRequests*., it follows that the instance model derived the key from a message payload in the instantiation of one of the handlers. Let  $\tilde{m}$  denote this message payload.

As the server instance model cannot derive the symmetric key without processing a message from which it can derive the symmetric key, and as the server algorithm only create the original request  $m$  as the only message with the symmetric key as a payload, it follows that  $u$  must have created  $\tilde{m}$ , as no other process can derive the symmetric key from  $m$ .

However, when receiving  $m$ ,  $u$  will use the symmetric key only as an encryption key, and in particular, will not create a message where the symmetric key is a payload (Precondition (VIII)).

Thus, the symmetric key cannot be derived by the instance of the server model, which is a contradiction to the statement that the instance model of  $s$  must have created the response  $m'$ .