# An Extensive Formal Security Analysis of the OpenID Financial-grade API

Daniel Fett
yes.com AG
mail@danielfett.de

Pedram Hosseyni
University of Stuttgart, Germany
pedram.hosseyni@sec.uni-stuttgart.de

Ralf Küsters
University of Stuttgart, Germany
ralf.kuesters@sec.uni-stuttgart.de

*Abstract*—Forced by regulations and industry demand, banks worldwide are working to open their customers' online banking accounts to third-party services via web-based APIs. By using these so-called *Open Banking* APIs, third-party companies, such as FinTechs, are able to read information about and initiate payments from their users' bank accounts. Such access to financial data and resources needs to meet particularly high security requirements to protect customers.

One of the most promising standards in this segment is the *OpenID Financial-grade API (FAPI)*, currently under development in an open process by the OpenID Foundation and backed by large industry partners. The FAPI is a profile of OAuth 2.0 designed for high-risk scenarios and aiming to be secure against very strong attackers. To achieve this level of security, the FAPI employs a range of mechanisms that have been developed to harden OAuth 2.0, such as *Code and Token Binding* (including mTLS and OAUTB), *JWS Client Assertions*, and *Proof Key for Code Exchange*.

In this paper, we perform a rigorous, systematic formal analysis of the security of the FAPI, based on an existing comprehensive model of the web infrastructure—the *Web Infrastructure Model (WIM)* proposed by Fett, Küsters, and Schmitz. To this end, we first develop a precise model of the FAPI in the WIM, including different profiles for read-only and read-write access, different flows, different types of clients, and different combinations of security features, capturing the complex interactions in a web-based environment. We then use our model of the FAPI to precisely define central security properties. In an attempt to prove these properties, we uncover partly severe attacks, breaking authentication, authorization, and session integrity properties. We develop mitigations against these attacks and finally are able to formally prove the security of a fixed version of the FAPI.

Although financial applications are high-stakes environments, this work is the first to formally analyze and, importantly, verify an Open Banking security profile.

By itself, this analysis is an important contribution to the development of the FAPI since it helps to define exact security properties and attacker models, and to avoid severe security risks before the first implementations of the standard go live.

Of independent interest, we also uncover weaknesses in the aforementioned security mechanisms for hardening OAuth 2.0. We illustrate that these mechanisms do not necessarily achieve the security properties they have been designed for.

## I. Introduction

Delivering financial services has long been a field exclusive to traditional banks. This has changed with the emergence of FinTech companies that are expected to deliver more than 20% of all financial services in 2020 [1]. Many FinTechs provide services that are based on access to a customers online banking account information or on initiating payments from a customers bank account.

For a long time, screen scraping has been the primary means of these service providers to access the customer's data at the bank. Screen scraping means that the customer enters online banking login credentials at the service provider's website, which then uses this data to log into the customer's online banking account by emulating a web browser. The service provider then retrieves account information (such as the balance or recent activities) and can trigger, for example, a cash transfer, which may require the user to enter her second-factor authentication credential (such as a TAN) at the service provider's web interface.

Screen scraping is inherently insecure: first of all, the service provider gets to know all login credentials, including the second-factor authentication of the customer. Also, screen scraping is prone to errors, for example, when the website of a bank changes.

Over the last years, the terms *API banking* and *Open Banking* have emerged to mark the introduction of standardized interfaces to financial institutions' data. These interfaces enable third parties, in particular FinTech companies, to access users' bank account information and initiate payments through well-defined APIs. All around the world, API banking is being promoted by law or by industry demand: In Europe, the *Payment Services Directive 2 (PSD2)* regulation mandates all banks to introduce Open Banking APIs by September 2019 [2]. The U.S. Department of the Treasury recommends the implementation of such APIs as well [3]. In South Korea, India, Australia, and Japan, open banking is being pushed by large financial corporations [4].

One important open banking standard currently under development for this scenario is the *OpenID Financial-grade API (FAPI)*.[1] The FAPI [5] is a profile (i.e., a set of concrete protocol flows with extensions) of the *OAuth 2.0 Authorization Framework* and the identity layer *OpenID Connect* to provide a secure authorization and authentication scheme for high-risk scenarios. The FAPI is under development at the OpenID Foundation and supported by many large corporations, such as Microsoft and the largest Japanese consulting firm, Nomura Research Institute. The OpenID Foundation is also cooperating

---

[1]In its current form, the FAPI does not (despite its name) define an API itself, but defines a security profile for the access to APIs.

with other banking standardization groups: The UK Open Banking Implementation Entity, backed by nine major UK banks, has adopted the FAPI security profile.

The basic idea behind the FAPI is as follows: The owner of the bank account (*resource owner*, also called user in what follows) visits some website or uses an app which provides some financial service. The website or app is called a *client* in the FAPI terminology. The client redirects the user to the *authorization server*, which is typically operated by the bank. The authorization server asks for the user's bank account credentials. The user is then redirected back to the client with some token. The client uses this token to obtain bank account information or initiate a payment at the *resource server*, which is typically also operated by the bank.

The FAPI aims to be secure against much stronger attackers than its foundations, OAuth 2.0 and OpenID Connect: the FAPI assumes that sensitive tokens leak to an attacker through the user's browser or operating system, and that endpoint URLs can be misconfigured. On the one hand, both assumptions are well motivated by real-world attacks and the high stakes nature of the environment where the FAPI is to be used. On the other hand, they directly break the security of OAuth 2.0 and OpenID Connect.

To provide security against such strong attackers, the FAPI employs a range of OAuth 2.0 security extensions beyond those used in plain OAuth 2.0 and OpenID Connect: the FAPI uses the so-called Proof Key for Code Exchange (PKCE)[2] extension to prevent unauthorized use of tokens. For client authentication towards the authorization server, the FAPI employs *JWS Client Assertions* or *mutual TLS*. Additionally, *OAuth token binding*[3] or *certificate-bound access tokens*[4] can be used as holder-of-key mechanisms. To introduce yet another new feature, the FAPI is the first standard to make use of the so-called JWT Secured Authorization Response Mode (JARM).

The FAPI consists of two main so-called *parts*, here also called modes, that stipulate different security profiles for read-only access to resource servers (e.g., to retrieve bank account information) and read-write access (e.g., for payment initiation). Both modes can be used by *confidential* clients, i.e., clients that can store and protect secrets (such as web servers), and by *public* clients that cannot securely store secrets, such as JavaScript browser applications. Combined with the new security features, this gives rise to many different settings and configurations in which the FAPI can run (see also Figure 3).

This, the expected wide adoption, the exceptionally strong attacker model, and the new security features make the FAPI a particularly interesting, challenging, and important subject for a detailed security analysis. While the security of (plain) OAuth 2.0 and OpenID Connect has been studied formally and informally many times before [6]–[21], there is no such analysis for the FAPI—or any other open banking API—so far. In particular, there are no results in the strong attacker

model adopted for the FAPI, and there has been no formal security analysis of the additional OAuth security mechanisms employed by the FAPI (PKCE, JWS Client Assertions, mTLS Client Authentication, OAuth Token Binding, Certificate-Bound Access Tokens, JARM), which is of practical relevance in its own right.

In this paper, we therefore study the security of the FAPI in-depth, including the OAuth security extensions. Based on a detailed formal model of the web, we formalize the FAPI with its various configurations as well as its security properties. We discover four previously unknown and severe attacks, propose fixes, and prove the security of the fixed protocol based on our formal model of the FAPI, again considering the various configurations in which the FAPI can run. Importantly, this also sheds light on new OAuth 2.0 security extensions. In detail, our contributions are as follows:

*Contributions of this Paper:* We build a **detailed formal model of the FAPI** based on a comprehensive formal model of the web infrastructure proposed by Fett et al. in [22], which we refer to as the Web Infrastructure Model (WIM). The WIM has been successfully used to find vulnerabilities in and prove the security of several web applications and standards [6], [7], [22]–[24]. It captures a wide set of web features from DNS to JavaScript in unrivaled detail and comprehensiveness. In particular, it accounts for the intricate inner workings of web browsers and their interactions with the web environment. The WIM is ideally suited to identify logical flaws in web protocols, detect a range of standard web vulnerabilities (like cross-site request forgery, session fixation, misuse of certain web browser features, etc.), and even to find new classes of web attacks.

Based on the generic descriptions of web servers in the WIM, our models for FAPI clients and authorization servers contain all important features currently proposed in the FAPI standards. This includes the flows from both parts of the FAPI, as well as the different options for client authentication, holder-of-key mechanisms, and token binding mentioned above.

Using this model of the FAPI, we define precise **security properties** for authorization, authentication, and session integrity. Roughly speaking, the authorization property requires that an attacker is unable to access the resources of another user at a bank, or act on that user's behalf towards the bank. Authentication means that an attacker is unable to log in at a client using the identity of another user. Session integrity means that an attacker is unable to force a user to be logged in at a client under the attackers identity, or force a user to access (through the client) the attacker's resources instead of the user's own resources (session fixation).

During our first attempts to prove these properties, we **discovered four unknown attacks** on the FAPI. With these attacks, adversaries can gain access to the bank account of a user, break session integrity, and, interestingly, circumvent certain OAuth security extensions, such as PKCE and Token Binding, employed by the FAPI.

We notified the OpenID FAPI Working Group of the attacks and vulnerabilities found by our analysis and are working together with them to fix the standard. To this end, we first

---

[2]Pronounced *pixie*, RFC 7636.
[3]https://tools.ietf.org/html/draft-ietf-oauth-token-binding-07
[4]https://tools.ietf.org/html/draft-ietf-oauth-mtls-11

**developed mitigations against the vulnerabilities**. We then, as another main contribution of our work and to support design decisions during the further development of the FAPI, implemented the fixes in our formal model and provided the **first formal proof of the security of the FAPI** (with our fixes applied) within our model of the FAPI, including all configurations of the FAPI and the various ways in which the new OAuth security extensions are employed in the FAPI (see Figure 3). This makes the FAPI the only open banking API to enjoy a thorough and detailed formal security analysis.

Our findings also show that (1) several **OAuth 2.0 security extensions** do not necessarily achieve the security properties they have been designed for and that (2) combining these extensions in a secure way is far from trivial. These results are relevant for all web applications and standards which employ such extensions.

*Structure of this Paper:* We first, in Section II, recall OAuth 2.0 and OpenID Connect as the foundations of the FAPI. We also introduce the new defense mechanisms that set the FAPI apart from "traditional" OAuth 2.0 and OpenID Connect flows. This sets the stage for Section III where we go into the details of the FAPI and explain its design and features. In Section IV, we present the attacks on the FAPI (and the new security mechanisms it uses), which are the results of our initial proof attempts, and also present our proposed fixes. The model of the FAPI and the analysis are outlined in Section V, along with a high-level introduction to the Web Infrastructure Model we use as the basis for our formal model and analysis of the FAPI. We conclude in Section VI. Full details and proofs are provided in the appendices.

## II. OAuth and New Defense Mechanisms

The *OpenID Financial-grade API* builds upon the OAuth 2.0 Authorization Framework [25]. Compared to the original OAuth 2.0 protocol, the FAPI aims at providing a much higher degree of security. For achieving this, the FAPI security profiles incorporate mechanisms defined in *OpenID Connect* [26] (which itself builds upon OAuth 2.0), and importantly, security extensions for OAuth 2.0 developed only recently by the IETF and the OpenID Foundation.

In the following, we give a brief overview of both OAuth 2.0 and OpenID Connect, and their security extensions used (among others) within the FAPI, namely *Proof Key for Code Exchange*, *JWS Client Assertions*, *OAuth 2.0 Mutual TLS for Client Authentication and Certificate Bound Access Tokens*, *OAuth 2.0 Token Binding* and the *JWT Secured Authorization Response Mode*. The FAPI itself is presented in Section III.

### A. Fundamentals of OAuth 2.0 and OpenID Connect

OAuth 2.0 and OpenID Connect are widely used for various authentication and authorization tasks. In what follows, we first explain OAuth 2.0 and then briefly OpenID Connect, which is based on OAuth 2.0.

*1) OAuth 2.0:* On a high level, OAuth 2.0 allows a *resource owner*, or user, to enable a *client*, a website or an application, to access her resources at some *resource server*. In order for
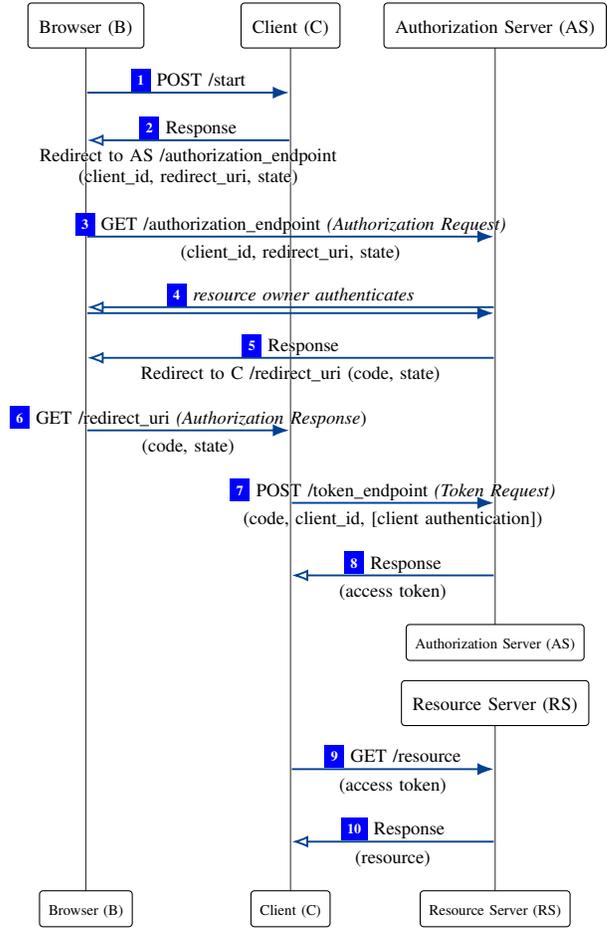


**Figure 1.** Overview of the OAuth Authorization Code Flow

the user to grant the client access to her resources, the user has to authenticate herself at an *authorization server*.

For example, in the context of the FAPI, resources include the user's account information (like balance and previous transactions) at her bank or the initiation of a payment transaction (cash transfer). The client can be a FinTech company which wants to provide a financial service to the user via access to the user's bank account. More specifically, the client might be the website of such a company (*web server client*) or the company's app on the user's device. The resource and authorization servers would typically be run by the user's bank. One client can make use of several authorization and resource servers.

RFC 6749 [25] defines multiple modes of operation for OAuth 2.0, so-called *grant types*. We here focus on the *authorization code grant* since the other grant types are not used in the FAPI.

Figure 1 shows the authorization code grant, which works as follows: The user first visits the client's website or opens the client's app on her smartphone and selects to log in or to give the client access to her resources (Step 1). The client then redirects the user to the so-called *authorization endpoint*

3

at the authorization server (AS) in Steps 2 and 3. (Endpoints are URIs used in the OAuth flow.) In this redirection, the client passes several parameters to the AS, for example, the *client id* which identifies the client at the AS, a *state* value that is used for CSRF protection,[5] a *scope* parameter (not shown in Figure 1) that describes the permissions requested by the client, and a redirection URI explained below. Note that if the client's app is used, the redirection from the app to the AS (Step 2) is done by opening the website of the AS in a browser window. The AS authenticates the user (e.g., by the user entering username and password) in Step 4 and asks for her consent to give the client access to her resources. The AS then creates a so-called *authorization code* (typically a nonce) and redirects the user back to the so-called *redirection endpoint* of the client via the user's browser in Steps 5 and 6. (If the client's app is used, a special redirect URI scheme, e.g., some-app://, is used which causes the operating system to forward the URI to the client's app.) At the AS, one or more redirection endpoints for a client are preregistered.[6] In Step 2, the client chooses one of these preregistered URIs. The authorization response (Step 5) is a redirection to this URI, with the authorization code, the state value from the request, and optionally further values appended as URI parameters.

When receiving the request resulting from the redirection in Step 6, the client first checks that the state value is the same as the one in the authorization request, typically by looking it up in the user's session with the client. If it is not the same, then the client suspects that an attacker tried to inject an authorization code into the client's session (cross-site request forgery, CSRF) and aborts the flow (see also Footnote 5). Otherwise, the client now exchanges the code for an *access token* at the so-called *token endpoint* of the AS in Steps 7 and 8. For this purpose, the client might be required to authenticate to the AS (see below). With this access token, the client can finally access the resources at the resource server (RS), as shown in Steps 9 and 10.

The RS can use different methods to check the validity of an access token presented by a client. The access token can, for example, be a document signed by the AS containing all necessary information. Often, the access token is not a structured document but a nonce. In this case, the RS uses Token Introspection [27], i.e., it sends the access token to the *introspection endpoint* of the AS and receives the information associated with the token from the AS. An RS typically has only one (fixed) AS, which means that when the RS receives an access token, it sends the introspection request to this AS.

*Public and Confidential Clients:* Depending on whether a client can keep long-term secrets, it is either called a *public* or a *confidential* client. If the client is not able to maintain

secrets, as is typically the case for applications running on end-user devices, the client is not required to authenticate itself at the token endpoint of the AS. These kinds of clients are called *public* clients. Clients able to maintain secrets, such as web server clients, must authenticate to the token endpoint (in Step 7 of Figure 1) and are called *confidential* clients.

For confidential clients, client authentication ensures that only a legitimate client can exchange the authorization code for an access token. OAuth 2.0 allows for several methods for client authentication at the token endpoint, including sending a password or proving possession of a secret [25, Section 2.3]. For public clients, other measures are available, such as PKCE (see below), to obtain a sufficient level of security.

*2) OpenID Connect:* OAuth 2.0 is built for *authorization* only, i.e., the client gets access to the resources of the user only if the user consented to this access. It does not per se provide *authentication*, i.e., proving the identity of the user to the client. This is what OpenID Connect [26] was developed for. It adds an *id token* to OAuth 2.0 which is issued by the AS and contains identity information about the end-user. ID tokens can be issued in the response from the authorization endpoint (Step 5 of Figure 1) and/or at the token endpoint (Step 8 of Figure 1). They are signed by the AS and can be bound to other parameters of the response, such as the hash of authorization codes or access tokens. Therefore, they can also be used to protect responses against modification.

### B. Proof Key for Code Exchange

The *Proof Key for Code Exchange* (PKCE) extension (RFC 7636) was initially created for OAuth public clients and independently of the FAPI. Its goal is to protect against the use of intercepted authorization codes. Before we explain how it works, we introduce the attack scenario against which PKCE should protect according to RFC 7636.

This attack starts with the leakage of the authorization code after the browser receives it in the response from the authorization endpoint (Step 5 of Figure 1). A multitude of problems can lead to a leak of the code, even if TLS is used to protect the network communication:

- On mobile operating systems, multiple apps can register themselves onto the same custom URI scheme (e.g., some-app://redirection-response). When receiving the authorization response, the operating system may forward the response (and the code) to a malicious app instead of the honest app (see [28, Section 1] and [29, Section 8.1]).
- Mix-up attacks, in which a different AS is used than the client expects (see [6] for details), can be used to leak an authorization code to a malicious server.
- As highlighted in [7], a Referer header can leak the code to an adversary.
- The code can also appear in HTTP logs that can be disclosed (accidentally) to third parties or (intentionally) to administrators.

In a setting with a public client (i.e., without client authentication at the token endpoint), an authorization code leaked to

---

[5]The state value is a nonce. The client later ensures that it receives the same nonce in the authorization response. Otherwise, an attacker could authenticate to the AS with his own identity and use the corresponding authorization response for logging in an honest user under the attacker's identity with a CSRF attack. This attack is also known as *session swapping*.

[6]Without preregistration, a malicious client starting a login flow with the client id of an honest client could receive a code associated with the honest client.

the attacker can be redeemed directly by the attacker at the authorization server to obtain an access token.

RFC 7636 aims to protect against such attacks even if not only the authorization response leaks but also the authorization request as well. Such leaks can happen, for example, from HTTP logs (Precondition 4b of Section 1 of RFC 7636) or unencrypted HTTP connections.

PKCE works as follows: Before sending the authorization request, the client creates a random value called *code verifier*. The client then creates the *code challenge* by hashing the verifier[7] and includes the challenge in the authorization request (Step [2] of Figure 1). The AS associates the generated authorization code with this challenge. Now, when the client redeems the code in the request to the token endpoint (Step [7] of Figure 1), it includes the code verifier in the token request. This message is sent directly to the AS and protected by TLS, which means that the verifier cannot be intercepted. The idea is that if the authorization code leaked to the attacker, the attacker still cannot redeem the code to obtain the access token since he does not know the code verifier.

### C. Client Authentication using JWS Client Assertions

As mentioned above, the goal of client authentication is to bind an authorization code to a certain confidential client such that only this client can redeem the code at the AS. One method for client authentication is the use of JWS Client Assertions [26, Section 9], which requires proving possession of a key instead of sending a password directly to the authorization server, as in plain OAuth 2.0.

To this end, the client first generates a short document containing its client identifier and the URI of the token endpoint. Now, depending on whether the *client secret* is a private (asymmetric) or a symmetric key, the client either signs or MACs this document. It is then appended to the token request (Step [7] of Figure 1). As the document contains the URI of the receiver, attacks in which the attacker tricks the client into using a wrong URI are prevented, as the attacker cannot reuse the document for the real endpoint (cf. Section III-C4). Technically, the short document is encoded as a JSON Web Token (JWT) [30] to which its signature/MAC is attached to create a so-called JSON Web Signature (JWS) [31].

### D. OAuth 2.0 Mutual TLS

*OAuth 2.0 Mutual TLS for Client Authentication and Certificate Bound Access Tokens* (mTLS) [32] provides a method for both client authentication and token binding.

OAuth 2.0 Mutual TLS Client Authentication makes use of *TLS client authentication*[8] at the token endpoint (in Step [7] 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. To this end, the client proves

that it knows the private key belonging to a certificate that is either (a) self-signed and preconfigured 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 [9] 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.[9]

### E. OAuth 2.0 Token Binding

*OAuth 2.0 Token Binding* (OAUTB) [33] is used to bind access tokens and/or authorization codes to certain TLS connections. It is based on the *Token Binding* protocol [34]–[37] and can be used with all TLS versions. In the following, we first sketch token binding in general before we explain OAuth 2.0 Token Binding.

*1) Basics:* For simplicity of presentation, in the following, we assume that a browser connects to a web server. The protocol remains the same if the browser is replaced by another server. (In the context of OAuth 2.0, in some settings in fact the client takes the role of the browser as explained below.)

At its core, token binding works as follows: When a web server indicates (during TLS connection establishment) that it wants to use token binding, the browser making the HTTP request over this TLS connection creates a public/private key pair for the web server's origin. It then sends the public key to the server and proves possession of the private key by using it to create a signature over a value unique to the current TLS connection. Since the browser re-uses the same key pair for future connections to the same origin, the web server will be able to unambiguously recognize the browser in future visits.

Central for the security of token binding is that the private key remains secret inside the browser. To prevent replay attacks, the browser has to prove possession of the private key by signing a value that is unique for each TLS session. To this end, token binding uses the *Exported Keying Material* (EKM) of the TLS connection, a value derived from data of the TLS handshake between the two participants, as specified in [37]. As long as at least one party follows the protocol, the EKM will be unique for each TLS connection.

We can now illustrate the usage of token binding in the context of a simplified protocol in which a browser $B$ requests a token from a server $S$: First, $B$ initiates a TLS connection to $S$, where $B$ and $S$ use TLS extensions [35] to negotiate the use of token binding and technical details thereof. Browser $B$ then creates a public/private key pair $(k_{B,S}, k'_{B,S})$ for the origin of $S$, unless such a key pair exists already. The public key $k_{B,S}$ (together with technical details about the key, such as its bit length) is called *Token Binding ID* (for the specific origin).

---

[7]If it is assumed that the authorization request never leaks to the attacker, it is sufficient and allowed by RFC 7636 to use the verifier as the challenge, i.e., without hashing.

[8]As noted in [32], Section 5.1 this extension supports all TLS versions with certificate-based client authentication.

[9]As mentioned above, the RS can read this information either directly from the access token if it is a signed document, or uses token introspection to retrieve the data from the AS.

When sending the first HTTP request over the established TLS connection, $B$ includes in an HTTP header the so-called *Token Binding Message*:

$$\text{TB-Msg}[k_{B,S}, \text{sig}(EKM, k'_{B,S})] \tag{1}$$

It contains both the Token Binding ID (i.e., essentially $k_{B,S}$) and the signed EKM value from the TLS connection, as specified in [38]. The server $S$ checks the signature using $k_{B,S}$ as included in this message and then creates a token and associates it with the Token Binding ID as the unique identifier of the browser.

When $B$ wants to redeem the token in a new TLS connection to $S$, $B$ creates a new Token Binding Message using the same Token Binding ID, but signs the new EKM value:

$$\text{TB-Msg}[k_{B,S}, \text{sig}(\overline{EKM}, k'_{B,S})] \tag{2}$$

As the EKM values are unique to each TLS connection, $S$ concludes that the sender of the message knows the private key of the Token Binding ID, and as the sender used the same Token Binding ID as before, the same party that requested the token in the first request is using it now.

The above describes the simple situation that $B$ wants to redeem the token received from $S$ again at $S$, i.e., from the same origin. In this case, we call the token binding message in (1) a *provided* token binding message. If $B$ wants to redeem the token received from $S$ at another origin, say at $C$, then instead of just sending the provided token message in (1), $B$ would in addition also send the so-called *referred* token binding message, i.e., instead of (1) B would send

$$\begin{aligned}
&\text{TB-prov-Msg}[k_{B,S}, \text{sig}(EKM, k'_{B,S})], \\
&\text{TB-ref-Msg}[k_{B,C}, \text{sig}(EKM, k'_{B,C})].
\end{aligned} \tag{3}$$

Note that the EKM is the same in both messages, namely the EKM value of the TLS connection between $B$ and $S$ (rather than between $B$ and $C$, which has not happened yet anyway). Later when $B$ wants to redeem the token at $C$, $B$ would use $k_{B,C}$ in its (provided) token message to $C$.

*2) Token Binding for OAuth:* In the following, we explain how token binding is used in OAuth in the case of app clients. The case of web server clients is discussed below.

The flow is shown in Figure 2. Note that in this case, token binding is used between the OAuth client and the authorization and resource servers; the browser in Figure 1 is not involved.

The client has two token binding key pairs, one for the AS and one for the RS (if these key pairs do not already exist, the client creates them during the flow). When sending the authorization request (Step 2 of Figure 2), the client includes the hash of the Token Binding ID it uses for the AS as a PKCE challenge (cf. Section II-B). When exchanging the code for an access token in Step 7, the client proves possession of the private key of this Token Binding ID, and the AS only accepts the request when the hash of the Token Binding ID is the same as the PKCE challenge. Therefore, the code can only be exchanged by the participant that created the authorization request. Note that for this purpose the AS only takes the *provided* token binding message sent to the AS in Step 7



**Figure 2.** OAUTB for App Clients

into account. However, the AS also checks the validity of the *referred* token binding message (using the same EKM value) and associates $k_{C,RS}$ with the token issued by the AS in Step 8.

The token binding ID $k_{C,RS}$ is used in Step 9 by the client to redeem the token at the RS. The RS then checks if this is the same token binding ID that is associated with the access token. This information can be contained in the access token if it is structured and readable by the RS or via token introspection.

Altogether, Token Binding for OAuth (in the case of app clients) is supposed to bind both the authorization code and the access token to the client. That is, only the client who initiated the flow (in Step 2) can redeem the authorization code at the AS and the corresponding access token at the RS, and hence, get access to the resource at the RS.

*3) Binding Authorization Codes for Web Server Clients:* In the case that the client is a web server, the binding of the authorization code to the client is already done by client authentication, as a web server client is always confidential (cf. Section II-A1). Therefore, the client does *not* include the hash

of a Token Binding ID in the authorization request (Step [2] of Figure 2). Instead, the mechanism defined in OAUTB aims at binding the authorization code to the browser/client pair. (The binding of the access token to the client is done in the same way as for an app client).

More precisely, for web server clients, the authorization code is bound to the token binding ID that the browser uses for the client. For this purpose, the client includes an additional HTTP header in the first response to the browser (Step [2] of Figure 2), which signals the browser that it should give the token binding ID it uses for the client to the authorization server. When sending the authorization request to the authorization server in Step [3], the browser thus includes a provided and a referred token binding message, where the referred message contains the token binding ID, that the browser later uses for the client (say, $k_{B,C}$). When generating the authorization code, the authorization server associates the code with $k_{B,C}$.

When redirecting the code to the client in Step [6], the browser includes a token binding message for $k_{B,C}$, thereby proving possession of the private key.

When sending the token request in Step [7], the client includes $k_{B,C}$. We highlight that the client does not send a token binding message for $k_{B,C}$ since the client does not know the corresponding private key (only the browser does).

The authorization server checks if this key is the same token binding ID it associated the authorization code with, and therefore, can check if the code was redirected to the client by the same browser that made the authorization request. In other words, by this the authorization code is bound to the browser/client pair.

### F. JWT Secured Authorization Response Mode

The recently developed *JWT Secured Authorization Response Mode* (JARM) [39] aims at protecting the OAuth authorization response (Step [5] of Figure 1) by having the AS sign (and optionally encrypt) the response. The authorization response is then encoded as a JWT (see Section II-C). The JARM extension can be used with any OAuth 2.0 flow.

In addition to the regular parameters of the authorization response, the JWT also contains its issuer (identifying the AS) and its audience (client id). For example, if combined with the Authorization Code Flow, the response JWT contains the issuer, audience, authorization code, and state values.

By using JARM, the authorization response is integrity protected and injection of leaked authorization codes is prevented.

### III. THE OPENID FINANCIAL-GRADE API

The OpenID Financial-grade API [5] currently comprises two implementer's drafts. One defines a profile for read-only access, the other one for read-write access. Building on Section II, here we describe both profiles and the various configurations in which these profiles can run (see Figure 3). Furthermore, we explain the assumptions made within the FAPI standard and the underlying OAuth 2.0 extensions.
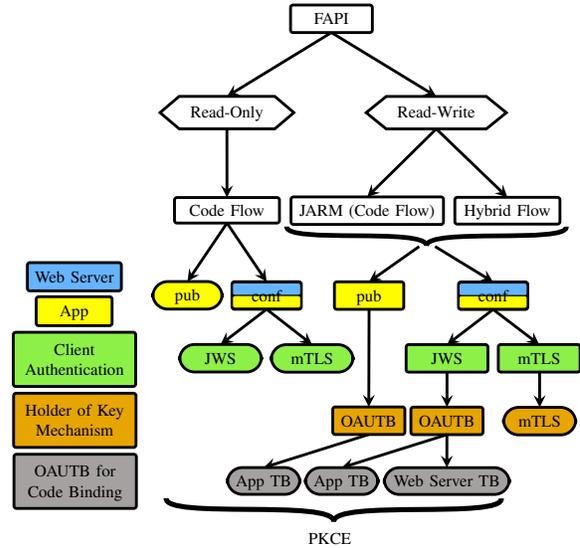


**Figure 3.** Overview of the FAPI. One path (terminated by a box with rounded corners) describes one possible configuration of the FAPI. The paths marked with PKCE use PKCE. JARM and Hybrid flows both allow for the configurations shown.

### A. Financial-grade API: Read-Only Profile

In the following, we explain the Read-Only flow as described in [40]. The Read-Only profile aims at providing a secure way for accessing data that needs a higher degree of protection than regular OAuth, e.g., for read access to financial data.

The Read-Only flow is essentially an **OAuth Authorization Code flow** (cf. Section II). Additionally, the client can request an ID Token (see Section II-A2) from the token endpoint by adding a *scope* parameter to the authorization request (Step [2] of Figure 1) with the value openid.

In contrast to regular OAuth and OpenID Connect, the client is required to have a different set of **redirection URIs** for each authorization server. This separation prevents mix-up attacks, where the authorization response (Step [6] in Figure 1) comes from a different AS than the client expects (see [6] and [41] for more details on mix-up attacks). When receiving the authorization response, the client checks if the response was received at the redirection URI specified in the authorization request (Step [2] in Figure 1).

One of the main additions to the regular OAuth flow is the use of **PKCE** as explained in Section II-B. The PKCE challenge is created by hashing a nonce.

The FAPI furthermore requires **confidential clients to authenticate** at the token endpoint (in Step [7] of Figure 1) using either *JWS Client Assertions* (cf. Section II-C) or *Mutual TLS* (cf. Section II-D). Public clients do not use client authentication.

### B. Financial-grade API: Read-Write Profile

The Read-Write profile [42] aims at being secure under stronger assumptions than the Read-Only profile, in order to

be suitable for scenarios such as write access to financial data. The full set of assumptions is described in Section III-C.

The flow can be either an **OpenID Connect (OIDC) Hybrid flow**, which means that both the authorization response (Step 5 in Figure 1) and the token response (Step 8 in Figure 1) contain an id token (see Section II-A2), or any other OAuth-based flow used together with **JARM** (see Section II-F). When using the Hybrid flow, the FAPI profile also requires that the hash of the state value is included in the first id token.

In addition to the parameters of the Read-Only flow, the authorization request prepared by the client (Step 2 of Figure 1) is required to contain a **request JWS**, which is a JWT, signed by the client, containing all request parameters together with the audience of the request (cf. Section II-C).

One of the main security features of the profile is the **binding of the authorization code and the access token** to the client, which is achieved by using either mTLS (cf. Section II-D) or OAUTB (OAuth 2.0 Token Binding, see Section II-E). A public client is required to use OAUTB, while a confidential client can use either OAUTB or mTLS.

If the client is a confidential client using mTLS, the request does not contain a PKCE challenge. When using OAUTB, the client uses a **variant of PKCE**, depending on whether the client is a web server client or an app client (cf. Section II-E).

In the case of a confidential client, the **client authentication at the token endpoint** is done in the same way as for the Read-Only flow, i.e., by using either JWS Client Assertions (cf. Section II-C) or Mutual TLS (cf. Section II-D).

### C. Overview of Assumptions and Mitigations

In the following, we explain the conditions under which the FAPI profiles and the OAuth extensions aim to be secure according to their specifications.

*1) Leak of Authorization Response:* As described in Section II-B in the context of PKCE, there are several scenarios in which the authorization response (Step 6 of Figure 1), and hence, the authorization code, can leak to the attacker (in clear), in particular in the case of app clients. In our model of the FAPI, we therefore assume that the authorization response is given to the attacker if the client is an app. At first glance, leakage of the authorization code is indeed mitigated by the use of PKCE since an attacker does not know the code verifier, and hence, cannot redeem the code at the AS. However, our attack described in Section IV-C shows that the protection provided by PKCE can be circumvented.

*2) Leak of Authorization Request:* The Read-Only profile of the FAPI explicitly states that the PKCE challenge should be created by hashing the verifier. The use of hashing should protect the PKCE challenge even if the authorization request leaks (e.g., by leaking HTTP logs, cf. Section II-B), and therefore, we assume in our model that the authorization request (Step 2 of Figure 1) leaks to the attacker.

*3) Leak of Access Token:* In the Read-Write profile, it is assumed that the access token might leak due to phishing [42, Section 8.3.5]. In our model, we therefore assume that the access token might leak in Step 5 of Figure 1. This problem is seemingly mitigated by using either mTLS or OAUTB, which bind the access token to the legitimate client, and hence, only the legitimate client should be able to redeem the access token at the RS even if the access token leaked. The FAPI specification states: "When the FAPI client uses MTLS or OAUTB, the access token is bound to the TLS channel, it is access token phishing resistant as the phished access tokens cannot be used." [42, Section 8.3.5]. However, our attack presented in Section IV-A shows that this is not the case.

*4) Misconfigured Token Endpoint:* An explicit design decision by the FAPI working group was to make the Read-Write profile secure even if the token request (Step 7 of Figure 1) leaks. The FAPI specification describes this attack as follows: "In this attack, the client developer is social engineered into believing that the token endpoint has changed to the URL that is controlled by the attacker. As the result, the client sends the code and the client secret to the attacker, which will be replayed subsequently." [42, Section 8.3.2].

Therefore, we make this assumption also in our FAPI model. Seemingly, this problem is mitigated by code binding through client authentication or OAUTB, which means that the attacker cannot use the stolen code at the legitimate token endpoint. "When the FAPI client uses MTLS or OAUTB, the authorization code is bound to the TLS channel, any phished client credentials and authorization codes submitted to the token endpoint cannot be used since the authorization code is bound to a particular TLS channel." [42, Section 8.3.2]. Note that in the FAPI the client does not authenticate by using the client secret as a password, but by proving possession (either using JWS Client Assertions or mTLS), which means that the attacker cannot reuse credentials.

However, our attack presented in Section IV-B shows that this intuition is misleading.

## IV. ATTACKS

As already mentioned in the introduction, in Section V we present our rigorous formal analysis of the FAPI based on the Web Infrastructure Model. Through this formal analysis of the FAPI with the various OAuth 2.0 extensions it uses, we not only found attacks on the FAPI but also on some of the OAuth 2.0 extensions, showing that (1) these extensions do not achieve the security properties they have been designed for and (2) that combining these extensions in a secure way is far from trivial. Along with the attacks, we also propose fixes to the standards. Our formal analysis presented in Section V considers the fixed versions.

We start by describing two attacks on Token Binding, followed by an attack on PKCE, and one vulnerability hidden in the assumptions of PKCE.

We emphasize that our attacks work even if all communication uses TLS and even if the attacker is merely a web attacker, i.e., does not control the network but only certain parties.

As already mentioned in the introduction, we notified the OpenID FAPI Working Group of the attacks found by our analysis and are working together with them to fix the standard.

## A. Cuckoo's Token Attack

As explained in Section III-C3, the Read-Write profile of the FAPI aims at providing security even if the attacker obtains an access token, e.g., due to phishing. Intuitively, this protection seems to be achieved by binding the access token to the client via mTLS (see Section II-D) or OAUTB (see Section II-E).

However, these mechanisms prevent the attacker only from directly using the access token in the same flow. As illustrated next, in a second flow, the attacker *can* inject the bound access token and let the client (to which the token is bound) use this token, which enables the attacker to access resources belonging to an honest identity.

This attack affects all configurations of the Read-Write profile (see Figure 3). Also, the Read-Only profile is vulnerable to this attack; this profile is, however, not meant to defend against stolen access tokens.

We note that the underlying principle of the attack should be relevant to other use-cases of token binding as well, i.e., whenever a token is bound to a participant, the involuntary use of a leaked token (by the participant to which the token is bound) should be prevented.
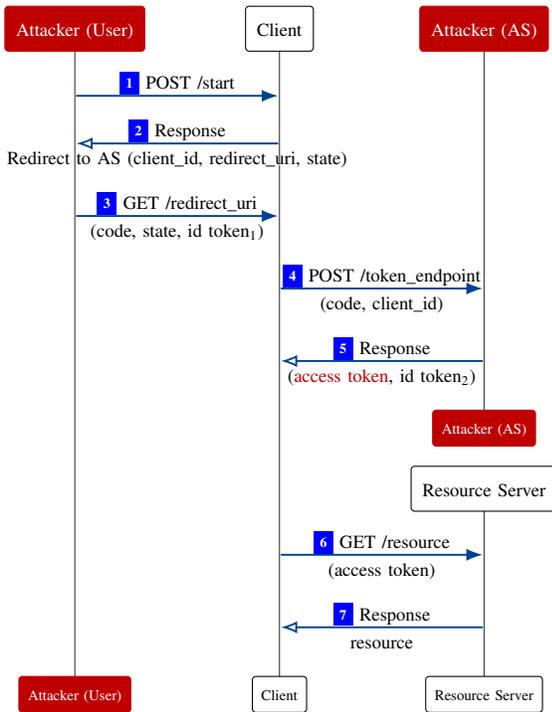


**Figure 4.** Cuckoo's Token Attack

Figure 4 depicts the attack for the OIDC Hybrid Flow, i.e., when both responses of the AS contain id tokens (see Section III-B). The attack works analogously for the code flow in combination with JARM (see Section III-B).

As explained, we assume that the attacker already obtained (phished) an access token issued by an honest AS to an honest client for accessing resources of an honest user. We also assume that the honest client supports the use of several ASs (a common setting in practice, as already mentioned in Section II), where in this case one of the ASs is dishonest.[10]

First, the attacker starts the flow at the client and chooses his own AS. Since he is redirected to his own AS in Step 2, he can skip the user authentication step and return an authorization response immediately. Apart from that, the flow continues normally until Step 4, where the client sends the code to the attacker AS. In Step 5, the attacker AS returns the previously phished access token together with the second id token.

Until here, all checks done by the client pass successfully, as the attacker AS adheres to the protocol. The only difference to an honest authorization server is that the attacker AS returns a phished access token. In Step 6, the resource server receives the (phished) access token and provides the client access to the honest resource owner's resources for the phished access token,[11] which implies that now the attacker has access to these resources through the client.

To prevent the use of leaked access tokens, the client should include, in the request to the RS, the identity of the AS the client received the access token from. The client can take this value from the second id token. Now, the RS would only continue the flow if its belief is consistent with the one of the RS. We apply an analogous fix for flows with JARM. These fixes are included in our model and shown to work in Section V.

## B. Access Token Injection with ID Token Replay

As described in Section III-C3, the Read-Write profile aims to be secure if an attacker acquires an access token for an honest user. The profile also aims to be secure even if the token endpoint URI is changed to an attacker-controlled URI (see Section III-C4). Now, interestingly, these two threat scenarios combined in this order are the base for the attack described in the following. In this attack, the attacker returns an access token at the misconfigured token endpoint. While the attack looks similar to the previous attack at first glance, here the attacker first interacts with the *honest* AS and later *replays an id token* at the token endpoint. Both attacks necessitate different fixes. The outcome, however, is the same, and, just as the previous attack, this attack affects all configurations of the Read-Write profile, even if JARM is used. We explain the attack using the Hybrid Flow.

Figure 5 shows how the attack proceeds. The attacker initiates the Read-Write flow at the client and follows the

---

[10]We highlight that we do not assume that the attacker controls the AS that issued the access token (i.e., the AS at which the honest user is registered). This means that the (honest) user uses an honest client and an honest authorization server.

[11]Which RS is used in combination with an AS depends on the configuration of the client, which is acquired through means not defined in OAuth. Especially in scenarios where this configuration is done dynamically, a dishonest AS might be used in combination with an honest RS. But also if the client is configured manually, as is often the case today, it might be misconfigured or social engineered into using specific endpoints. Recall from Section II-A that the access token might be a document signed by the (honest) AS containing all information the RS needs to process the access token. Alternatively, and more common, the RS performs token introspection, if the access token is just a nonce. The RS typically uses only one AS (in this case, the honest AS) to which it will send the introspection request.
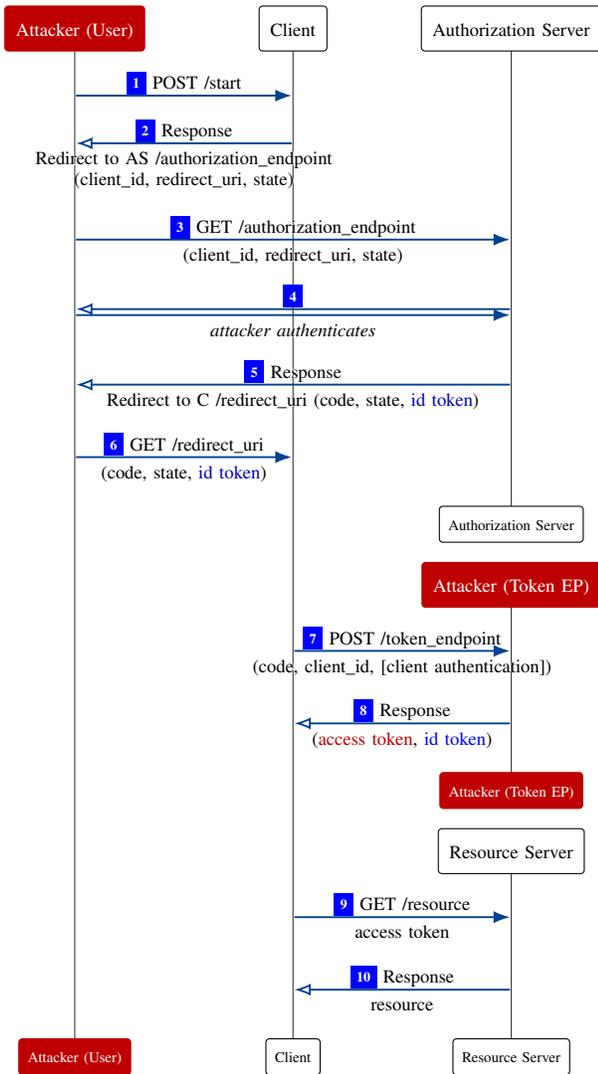
**Figure 5.** Access Token Injection with ID Token Replay Attack

regular flow until Step ⑥. As the authorization response was created by the honest AS, the state and all values of the id token are correct and the client accepts the authorization response.

In Step ⑦, the client sends the token request to the misconfigured token endpoint controlled by the attacker. The value of the code and the checks regarding client authentication and proof of possession of keys are not relevant for the attacker.

In Step ⑧, the attacker sends the token response containing the phished access token. As the flow is an OIDC Hybrid Flow, the attacker is required to return an id token. Here, he returns the same id token that he received in Step ⑤, which is signed by the honest AS. The client is required to ensure that both id tokens have the same subject and issuer values, which in this case holds true since they are identical.

The client sends the access token to the honest resource server, by which the attacker gets read-write access to the resource of the honest resource owner through the client.

As we show in our security analysis (see Section V), this

scenario is prevented if the second id token is required to contain the hash of the access token that is returned to the client, as the attacker cannot create id tokens with a valid signature of the AS. A similar fix also works for flows with JARM. The fixes are already included in our model.

### C. PKCE Chosen Challenge Attack

As detailed in Section III-C1, the FAPI uses PKCE in order to protect against leaked authorization codes. This is particularly important for public clients as these clients, unlike confidential ones, do not authenticate to an AS when trying to exchange the code for an access token.

Recall that the idea of PKCE is that a client creates a PKCE challenge (hash of a nonce), gives it to the AS, and when redeeming the authorization code at the AS, the client has to present the correct PKCE verifier (the nonce). This idea works when just considering an honest flow in which the code leaks to the attacker, who does not know the PKCE verifier. However, our attack shows that the protection can be circumvented by an attacker who pretends to be an honest client.

This attack affects public clients who use the Read-Only profile of the FAPI. It works as follows (see Figure 6): As in RFC 7636, two apps are installed on a user's device, an honest app and a malicious app. The honest app is a client of an honest AS with the client identifier *hon_client_id* and the redirection URI *hon_redir_uri*. The malicious app is not registered at the AS.

The Read-Only flow starts at the malicious app, which prompts the user to log in. Now, the malicious app prepares an authorization request containing the client id and a redirect URI of the honest client (Step ②). At this point, the malicious app also creates a PKCE verifier and includes the corresponding challenge in the authorization request.

The flow continues until the browser receives the authorization response in Step ⑤. As the redirection URIs are preregistered at the AS, the redirection URI in the authorization request was chosen from the set of redirect URIs of the honest app, and therefore, the authorization response is redirected to the honest client after the browser receives it.

As described in Sections II-B and III-C1, at this point, the authorization response with the authorization code might leak to the attacker (Step ⑥). The malicious app is now able to exchange the code (associated with the honest client) at the token endpoint in Steps ⑦ and ⑧, as it knows the correct PKCE verifier and, as the honest app is a public client, without authenticating to the AS.

To prevent this scenario, an honest AS must ensure that the PKCE challenge was created by the client with the id *hon_client_id*. To achieve this, for public clients in the Read-Only flow we use the same mechanism that the FAPI uses for public clients in the Read-Write flow, namely the authorization request should contain a signed JWT (see also Section II-C, although JWTs are now used in a different way). This ensures that the client stated in the request actually made the request, and hence, no other client should know the PKCE verifier. Note that by using signed JWTs for public clients the FAPI assumes
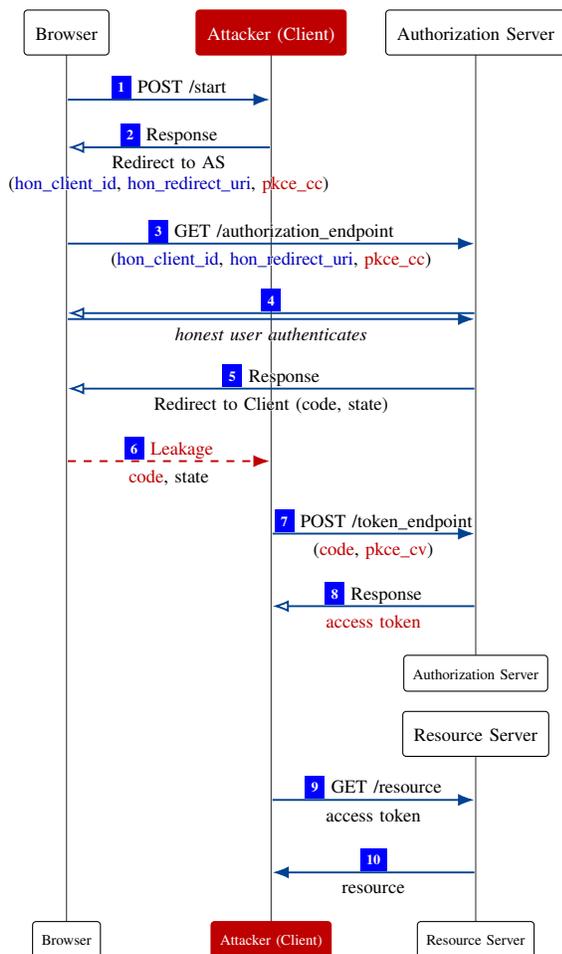
**Figure 6.** PKCE Chosen Challenge Attack

that public clients can store some secrets (which might, for example, be protected by user passwords). Our fix is already included in the model and our analysis (Section V) shows that it works.

### D. Authorization Request Leak Attacks

As explained in Section III-C2, the PKCE challenge is created such that PKCE is supposed to work even if the authorization request leaks (see also Section II-B).

However, if a leak of the authorization request occurs not only the PKCE challenge leaks to the attacker but also the state value, since both values are contained in the authorization request. Our attack shows that an attacker who knows the state value can circumvent the CSRF protection the state value was supposed to provide. As a result of the attack, the honest user is logged in under the identity of the attacker and uses the resources of the attacker, which breaks session integrity. The details of this attack are presented in Appendix A.

This is a well-known class of attacks for plain OAuth flows [43], but it is important to highlight that the protections designed into the FAPI do not sufficiently protect most flows

against such attacks, even though PKCE explicitly foresees the attack vector.

To prevent this attack, one essentially has to prevent CSRF forgery in this context. However, this is non-trivial because of the very strong attacker model considered by the OpenID FAPI Working Group: leaks and misconfigurations are assumed to occur at various places. As further explained in Appendix A, just assuming that the authorization request does not leak to the attacker would not fix the problem in general; one at least would have to assume that the authorization response does not leak either. Making these assumptions, however, of course contradicts the OpenID FAPI Working Group's intention, namely providing security even in the presence of very strong attackers.

Fortunately, we can prove that regular FAPI web server clients which use OAUTB are not vulnerable to this attack even in the presence of the strong attackers assumed by the OpenID FAPI Working Group and throughout this paper. More specifically, we can prove session integrity of the FAPI for such clients (and strong attackers), which in particular excludes the above attack (see Section V). For all other types of clients, our attack works, and there does not seem to be a fix which would not massively change the flows, and hence, the standards, as argued in Appendix A. In this sense, our results for session integrity appear to be the best we can obtain for the FAPI.

## V. FORMAL SECURITY ANALYSIS

In this section, we present our formal analysis of the FAPI. We start by very briefly recalling the Web Infrastructure Model (WIM), followed by a sketch of our formal model of the FAPI, which as already mentioned uses the WIM as its basic web infrastructure model. We then introduce central security properties the FAPI is supposed to satisfy, along with our main theorem stating that these properties are satisfied.

Since we cannot present the full formal details here, we provide the complete analysis in Appendices G–L. This includes the precise formalization of clients, authorization servers, and resource servers, as well as full detailed proofs.

### A. The Web Infrastructure Model

The Web Infrastructure Model (WIM) was introduced by Fett, Küsters, and Schmitz in [22] (therefore also called the FKS model) and further developed in subsequent work. The appendix of [44] provides a detailed description of the model; a comparison with other models and a discussion of its scope and limitations can be found in [22]–[24]. We here only give a brief overview of the WIM following the description in [7], with some more details presented in Appendix B. As explained there, we slightly extend the WIM, among others to model OAUTB. We choose the WIM for our work because, as mentioned in the introduction, the WIM is the most comprehensive model of the web infrastructure to date.

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. Among

others, HTTP(S) requests and responses,[12] including several headers, such as cookie, location, referer, authorization, strict transport security (STS), and origin headers, are modeled. The model of web browsers captures the concepts of windows, documents, and iframes, including the complex navigation rules, as well as modern technologies, such as web storage, web messaging (via postMessage), and referrer policies. JavaScript is modeled in an abstract way by so-called *scripts* which can be sent around and, among others, can create iframes, access other windows, and initiate XMLHttpRequests.

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. 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. The WIM follows the Dolev-Yao approach (see, e.g., [45]). That is, messages are expressed as formal terms over a signature $\Sigma$ which contains constants (for addresses, strings, nonces) as well as sequence, projection, and function symbols (e.g., for encryption/decryption and signatures).

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 computed (formally, derived in the usual Dolev-Yao style) from the input event and the state.

The so-called *attacker process* 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 Dolev-Yao process could possibly perform. Attackers can corrupt other parties, browsers, and servers.

A *script* models JavaScript running in a browser. Scripts are defined similarly to Dolev-Yao processes, but run in and interact with the browser. Similar to an attacker process, an *attacker script* can (non-deterministically) perform every action a script can possibly perform within a browser.

A *system* is a set of processes. A *configuration* of a system is a tuple of the form $(S, E, N)$ where $S$ maps every process of the system to its state, $E$ is the pool of waiting events, and $N$ is a sequence of unused nonces. In what follows, $s_0^p$ denotes the initial state of process $p$. 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.

---

[12]We note that the WIM models TLS at a high level of abstraction such that messages are exchanged in a secure way.

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).

In our FAPI model, we need to specify only the behavior of servers and scripts. These are not defined by the WIM since they depend on the specific application, unless they become corrupted, in which case they behave like attacker processes and attacker scripts. We assume the presence of a strong network attacker which also controls all DNS servers (but we assume a working PKI).

## B. Sketch of the Formal FAPI Model

A **FAPI web system (with a network attacker)**, denoted by $\mathcal{FAPI}$, is a web system (as explained in Section V-A) and can contain an unbounded finite number of clients, authorization servers, resource servers, browsers, and a network attacker. Note that a network attacker is the most powerful attacker, which subsumes all other attackers. Except for the attacker, all processes are initially honest and can become (dynamically) corrupted by the attacker at any time.

In a FAPI web system, clients, authorization servers, and resource servers act according to the specification of the FAPI presented in Section III. (As mentioned in Section V-A, the behavior of browsers is fixed by the standards. Their modeling is independent of the FAPI and already contained in the WIM.) Our models for clients and servers follow the latest recommendations regarding the security of OAuth 2.0 [41] to mitigate all previously known attacks. The model also contains the fixes pointed out in Section IV, as otherwise, we would not be able to prove the desired security properties (see below).

The primary goal of the FAPI is to provide a high degree of security. Its flows are intended to be secure even if information leaks to an attacker. As already outlined in Section III-C, we model this by sending the authorization response (in the case of an app client), the access token (in the case of a Read-Write flow), and the authorization request to an arbitrary (non-deterministically chosen) IP address. Furthermore, in the Read-Write profile, the token request can be sent to an arbitrary URI.

Importantly, one FAPI web system contains all possible settings in which the FAPI can run, as depicted in Figure 3, in particular, we consider all OAuth 2.0 extensions employed in the FAPI. More precisely, every client in a FAPI web system runs one of the possible configurations (i.e., it implements on one path in Figure 3). Different clients may implement different configurations. Every authorization and resource server in a FAPI web system supports all configurations at once. When interacting with a specific client, a server just chooses the configuration the client supports. In our model, the various endpoints (authorization, redirection, token), the information which client supports which FAPI configuration,

client credentials, etc. are preconfigured and contained in the initial states of the processes. How this information is acquired is out of the scope of the FAPI.

We emphasize that when proving security properties of the FAPI, we prove these properties for all FAPI web systems, where different FAPI web systems can differ in the number of clients and servers, and their preconfigured information.

Furthermore, we note that there is no notion of time in the WIM, hence, tokens do not expire. This is a safe overapproximation as it gives the attacker more power.

To give a feel for our formal FAPI model, an excerpt of the model is provided in Appendix C.

## C. Security Properties and Main Theorem

In the following, we define the security properties the FAPI should fulfill, namely authorization, authentication, and session integrity. These properties have been central to also OAuth 2.0 and OpenID Connect [6], [7]. But as mentioned, the FAPI has been designed to fulfill these properties under stronger adversaries, therefore using various OAuth extensions. While our formulations of these properties are inspired by those for OAuth 2.0 and OpenID Connect, they had to be adapted and extended for the FAPI, e.g., to capture properties of resource servers, which previously have not been modeled. We also state our main theorem.

We give an overview of each security property. For the authorization property, we provide an in-depth explanation, together with the formal definition. Appendix D contains a proof sketch for the authorization property. Full details and proofs of all properties are given in Appendix K.

*1) Authorization:* Informally speaking, for authorization we require that an attacker cannot access resources belonging to an honest user (browser). A bit more precise, we require that in all runs $\rho$ of a FAPI web system $\mathcal{FAPI}$ if an honest resource server receives an access token that is associated with an honest client, an honest authorization server, and an identity of an honest user, then access to the corresponding resource is not provided to the attacker in any way. We highlight that this does not only mean that the attacker cannot access the resource directly at the resource server, but also that the attacker cannot access the resource through a client.

In order to formalize this property, we first need to define what it means for an access token to be associated with a client, an AS, and a user identity (see below for an explanation of this definition).

*Definition 1 (Access Token associated with C, AS and ID).* Let $c$ be a client with client id *clientId* issued to $c$ by the authorization server *as*, and let $id \in \mathsf{ID}^{as}$, where $\mathsf{ID}^{as}$ denotes the set of identities governed by as. We say that an *access token t is associated with c, as and id* in state $S$ of the configuration $(S,E,N)$ of a run $\rho$ of a FAPI web system, if there is a sequence $s \in S(as).\mathtt{accessTokens}$ such that $s \equiv \langle id, clientId, t, \mathtt{r} \rangle$, $s \equiv \langle \mathtt{MTLS}, id, clientId, t, key, \mathtt{rw} \rangle$ or $s \equiv \langle \mathtt{OAUTB}, id, clientId, t, key', \mathtt{rw} \rangle$, for some *key* and *key'*.

Intuitively, an access token $t$ is associated with a client $c$, authorization server *as*, and user identity *id*, if $t$ was created by the authorization server *as* and if the AS has created $t$ for the client $c$ and the identity *id*.

More precisely, the access token is exchanged for an authorization code (at the token endpoint of the AS), which is issued for a specific client. This is also the client to which the access token is associated with. The user identity with which the access token is associated is the user identity that authenticated at the AS (i.e., logged in at the website of the AS). In the model, the AS associates the access token with the client identifier and user identity by storing a sequence containing the identity, the client identifier and the access token (i.e., $\langle id, clientId, t, \mathtt{r} \rangle$, $\langle \mathtt{MTLS}, id, clientId, t, key, \mathtt{rw} \rangle$ or $\langle \mathtt{OAUTB}, id, clientId, t, key', \mathtt{rw} \rangle$). Furthermore, the last entry of the sequence indicates if the client is using the Read-Only or the Read-Write flow. In addition to this, for the Read-Write flow, the AS stores whether the access token is bound via mTLS or OAUTB (along with the corresponding key with which the access token is associated).

We can now define authorization formally, again the explanation of this definition follows below.

*Definition 2 (Authorization Property).* We say that the FAPI web system with a network attacker $\mathcal{FAPI}$ *is secure w.r.t. authorization* iff for every run $\rho$ of $\mathcal{FAPI}$, every configuration $(S,E,N)$ in $\rho$, every authorization server $as \in \mathsf{AS}$ that is honest in $S$ with $s_0^{as}.\mathtt{resource\_servers}$ being domains of honest resource servers used by *as*, every identity $id \in \mathsf{ID}^{as}$ for which the corresponding browser, say $b$, is honest in $S$, every client $c \in \mathsf{C}$ that is honest in $S$ with client id *clientId* issued to $c$ by *as*, every resource server $rs \in \mathsf{RS}$ that is honest in $S$ such that $id \in s_0^{rs}.\mathtt{ids}$ (set of IDs handled by *rs*), $s_0^{rs}.\mathtt{authServ} \in \mathrm{dom}(as)$ (set of domains controlled by *as*) and with $dom_{rs} \in s_0^{as}.\mathtt{resource\_servers}$ (with $dom_{rs} \in \mathrm{dom}(rs)$), every access token $t$ associated with $c$, *as* and *id* and every resource access nonce $r \in s_0^{rs}.\mathtt{rNonce}[id] \cup s_0^{rs}.\mathtt{wNonce}[id]$ it holds true that:

If $r$ is contained in a response to a request $m$ sent to *rs* with $t \equiv m.header[\mathtt{Authorization}]$, then $r$ is not derivable from the attackers knowledge in $S$.

As outlined above, the authorization property states that if the honest resource server receives an access token associated with a client identifier, authorization server, and user identifier, then the corresponding resource access is not given to the attacker. Access to resources is modeled by nonces called *resource access nonces*. For each user identity, there is one set of nonces representing read access, and another set representing write access. In our model of the FAPI, when a resource server receives an access token associated with a user from a client, the resource server returns to the client one of the resource access nonces of the user, which in turn the client forwards to the user's browser. The above security property requires that the attacker does not obtain such a resource access nonce (under the assumptions state in the property). This captures that there should be no direct or indirect way for the attacker

to access the corresponding resource. In particular, the attacker should not be able to use a client such that he can access the resource through the client.

For the authorization property to be meaningful, we require that the involved participants are honest. For example, we require that the authorization server at which the identity is registered is honest. If this is not the case (i.e., the attacker controls the AS), then the attacker could trivially access resources. The same holds true for the client for which the access token is issued: If the user chooses a client that is controlled by the attacker, then the attacker can trivially access the resource (as the user authorized the attacker client to do so). In our model of the FAPI, the client (non-deterministically) chooses a resource server that the authorization server supports (this can be different for each login flow). As in the Read-Only flow, the access token would trivially leak to the attacker if the resource server is controlled by the attacker, we require that the resource servers that the AS supports are honest. Furthermore, in the WIM, the behavior of the user is subsumed in the browser model, therefore, we require that the browser that is responsible for the user identity that is involved in the flow should be honest. Otherwise, the attacker could trivially obtain the credentials of the user.

*2) Authentication:* Informally speaking, the authentication property states that an attacker should not be able to log in at a client under the identity of an honest user. More precisely, we require that in all runs $\rho$ of a FAPI web system $\mathcal{FAPI}$ if in $\rho$ a client considers an honest user (browser) whose ID is governed by an honest AS to be logged in (indicated by a service token which a user can use at the client), then the adversary cannot obtain the service token.

*3) Session Integrity:* There are two session integrity properties that capture that an honest user should not be logged in under the identity of the attacker and should not use resources of the attacker. As shown in Section IV-D, session integrity is not given for all configurations available in the FAPI. Therefore, we show a limited session integrity property that captures session integrity for web server clients that use OAUTB.

Nonetheless, our session integrity property here is stronger than those used in [6], [7] in the sense that we define (and prove) session integrity not only in the presence of web attackers, but also for the much stronger network attacker. (This is enabled by using the __Secure- prefix for cookies.)

*Session Integrity for Authorization for Web Server Clients with OAUTB:* Intuitively, this property states that for all runs $\rho$ of a FAPI web system $\mathcal{FAPI}$, if an honest user can access the resource of some identity $u$ (registered at AS $as$) through the honest web server client $c$, where $c$ uses OAUTB as the holder of key mechanism, then (1) the user started the flow at $c$ and (2) if $as$ is honest, the user authenticated at the $as$ using the identity $u$.

*Session Integrity for Authentication for Web Server Clients with OAUTB:* Similar to the previous property, this property states that for all runs $\rho$ of a FAPI web system $\mathcal{FAPI}$, if an honest user is logged in at the honest client $c$ under some identity $u$ (registered at AS $as$), with $c$ being a web server

client using OAUTB as the holder of key mechanism, then (1) the user started the flow at $c$ and (2) if $as$ is honest, the user authenticated at the $as$ using the identity $u$.

By *Session Integrity for Web Server Clients with OAUTB* we denote the conjunction of both properties.

Now, our main theorem says that these properties are satisfied for all FAPI web systems.

*Theorem 1.* Let $\mathcal{FAPI}$ be a FAPI web system with a network attacker. Then, $\mathcal{FAPI}$ is secure w.r.t. authorization and authentication. Furthermore, $\mathcal{FAPI}$ is secure w.r.t. session integrity for web server clients with OAUTB.

We emphasize that the FAPI web systems take into account the strong attacker the FAPI is supposed to withstand as explained in Section III-C. Such attackers immediately break plain OAuth 2.0 and OpenID Connect. This, together with the various OAuth 2.0 security extensions which the FAPI uses and combines in different ways, and which have not formally been analyzed before, makes the proof challenging.

## VI. CONCLUSION

In this paper, we performed the first formal analysis of an Open Banking API, namely the OpenID Financial-grade API. Based on the Web Infrastructure Model, we built a comprehensive model comprising all protocol participants (clients, authorization servers, and resource servers) and all important options employed in the FAPI: clients can be app clients or web server clients and can make use of either the Read-Only or the Read-Write profile. We modeled all specified methods for authenticating at the authorization server and both mechanisms for binding tokens to the client, namely, Mutual TLS and OAuth 2.0 Token Binding. We also modeled PKCE, JWS Client Assertions, and the JWT Secured Authorization Response Mode (JARM).

Based on this model, we then defined precise security properties for the FAPI, namely authorization, authentication, and session integrity. While trying to prove these properties for the FAPI, we found several vulnerabilities that can enable an attacker to access protected resources belonging to an honest user or perform attacks on session integrity. We developed fixes against these attacks and formally verified the security of the (fixed) OpenID FAPI.

This is an important result since the FAPI enjoys wide industry support and is a promising candidate for the future lead in open banking APIs. Financial-grade applications entail very high security requirements that make a thorough formal security analysis, as performed in this paper, indispensable.

Our work also constitutes the very first analysis of various OAuth security extensions, namely PKCE, OAuth mTLS, OAUTB, JARM, and JWS Client Assertions.

REFERENCES

[1] "Blurred Lines: How FinTech Is Shaping Financial Services," 2016. PwC Global Fin-Tech Report.

[2] European Union, "DIRECTIVE (EU) 2015/2366 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL." https://eur-lex.europa.eu/legal-content/EN/TXT/HTML/?uri=CELEX:32015L2366&from=DE.

[3] S. T. Mnuchin and C. S. Phillips, "A Financial System That Creates Economic Opportunities – Nonbank Financials, Fintech, and Innovation." https://home.treasury.gov/sites/default/files/2018-08/A-Financial-System-that-Creates-Economic-Opportunities---Nonbank-Financials-Fintech-and-Innovation_0.pdf.

[4] M. Leszcz, "The UK Open Banking Implementation Entity Adopts the OpenID Foundation Financial-Grade API (FAPI) Specification & Certification Program." https://openid.net/2018/07/12/the-uk-open-banking-implementation-entity-adopts-the-openid-foundation-financial-grade-api-fapi-specification-certification-program/.

[5] OpenID Financial-grade API Working Group, "OpenID Foundation Financial-grade API (FAPI)." Aug. 23, 2018. https://bitbucket.org/openid/fapi/src/ceb0f829bc532e9c540efaa94f6f96d007371ca2/.

[6] D. Fett, R. Küsters, and G. Schmitz, "A Comprehensive Formal Security Analysis of OAuth 2.0," in Proceedings of the 23nd ACM SIGSAC Conference on Computer and Communications Security (CCS 2016), pp. 1204–1215, ACM, 2016.

[7] D. Fett, R. Küsters, and G. Schmitz, "The Web SSO Standard OpenID Connect: In-Depth Formal Security Analysis and Security Guidelines," in IEEE 30th Computer Security Foundations Symposium (CSF 2017), IEEE Computer Society, 2017.

[8] A. Kumar, "Using automated model analysis for reasoning about security of web protocols," in Proceedings of the 28th Annual Computer Security Applications Conference on - ACSAC'12, Association for Computing Machinery (ACM), 2012.

[9] C. Bansal, K. Bhargavan, and S. Maffeis, "Discovering Concrete Attacks on Website Authorization by Formal Analysis," in 25th IEEE Computer Security Foundations Symposium, CSF 2012 (S. Chong, ed.), pp. 247–262, IEEE Computer Society, 2012.

[10] C. Bansal, K. Bhargavan, A. Delignat-Lavaud, and S. Maffeis, "Discovering Concrete Attacks on Website Authorization by Formal Analysis," Journal of Computer Security, vol. 22, no. 4, pp. 601–657, 2014.

[11] R. Wang, Y. Zhou, S. Chen, S. Qadeer, D. Evans, and Y. Gurevich, "Explicating SDKs: Uncovering Assumptions Underlying Secure Authentication and Authorization," in Proceedings of the 22th USENIX Security Symposium, Washington, DC, USA, August 14-16, 2013, pp. 399–314, USENIX Association, 2013.

[12] S. Pai, Y. Sharma, S. Kumar, R. M. Pai, and S. Singh, "Formal Verification of OAuth 2.0 Using Alloy Framework," in CSNT '11 Proceedings of the 2011 International Conference on Communication Systems and Network Technologies, pp. 655–659, Proceedings of the International Conference on Communication Systems and Network Technologies, 2011.

[13] S. Chari, C. S. Jutla, and A. Roy, "Universally Composable Security Analysis of OAuth v2.0," IACR Cryptology ePrint Archive, vol. 2011, p. 526, 2011.

[14] S.-T. Sun and K. Beznosov, "The Devil is in the (Implementation) Details: An Empirical Analysis of OAuth SSO Systems," in ACM Conference on Computer and Communications Security, CCS'12 (T. Yu, G. Danezis, and V. D. Gligor, eds.), pp. 378–390, ACM, 2012.

[15] W. Li and C. J. Mitchell, "Security issues in OAuth 2.0 SSO implementations," in Information Security - 17th International Conference, ISC 2014, Hong Kong, China, October 12-14, 2014. Proceedings, pp. 529–541, 2014.

[16] R. Yang, G. Li, W. C. Lau, K. Zhang, and P. Hu, "Model-based Security Testing: An Empirical Study on OAuth 2.0 Implementations," in Proceedings of the 11th ACM on Asia Conference on Computer and Communications Security, AsiaCCS 2016, Xi'an, China, May 30 - June 3, 2016, pp. 651–662, ACM, 2016.

[17] E. Shernan, H. Carter, D. Tian, P. Traynor, and K. R. B. Butler, "More Guidelines Than Rules: CSRF Vulnerabilities from Noncompliant OAuth 2.0 Implementations," in Detection of Intrusions and Malware, and Vulnerability Assessment - 12th International Conference, DIMVA 2015, Milan, Italy, July 9-10, 2015, Proceedings, vol. 9148 of Lecture Notes in Computer Science, pp. 239–260, Springer, 2015.

[18] E. Y. Chen, Y. Pei, S. Chen, Y. Tian, R. Kotcher, and P. Tague, "OAuth Demystified for Mobile Application Developers," in Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security - CCS '14, pp. 892–903, 2014.

[19] M. Shehab and F. Mohsen, "Towards Enhancing the Security of OAuth Implementations in Smart Phones," in 2014 IEEE International Conference on Mobile Services, Institute of Electrical & Electronics Engineers (IEEE), 6 2014.

[20] W. Li and C. J. Mitchell, "Analysing the Security of Google's Implementation of OpenID Connect," in Proceedings of the 13th International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment (DIMVA), vol. 9721, pp. 357–376, 2016.

[21] V. Mladenov, C. Mainka, J. Krautwald, F. Feldmann, and J. Schwenk, "On the security of modern Single Sign-On Protocols: Second-Order Vulnerabilities in OpenID Connect," CoRR, vol. abs/1508.04324v2, 2016.

[22] D. Fett, R. Küsters, and G. Schmitz, "An Expressive Model for the Web Infrastructure: Definition and Application to the BrowserID SSO System," in 35th IEEE Symposium on Security and Privacy (S&P 2014), pp. 673–688, IEEE Computer Society, 2014.

[23] D. Fett, R. Küsters, and G. Schmitz, "SPRESSO: A Secure, Privacy-Respecting Single Sign-On System for the Web," in Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security, Denver, CO, USA, October 12-6, 2015, pp. 1358–1369, ACM, 2015.

[24] D. Fett, R. Küsters, and G. Schmitz, "Analyzing the BrowserID SSO System with Primary Identity Providers Using an Expressive Model of the Web," in Computer Security - ESORICS 2015 - 20th European Symposium on Research in Computer Security, Vienna, Austria, September 21-25, 2015, Proceedings, Part I, vol. 9326 of Lecture Notes in Computer Science, pp. 43–65, Springer, 2015.

[25] D. Hardt (ed.), "RFC6749 – The OAuth 2.0 Authorization Framework." IETF. Oct. 2012. https://tools.ietf.org/html/rfc6749.

[26] N. Sakimura, J. Bradley, M. Jones, B. de Medeiros, and C. Mortimore, "OpenID Connect Core 1.0 incorporating errata set 1." OpenID Foundation. Nov. 8, 2014. http://openid.net/specs/openid-connect-core-1_0.html.

[27] J. Richer (ed.), "RFC7662 – OAuth 2.0 Token Introspection." IETF. Oct. 2015. https://tools.ietf.org/html/rfc7662.

[28] N. Sakimura (Ed.), J. Bradley, and N. Agarwal, "Proof Key for Code Exchange by OAuth Public Clients." RFC 7636 (Proposed Standard), Sept. 2015.

[29] W. Denniss and J. Bradley, "OAuth 2.0 for Native Apps," RFC, vol. 8252, pp. 1–21, 2017.

[30] M. Jones, J. Bradley, and N. Sakimura, "RFC7519 – JSON Web Token (JWT)." IETF. May 2015. https://tools.ietf.org/html/rfc7519.

[31] M. Jones, J. Bradley, and N. Sakimura, "RFC7515 – JSON Web Signature (JWS)." IETF. May 2015. https://tools.ietf.org/html/rfc7515.

[32] B. Campbell, J. Bradley, N. Sakimura, and T. Lodderstedt, "OAuth 2.0 Mutual TLS Client Authentication and Certificate Bound Access Tokens," Internet-Draft draft-ietf-oauth-mtls-09, Internet Engineering Task Force, June 2018. Work in Progress.

[33] M. Jones, B. Campbell, J. Bradley, and W. Denniss, "OAuth 2.0 Token Binding - draft-ietf-oauth-token-binding-07." https://www.ietf.org/id/draft-ietf-oauth-token-binding-07.txt.

[34] A. Popov, M. Nystrom, D. Balfanz, A. Langley, and J. Hodges, "The Token Binding Protocol Version 1.0." RFC 8471, Oct. 2018.

[35] A. Popov, M. Nystrom, D. Balfanz, and A. Langley, "Transport Layer Security (TLS) Extension for Token Binding Protocol Negotiation." RFC 8472, Oct. 2018.

[36] A. Popov, M. Nystrom, D. Balfanz, A. Langley, N. Harper, and J. Hodges, "Token Binding over HTTP." RFC 8473, Oct. 2018.

[37] E. Rescorla, "Keying Material Exporters for Transport Layer Security (TLS)." RFC 5705, Mar. 2010.

[38] A. Popov, M. Nystrom, D. Balfanz, A. Langley, N. Harper, and J. Hodges, "Token Binding over HTTP," internet-draft, Internet Engineering Task Force, June 2018. Work in Progress.

[39] T. Lodderstedt (ed.), "JWT Secured Authorization Response Mode for OAuth 2.0 (JARM)." Aug. 23, 2018. https://bitbucket.org/openid/fapi/src/ceb0f829bc532e9c540efaa94f6f96d007371ca2/Financial_API_JWT_Secured_Authorization_Response_Mode.md.

[40] OpenID Financial-grade API Working Group, "Financial API - Part 1: Read-Only API Security Profile." Aug. 23, 2018. https://bitbucket.org/openid/fapi/src/ceb0f829bc532e9c540efaa94f6f96d007371ca2/Financial_API_WD_001.md.

[41] T. Lodderstedt, J. Bradley, A. Labunets, and D. Fett, "OAuth 2.0 Security Best Current Practice," 10 2018. https://tools.ietf.org/html/draft-ietf-oauth-security-topics.

[42] OpenID Financial-grade API Working Group, "Financial API - Part 2: Read and Write API Security Profile." Aug. 23, 2018. https://bitbucket.org/openid/fapi/src/ceb0f829bc532e9c540efaa94f6f96d007371ca2/Financial_API_WD_002.md.

[43] T. Lodderstedt (ed.), M. McGloin, and P. Hunt, "RFC6819 – OAuth 2.0 Threat Model and Security Considerations." IETF. Jan. 2013. https://tools.ietf.org/html/rfc6819.

[44] D. Fett, R. Küsters, and G. Schmitz, "The Web SSO Standard OpenID Connect: In-Depth Formal Analysis and Security Guidelines," Tech. Rep. arXiv:1704.08539, arXiv, 2017. Available at http://arxiv.org/abs/1704.08539.

[45] M. Abadi and C. Fournet, "Mobile Values, New Names, and Secure Communication," in *Proceedings of the 28th ACM Symposium on Principles of Programming Languages (POPL 2001)*, pp. 104–115, ACM Press, 2001.

[46] A. Barth and M. West, "Cookies: HTTP State Management Mechanism." https://httpwg.org/http-extensions/rfc6265bis.html.

[47] T. Dierks and E. Rescorla, "The Transport Layer Security (TLS) Protocol Version 1.2." RFC 5246 (Proposed Standard), Aug. 2008. Updated by RFCs 5746, 5878, 6176.

[48] A. Popov, M. Nystroem, D. Balfanz, A. Langley, and J. Hodges, "The Token Binding Protocol Version 1.0," Internet-Draft draft-ietf-tokbind-protocol-19, IETF Secretariat, 5 2018.

[49] K. Bhargavan, A. Delignat-Lavaud, A. Pironti, A. Langley, and M. Ray, "Transport Layer Security (TLS) Session Hash and Extended Master Secret Extension," *RFC*, vol. 7627, pp. 1–15, 2015.

[50] K. Bhargavan, A. Delignat-Lavaud, C. Fournet, A. Pironti, and P. Strub, "Triple Handshakes and Cookie Cutters: Breaking and Fixing Authentication over TLS," in *2014 IEEE Symposium on Security and Privacy, SP 2014, Berkeley, CA, USA, May 18-21, 2014*, pp. 98–113, IEEE Computer Society, 2014.

[51] N. Sakimura, J. Bradley, M. Jones, and E. Jay, "OpenID Connect Discovery 1.0 incorporating errata set 1." OpenID Foundation. Nov. 8, 2014. http://openid.net/specs/openid-connect-discovery-1_0.html.

[52] N. Sakimura, J. Bradley, and M. Jones, "OpenID Connect Dynamic Client Registration 1.0 incorporating errata set 1." OpenID Foundation. Nov. 8, 2014. http://openid.net/specs/openid-connect-registration-1_0.html.

**Figure 7.** Leakage of Authorization Request Attack

AUTHORIZATION REQUEST LEAK ATTACK – DETAILS

We here provide further details about the authorization request leak attack, which was only sketched in Section IV-D.

A concrete instantiation of this attack is shown in Figure 7, where the scenario is based on the Read-Only flow of a public client. As explained below, similar attacks also work for all other configurations of the FAPI (except for web server clients which use OAUTB, for which, as mentioned, we show that they are not susceptible in Section V).

In the Authorization Request Leak Attack, the client sends the authorization request to the browser in Step 2, where it leaks to the attacker in Step 3. From here on, the attacker behaves as the browser and logs himself in (Step 5), hence, the authorization code received in Step 6 is associated with the identity of the attacker.

The state value used in the authorization request aims at preventing Cross-Site Request Forgery (CSRF) attacks. However, as the state value leaks, this protection does not work. For showing that this is the case, we assume that a CSRF attack happens. If, for example, the user is visiting a website that is controlled by the attacker, then the attacker can send, from the browser of the user, a request to the AS containing the code and the state value (Step 8). As the state received by the client is the same that it included in the authorization request, the client continues the flow and uses the code to retrieve an access token in Steps 9 and 10.

This access token is associated with the attacker, which means that the honest user is accessing resources belonging to the attacker.

As a result, the honest user can be logged in under the identity of the attacker if the authorization server returns an id token. In the case of the Read-Write flow, the honest user can modify resources of the attacker: for example, she might upload personal documents to the account of the attacker.

As noted above, this attack might happen for all configurations, except for the Read-Write flow when the client is a web server client using OAUTB (see Figure 3).

In all other configurations, this attack can happen as the attacker can behave exactly like the browser of the honest user, i.e., after receiving the authorization request, the attacker can send this request to the AS, log in under his own identity, and would then receive a response that the client accepts. The only flow in which this is different is the Read-Write flow where the client is a web server and uses OAUTB, as here, the browser (and therefore, also the attacker) needs to prove possession of a key

pair (i.e., the key pair used for the client). As the attacker cannot prove possession of the private key of the key pair which the browser uses for the client, the AS would then stop the flow. (In the other flows, the AS does not check if the response was sent by the browser that logged in the user.)

If we say that the FAPI is not required to be secure if the authorization request leaks (i.e., if we remove the assumption that the authorization request leaks), then the flow is still not secure, as the authorization response might still leak to the attacker (see Section III-C1), which also contains the state value. More precisely, the authorization response might leak in the case of app clients due to the operating system sending the response to the attacker app (for details, see Section II-B). After receiving the authorization response, the attacker app knows the state value and can start a new flow using this value. The attacker can then continue from Step ③ (Figure 7), and when receiving the authorization response (which is a URI containing the OAuth parameters), he could, using his own app that runs on the device of the victim, call the legitimate client app with this URI (i.e., with the code that is associated with the identity of the attacker and the state value with which the client started the flow). The effect of this is that the legitimate app, at which the honest user started the flow, would continue the flow using an authorization code associated with the attacker. Therefore, the honest user would either be logged in with the identity of the attacker or use the resources of the attacker.

We note that even encrypting the state value contained in the authorization request does not solve the problem, as the attacker is using the whole authorization request. (Strictly speaking, he acts as the browser of the honest user).

We here provide more details about the Web Infrastructure Model.

*a) Signature and Messages:* As mentioned, the WIM follows the Dolev-Yao approach where messages are expressed as formal terms over a signature $\Sigma$. For example, in the WIM an HTTP request is represented as a term $r$ containing a nonce, an HTTP method, a domain name, a path, URI parameters, request headers, and a message body. For instance, an HTTP request for the URI http://ex.com/show?p=1 is represented as $r := \langle \text{HTTPReq}, n_1, \text{GET}, \text{ex.com}, /\text{show}, \langle\langle \text{p}, 1\rangle\rangle, \langle\rangle, \langle\rangle \rangle$ where the body and the list of request headers is empty. An HTTPS request for $r$ is of the form $\text{enc}_\text{a}(\langle r, k'\rangle, \text{pub}(k_{\text{ex.com}}))$, 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 Dolev-Yao 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 $\text{dec}_\text{a}(\text{enc}_\text{a}(x, \text{pub}(y)), y) = x$. With this, we have that, for example, $\text{dec}_\text{a}(\text{enc}_\text{a}(\langle r, k'\rangle, \text{pub}(k_{\text{ex.com}})), k_{\text{ex.com}}) \equiv \langle r, k'\rangle$, i.e., these two terms are equivalent w.r.t. the equational theory.

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

*c) Running a system:* As mentioned, a run of a system is a sequence of configurations. The transition from one configuration to the next configuration in a run is called a *processing step*. We write, for example, $Q = (S, E, N) \to (S', E', N')$ to denote the transition from the configuration $(S, E, N)$ to the configuration $(S', E', N')$, where $S$ and $S'$ are the states of the processes in the system, $E$ and $E'$ are pools of waiting events, and $N$ and $N'$ are sequences of unused nonces.

*d) 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 selected web pages when needed. 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 XMLHttpRequests 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 relevant web standards.

A browser can output messages on the network of different types, namely DNS and HTTP(S) (including XMLHttpRequests), 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 so-called trigger message upon which the browser non-deterministically chooses an action, for instance, to trigger a script in some document. The script now outputs a command, as described above, which is then further processed by the browser. Browsers can also become corrupted, i.e., be taken over by web and network attackers. Once corrupted, a browser behaves like an attacker process.

As detailed in Appendix G, we extended the browser model of the WIM slightly in order to incorporate OAUTB in the browser model. We furthermore added the behavior of the *__Secure-* prefix of cookies to the model, which specifies that such cookies shall only be accepted when they are transmitted over secure channels [46]. Note that for the FAPI, mTLS is only needed between clients and servers. Therefore, mTLS has been modeled on top of the WIM, i.e., as part of the modeling of FAPI clients and servers. The servers we modeled for the FAPI of course also support OAUTB.

# Appendix C
## Excerpt of Client Model

In this section, we provide a brief excerpt of the client model in order to give an impression of the formal model. See Appendix I for the full formal model of the FAPI.

The excerpt given in Algorithm 1 shows how the client prepares and sends the token request to the authorization server, i.e., the part in which the client sends the authorization code in exchange for an access token (and depending on the flow, also an id token).

This function is called by the client. The first two inputs are the session identifier of the session (i.e., the session of the resource owner at the client) and the authorization code that the client wants to send to the AS. The value *responseValue* contains information related to mTLS or OAUTB (if used for the current flow). The last input is the current state of the client.

In Lines 5 to 8, the client chooses either the token endpoint of the AS or some URL that was chosen non-deterministically. This models the assumption shown in Section III-C4, which requires the Read-Write profile of the FAPI to be secure even if the token endpoint is misconfigured.

Starting from Line 15, the function chooses the parameters of the request that depend on the flow and configuration (see Figure 3).

If the client uses the Read-Only profile, the token request always contains the PKCE verifier (Line 15). For a confidential client (which means that the client has to authenticate at the token endpoint), the client either authenticates using JWS Client Assertions (Line 20, see also Section II-C), or with mTLS (Line 26; for details on our model of mTLS refer to Appendix E).

If the client uses the Read-Write profile, the client uses either mTLS (again Line 26) or OAUTB (Line 32; for details on our model of OAUTB refer to Appendix F).

**Algorithm 1** Relation of a Client $R^c$ – Request to token endpoint.

---

1: **function** SEND_TOKEN_REQUEST(*sessionId*, *code*, *responseValue*, $s'$)
2:   **let** *session* := $s'$.sessions[*sessionId*]
3:   **let** *identity* := *session*[identity]
4:   **let** *issuer* := $s'$.issuerCache[*identity*]
5:   **if** *session*[misconfiguredTEp] $\equiv \top$ **then**
6:     **let** *url* := *session*[token_ep]
7:   **else**
8:     **let** *url* := $s'$.oidcConfigCache[*issuer*][token_ep]
9:   **let** *credentials* := $s'$.clientCredentialsCache[*issuer*]
10:   **let** *clientId* := *credentials*[client_id]
11:   **let** *clientType* := *credentials*[client_type]
12:   **let** *profile* := *credentials*[profile]
13:   **let** *isApp* := *credentials*[is_app]
14:   **let** *body* := [grant_type:authorization_code, code:*code*, redirect_uri:*session*[redirect_uri], client_id:*clientId*]
15:   **if** *profile* $\equiv$ r **then**
16:     **let** *body*[pkce_verifier] := *session*[pkce_verifier]
17:   **if** *profile* $\equiv$ r $\wedge$ *clientType* $\equiv$ pub **then**
18:     **let** *message* := $\langle$HTTPReq, $\nu_2$, POST, *url*.domain, *url*.path, *url*.parameters, $\bot$, *body*$\rangle$
19:     **call** HTTPS_SIMPLE_SEND([responseTo:TOKEN, session:*sessionId*], *message*, $s'$)
20:   **else if** *profile* $\equiv$ r $\wedge$ *clientType* $\equiv$ conf_JWS **then**
21:     **let** *clientSecret* := *credentials*[client_secret]
22:     **let** *jwt* := [iss:*clientId*, aud:*url*.domain]
23:     **let** *body*[assertion] := mac(*jwt*, *clientSecret*)
24:     **let** *message* := $\langle$HTTPReq, $\nu_2$, POST, *url*.domain, *url*.path, *url*.parameters, $\bot$, *body*$\rangle$
25:     **call** HTTPS_SIMPLE_SEND([responseTo:TOKEN, session:*sessionId*], *message*, $s'$)
26:   **else if** *clientType* $\equiv$ conf_MTLS **then**     $\rightarrow$ both profiles
27:     **if** *responseValue*[type] $\not\equiv$ MTLS **then**
28:       **stop**
29:     **let** *body*[TLS_AuthN] := *responseValue*[mtls_nonce]
30:     **let** *message* := $\langle$HTTPReq, $\nu_2$, POST, *url*.domain, *url*.path, *url*.parameters, $\bot$, *body*$\rangle$
31:     **call** HTTPS_SIMPLE_SEND([responseTo:TOKEN, session:*sessionId*], *message*, $s'$)
32:   **else**     $\rightarrow$ rw with OAUTB
33:     **if** *responseValue*[type] $\not\equiv$ OAUTB **then**
34:       **stop**
35:     **let** *ekm* := *responseValue*[ekm]
36:     **let** *TB_AS* := $s'$.TBindings[*url*.host]     $\rightarrow$ priv. key
37:     **let** *TB_RS* := $s'$.TBindings[*session*[RS]]     $\rightarrow$ priv. key
38:     **let** *TB_Msg_prov* := [id:pub(*TB_AS*), sig:sig(*ekm*, *TB_AS*)]
39:     **let** *TB_Msg_ref* := [id:pub(*TB_RS*), sig:sig(*ekm*, *TB_RS*)]
40:     **let** *headers* := [Sec-Token-Binding:[prov:*TB_Msg_prov*, ref:*TB_Msg_ref*]]
41:     **if** *clientType* $\equiv$ conf_OAUTB **then**     $\rightarrow$ client authentication
42:       **let** *clientSecret* := *credentials*[client_secret]
43:       **let** *jwt* := [iss:*clientId*, aud:*url*.domain,]
44:       **let** *body*[assertion] := mac(*jwt*, *clientSecret*)
45:     **if** *isApp* $\equiv \bot$ **then**     $\rightarrow$ W.S. client: TBID used by browser
46:       **let** *body*[pkce_verifier] := *session*[browserTBID]
47:     **let** *message* := $\langle$HTTPReq, $\nu_2$, POST, *url*.domain, *url*.path, *url*.parameters, *headers*, *body*$\rangle$
48:     **call** HTTPS_SIMPLE_SEND([responseTo:TOKEN, session:*sessionId*], *message*, $s'$)

---

We here provide a proof sketch of Theorem 1 that is concerned with the authorization property. The complete formal proof of this theorem is given in Appendix L.

For proving the authorization property, we show that when a participant provides access to a resource, i.e., by sending a resource access nonce, this access is not provided to the attacker:

*a) Resource server does not provide the attacker access to resources:* We show that the resource server does not provide the attacker access to resources of an honest user.

In case of the Read-Only flow, we show that an access token associated with an honest client, an honest authorization server, and an honest identity does not leak to the attacker, and therefore, the attacker cannot obtain access to resources.

In case of the Read-Write flow, such an access token might leak to the attacker, but this token cannot be used by the attacker at the resource server due to Token Binding, either via OAUTB or mTLS.

*b) Web server client does not provide the attacker access to resources:* App clients are only usable via the device they are running on, i.e., they are not usable over the network (by which we mean that if, for example, the user wants to view one of her documents with an app client, she does this directly using the device). Therefore, we only look at the case of web server clients, as such a client can be used over the network, e.g., by the browser of the end-user or by the attacker.

In the following, we show that honest web server clients do not provide the attacker access to resources belonging to an honest identity. We show this for all possible configurations that could trick the client into doing so, e.g., with a misconfigured token endpoint or with an authorization server controlled by the attacker that returns a leaked access token.

The access to the resource is provided to the sender of the redirection request. To access a resource, this means that the attacker must have sent the request to the redirection endpoint of the client.

For a Read-Only flow, the token endpoint is configured correctly. This means that the attacker must include a code in the request such that the client can exchange it for an access token. We show that such a code (associated with an honest identity and the client) does not leak to an attacker.

For a Read-Write flow, the token endpoint can be misconfigured such that it is controlled by the attacker, and we also assume that access tokens leak to the attacker (see Section III-C).

We show that a leaked access token cannot be used at the client by the attacker. If only the token endpoint is controlled by the attacker, he must include an id token (when using the OIDC Hybrid flow, see below for the Authorization Code flow with JARM) in the token response such that it contains the hash of the access token and be signed by the honest authorization server (the hash of the access token was not included in the original draft and was included by us as a mitigation in Section IV-B). However, such an id token does not leak to the attacker, which prevents the use of leaked access tokens at misconfigured token endpoints. For the Authorization Code flow with JARM, the attacker would need a response JWS. As in the case of the Hybrid flow, we show that the response JWS needed by the client for accessing resources of an honest identity does not leak.

A leaked access token can also be used by the attacker if the client chooses an authorization server under the control of the attacker. Here, the id tokens are created by the attacker and accepted by the client. For preventing the use of this access token, the client includes the issuer of the second id token (or of the response JWS defined by JARM) in the request to the resource server, as detailed in Section IV-A. As each resource server has one preconfigured authorization server, the resource server does not provide access to a resource in this case.

The only remaining case is that the attacker includes a code associated with the honest user in the request to the redirection endpoint of the client. For the Hybrid flow, both id tokens contained in the authorization response and in the token response are required to have the same subject attribute and the same issuer value, which means that they are both signed by the authorization server. However, such an id token does not leak to the attacker, which means that the client will stop the flow when receiving the second id token contained in the token response. When using JARM, this would require the attacker to send a response JWS signed by the authorization server that contains the code that belongs to an honest client and an honest user identity. In the technical report, we show that such a response JWS does not leak to the attacker.

APPENDIX E
MODELING MTLS

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.

Figure 8 shows an overview of how we modeled mTLS. 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 authorization server. The authorization server 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 [47, 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 in Section II-D, the client uses the same certificate it used for the token request when sending the access token to the resource server. While the resource server has to check the possession of corresponding private keys, the validity of the certificate was already checked at the authorization server and can be ignored by the resource server. Therefore, in our model of the FAPI, the client does not send its client id to the resource server, but its public key, and the resource server encrypts the message with this public key.



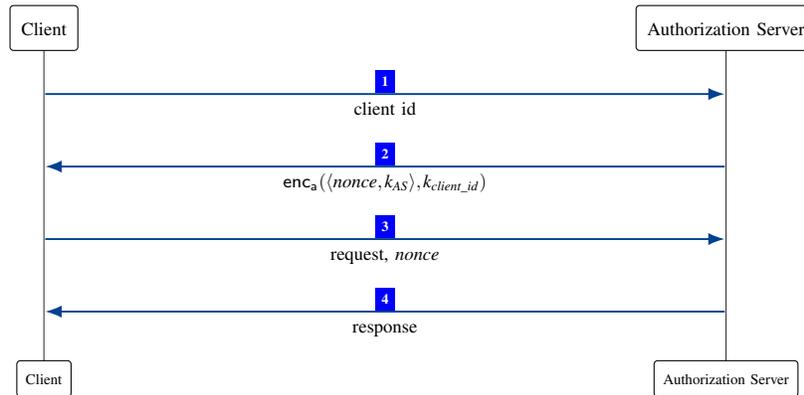**Figure 8.** Overview of mTLS

All messages are sent by the generic HTTPS server model (Appendix H), 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 9 shows the complete messages, i.e., with the encryption used for transmitting the messages.

**Figure 9.** Detailed view on mTLS

In the following, we describe more details of OAuth 2.0 Token Binding (see Section II-E) and the modeling within the Web Infrastructure Model.

## A. Exported Keying Material

The proof of possession of a private key is done by signing a so-called *Exported Keying Material* (EKM). This value is created with parameters of the TLS connection such that it is unique to the connection. Therefore, it is not possible for an attacker to reuse signed EKM values for an honest server. Before explaining how the EKM value is generated, we give a brief explanation of the relevant mechanisms used in TLS [47]:

- Client and Server Random: The client random and server random are values chosen by the client and server and transmitted in the TLS handshake.
- Pseudorandom Function: TLS specifies a pseudorandom function (PRF) that is used within the TLS protocol [47, Section 5].
- Premaster Secret: The premaster secret is a secret shared between both participants of the TLS connection, for example, created by a Diffie-Hellman key exchange. The length of the premaster secret varies depending on the method used for its creation.
- Master Secret: The master secret is created by applying the pseudorandom function to the premaster secret to generate a secret of a fixed length. It essentially has the value $\mathrm{PRF}(premaster\_secret, client\_random, server\_random)$ [47, Section 8.1]; we omitted constant values).

Using these building blocks, the EKM value [37] is essentially defined as

$$\mathrm{PRF}(master\_secret, client\_random, server\_random)$$

As noted in Section 7.5 of the Token Binding Protocol [48], the use of the Extended Master Secret TLS extension [49] is mandatory. This extension redefines the master secret as $\mathrm{PRF}(pre\_master\_secret, session\_hash)$, where *session_hash* is the hash of all TLS handshake messages of the session (again, omitting constant values). Without this extension, the *Triple Handshake Attack* [50] can be applied, which eventually leads to the creation of two TLS sessions with the same master secret. The basic idea is that if an honest client establishes a TLS connection to a malicious server (for example, when the token endpoint is misconfigured), the attacker can relay the random values chosen by the honest client and server, hence, creating TLS connections to both the client and server with the same master secret and therefore, with the same EKM value. By including the hash over the TLS handshake messages, relevant information about the session, like the certificate used by the server (which is exchanged in the handshake), influence the value of the master secret.

## B. Modeling Token Binding

The main difficulty of modeling OAuth 2.0 Token Binding is the high level of abstraction of TLS within the WIM, as already explained in Appendix E.

An overview of how we modeled OAUTB is shown in Figure 10. For compensating the absence of the TLS handshake, both participants choose nonces, as shown in Steps ① and ②. When the client sends the actual request, it creates and includes a Token Binding message, as shown in Step ③. As explained above, the client includes not only the two nonces, but also the public key of the authorization server in the EKM value for modeling the extended master secret.

For completeness, Figure 11 shows the entire messages, i.e., with the encryption done by the generic HTTPS server (Appendix H).

## C. Access Token issued from Authorization Endpoint

In the Read-Write profile, an access token issued from the authorization endpoint is required to be token bound. When using OAUTB, this means that the token must be bound to the Token Binding ID which the client uses for the resource server. However, the authorization request is sent to the authorization server by the browser, which cannot send Token Binding messages with an ID used by the client. Therefore, our model does not include access tokens being issued from the authorization endpoint, as this would mean that additional communication between the client and the authorization server is needed, which is not fully specified yet.

We note that this issue is also present in the case of confidential clients using mTLS. Here, the access token cannot be bound to the certificate of the client, as the authorization request is not sent directly by the client to the authorization server, but by the browser. More details can be found in Section 4.5 of [32].

**Figure 10.** Overview of OAUTB



**Figure 11.** Detailed view on OAUTB

Within the scope of this technical report, we adhere to the WIM as defined in [44], where it was used for modeling and analyzing OpenID Connect. In the following, we describe the additions to the model for the analysis of the FAPI.

*A. Functions*

In addition to the function symbols defined in Appendix B of [44], we add the function symbol for hashing hash(.) to the signature $\Sigma$. For computing and verifying message authentication codes, we add mac(.,.) and checkmac(.,.). Regarding the equational theory, we additionally define $\mathsf{mac}(x,y) = \mathsf{hash}(\langle x,y \rangle)$ and $\mathsf{checkmac}(\mathsf{mac}(x,y),y) = \top$. Furthermore, we extend the definition of extractmsg such that it also extracts messages out of $\mathsf{mac}(x,y)$, i.e., $\mathsf{extractmsg}(\mathsf{mac}(x,y)) = x$. As a short form of the mapping from identities to their governor, we define $gov(.) = governor(.)$.

*B. Cookies*

In Appendix C of [44], a cookie is defined as a term of the form $\langle name, content \rangle$ with $name \in \mathcal{T}_{\mathcal{N}}$. As the name is a term, it may also be a sequence consisting of two part. If the name it consists of two parts, we call the first part of the sequence (i.e., *name*.1) the *prefix* of the name.

In the following, define the *Secure* prefix (see [46]): When the __*Secure* prefix is set, 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.

For modeling this, we require that the AddCookie function should be called with (as an additional argument) the protocol with which the corresponding request was sent, i.e., we modify the function PROCESSRESPONSE (using the same number for the algorithm as in [44]) :

---
**Algorithm 8** Web Browser Model: Process an HTTP response.

---
1: **function** PROCESSRESPONSE(*response*, *reference*, *request*, *requestUrl*, $s'$)
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)
    ...

---

The modified function AddCookie looks as follows (again using the same number as in [44]):

*Definition 43.* For a sequence of cookies (with pairwise different names) *oldcookies* and a cookie $c$, the sequence AddCookie(*oldcookies*, $c$, *protocol*) is defined by the following algorithm: If $((c.\mathtt{name}.1 \equiv \_\_\mathtt{Secure}) \Rightarrow (protocol \equiv S))$, then: Let $m := oldcookies$. Remove any $c'$ from $m$ that has $c.\mathtt{name} \equiv c'.\mathtt{name}$. Append $c$ to $m$ and return $m$.

*C. Browser*

We use the same browser model as defined in Appendix D of [44], with small modifications. Therefore, we do not show the full model here but limit the description to the algorithms that differ from the original browser model.

*Main Differences:*

- Algorithm 2: New input argument *refTB*. This value is saved along with the other information needed for the current DNS request.
- Algorithm 3: Sends messages directly to a receiver, encrypted with a symmetric key. The receiver and key are given as an input. This algorithm is used for sending the follow-up request when using OAUTB.
- Algorithm 4: If the response contains tb_nonce, the original message is sent together with the OAUTB message. Furthermore, the actions specified when the Include-Referred-Token-Binding-ID header is set are implemented here.
- Algorithm 5: When receiving a DNS response, the browser checks if OAUTB needs to be applied. In this case, the browser first sends a request to the OAUTB endpoint to obtain a new nonce. The original message is saved and send together with the OAUTB message in a second request.

When sending the first OAUTB request, the nonce $v_{n1}$ is used as a reference (in Algorithm 5).

---

**Algorithm 2** Web Browser Model: Prepare headers, do DNS resolution, save message.

---

1: **function** HTTP_SEND(*reference*, *message*, *url*, *origin*, *referrer*, *referrerPolicy*, $s'$, *refTB*)
2:     **if** *message*.host $\in^{\langle\rangle}$ $s'$.sts **then**
3:         **let** *url*.protocol := S
4:     **let** *cookies* := $\langle\{\langle c.\text{name}, c.\text{content.value}\rangle | c \in^{\langle\rangle} s'.\text{cookies}\,[message.\text{host}]$
        ↪ $\wedge(c.\text{content.secure} \implies (url.\text{protocol} = \text{S}))\}\rangle$
5:     **let** *message*.headers[Cookie] := *cookies*
6:     **if** *origin* $\not\equiv \perp$ **then**
7:         **let** *message*.headers[Origin] := *origin*
8:     **if** *referrerPolicy* $\equiv$ noreferrer **then**
9:         **let** *referrer* := $\perp$
10:    **if** *referrer* $\not\equiv \perp$ **then**
11:       **if** *referrerPolicy* $\equiv$ origin **then**
12:          **let** *referrer* := $\langle$URL, *referrer*.protocol, *referrer*.host, $/, \langle\rangle, \perp\rangle$   → Referrer stripped down to origin.
13:       **let** *referrer*.fragment := $\perp$   → Browsers do not send fragment identifiers in the Referer header.
14:       **let** *message*.headers[Referer] := *referrer*
15:    **let** $s'$.pendingDNS[$\nu_8$] := $\langle reference, message, url, refTB\rangle$
16:    **stop** $\langle\langle s'.\text{DNSaddress}, a, \langle\text{DNSResolve}, message.\text{host}, \nu_8\rangle\rangle\rangle, s'$

---

**Algorithm 3** Web Browser Model: Send message to IP encrypted with given sym. key.

---

1: **function** OAUTB_CONT_SEND(*reference*, *message*, *url*, *key*, *f*, $s'$)
2:    **let** $s'$.pendingRequests := $s'$.pendingRequests
    ↪ $+^{\langle\rangle}$ $\langle reference, message, url, key, f\rangle$
3:    **let** *message* := $\text{enc}_\text{s}(message, key)$
4:    **stop** $\langle\langle f, a, message\rangle\rangle, s'$

---

**Algorithm 4** Web Browser Model: Process an HTTP response.

---

1: **function** PROCESSRESPONSE(*response*, *reference*, *request*, *requestUrl*, $s'$, *key*, *f*)
2:    **if** tb_nonce $\in$ *response*.body **then**
3:       **let** $\langle orig\_reference, orig\_message, orig\_url, refTB\rangle$ := $s'$.tokenBindingRequests[*request*.nonce]
4:       **let** *ekm* := $\text{hash}(\langle request.\text{nonce}, response.\text{body}[\text{tb\_nonce}], \text{keyMapping}[request.\text{host}]\rangle)$
5:       **let** *tb_prov_priv* := $s'$.tokenBindings[*orig_url*.host]
6:       **let** *TB_Msg_provided* := [id:pub(*tb_prov_priv*), sig:sig(*ekm*, *tb_prov_priv*)]
7:       **if** *refTB* $\not\equiv \perp$ **then**
8:         **let** *tb_ref_priv* := $s'$.tokenBindings[*refTB*]
9:         **let** *TB_Msg_referred* := [id:pub(*tb_ref_priv*), sig:sig(*ekm*, *tb_ref_priv*)]
10:      **else**
11:        **let** *TB_Msg_referred* := $\langle\rangle$
12:      **let** *orig_message*.headers[Sec-Token-Binding] := [prov:*TB_Msg_provided*, ref:*TB_Msg_referred*]
13:      **call** TB_CONT_SEND(*orig_reference*, *orig_message*, *orig_url*, *key*, *f*, $s'$)
14:    **if** Set-Cookie $\in$ *response*.headers **then**
15:      **for each** $c \in^{\langle\rangle}$ *response*.headers[Set-Cookie], $c \in$ Cookies **do**
16:        **let** $s'$.cookies[*request*.host]
        ↪ := AddCookie($s'$.cookies[*request*.host], $c$, *requestUrl*.protocol)
17:    **if** Strict-Transport-Security $\in$ *response*.headers $\wedge$ *requestUrl*.protocol $\equiv$ S **then**
18:      **let** $s'$.sts := $s'$.sts $+^{\langle\rangle}$ *request*.host
19:    **if** Referer $\in$ *request*.headers **then**
20:      **let** *referrer* := *request*.headers[Referer]
21:    **else**
22:      **let** *referrer* := $\perp$
23:    **if** Location $\in$ *response*.headers $\wedge$ *response*.status $\in \{303, 307\}$ **then**
24:      **let** *url* := *response*.headers[Location]
25:      **if** *url*.fragment $\equiv \perp$ **then**
26:        **let** *url*.fragment := *requestUrl*.fragment
27:      **if** Include-Referred-Token-Binding-ID $\in$ *response*.headers **then**   → Always use TB
28:        **if** *response*.headers[Include-Referred-Token-Binding-ID] $\equiv \top$ **then**
29:          **let** $s'$.useTB[*response*.host] := $\top$
30:          **let** $s'$.useTB[*response*.headers[Location]] := $\top$
31:      **let** $method'$ := *request*.method

---

32:   **let** $body' := request.\text{body}$
33:   **if** $\text{Origin} \in request.\text{headers}$ **then**
34:    **let** $origin := \langle request.\text{headers}[\text{Origin}], \langle request.\text{host}, url.\text{protocol} \rangle \rangle$
35:   **else**
36:    **let** $origin := \bot$
37:   **if** $response.\text{status} \equiv 303 \wedge request.\text{method} \notin \{\text{GET}, \text{HEAD}\}$ **then**
38:    **let** $method' := \text{GET}$
39:    **let** $body' := \langle \rangle$
40:   **if** $\exists \overline{w} \in \mathsf{Subwindows}(s')$ **such that** $s'.\overline{w}.\text{nonce} \equiv reference$ **then**   $\rightarrow$ Do not redirect XHRs.
41:    **let** $req := \langle \text{HTTPReq}, \nu_6, method', url.\text{host}, url.\text{path}, \langle \rangle, url.\text{parameters}, body' \rangle$
42:    **let** $referrerPolicy := response.\text{headers}[\text{ReferrerPolicy}]$
43:    **if** $\text{Include-Referred-Token-Binding-ID} \in response.\text{headers}$ **then**
44:     **if** $response.\text{headers}[\text{Include-Referred-Token-Binding-ID}] \equiv \top$ **then**
45:      **let** $refTBID := response.\text{host}$
46:    **else**
47:     **let** $refTBID := \bot$
48:    **call** $\mathsf{HTTP\_SEND}(reference, req, url, origin, referrer, referrerPolicy, s', refTBID)$
49:  **if** $\exists \overline{w} \in \mathsf{Subwindows}(s')$ **such that** $s'.\overline{w}.\text{nonce} \equiv reference$ **then**   $\rightarrow$ normal response
50:   **if** $response.\text{body} \not\sim \langle *, * \rangle$ **then**
51:    **stop** $\{\}, s'$
52:   **let** $script := \pi_1(response.\text{body})$
53:   **let** $scriptstate := \pi_2(response.\text{body})$
54:   **let** $referrer := request.\text{headers}[\text{Referer}]$
55:   **let** $d := \langle \nu_7, requestUrl, response.\text{headers}, referrer, script, scriptstate, \langle \rangle, \langle \rangle, \top \rangle$
56:   **if** $s'.\overline{w}.\text{documents} \equiv \langle \rangle$ **then**
57:    **let** $s'.\overline{w}.\text{documents} := \langle d \rangle$
58:   **else**
59:    **let** $\overline{i} \leftarrow \mathbb{N}$ **such that** $s'.\overline{w}.\text{documents}.\overline{i}.\text{active} \equiv \top$
60:    **let** $s'.\overline{w}.\text{documents}.\overline{i}.\text{active} := \bot$
61:    **remove** $s'.\overline{w}.\text{documents}.(\overline{i}+1)$ and all following documents
    $\hookrightarrow$ from $s'.\overline{w}.\text{documents}$
62:    **let** $s'.\overline{w}.\text{documents} := s'.\overline{w}.\text{documents} +^{\langle \rangle} d$
63:   **stop** $\{\}, s'$
64:  **else if** $\exists \overline{w} \in \mathsf{Subwindows}(s'), \overline{d}$ **such that** $s'.\overline{d}.\text{nonce} \equiv \pi_1(reference)$
  $\hookrightarrow \wedge s'.\overline{d} = s'.\overline{w}.\text{activedocument}$ **then**   $\rightarrow$ process XHR response
65:   **let** $headers := response.\text{headers} - \text{Set-Cookie}$
66:   **let** $s'.\overline{d}.\text{scriptinputs} := s'.\overline{d}.\text{scriptinputs} +^{\langle \rangle}$
   $\langle \text{XMLHTTPREQUEST}, headers, response.\text{body}, \pi_2(reference) \rangle$

---

**Algorithm 5** Web Browser Model: Main Algorithm

---

**Input:** $\langle a, f, m \rangle, s$
1: **let** $s' := s$
2: **if** $s.\text{isCorrupted} \not\equiv \bot$ **then**
3:  **let** $s'.\text{pendingRequests} := \langle m, s.\text{pendingRequests} \rangle$   $\rightarrow$ Collect incoming messages
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** $m \equiv \text{TRIGGER}$ **then**   $\rightarrow$ A special trigger message.
8:  **let** $switch \leftarrow \{\text{script}, \text{urlbar}, \text{reload}, \text{forward}, \text{back}\}$
9:  **let** $\overline{w} \leftarrow \mathsf{Subwindows}(s')$ **such that** $s'.\overline{w}.\text{documents} \neq \langle \rangle$
  $\hookrightarrow$ **if possible; otherwise stop**   $\rightarrow$ Pointer to some window.
10:  **let** $\overline{tlw} \leftarrow \mathbb{N}$ **such that** $s'.\overline{tlw}.\text{documents} \neq \langle \rangle$
  $\hookrightarrow$ **if possible; otherwise stop**   $\rightarrow$ Pointer to some top-level window.
11:  **if** $switch \equiv \text{script}$ **then**   $\rightarrow$ Run some script.
12:   **let** $\overline{d} := \overline{w} +^{\langle \rangle} \text{activedocument}$
13:   **call** $\mathsf{RUNSCRIPT}(\overline{w}, \overline{d}, s')$
14:  **else if** $switch \equiv \text{urlbar}$ **then**   $\rightarrow$ Create some new request.
15:   **let** $newwindow \leftarrow \{\top, \bot\}$
16:   **if** $newwindow \equiv \top$ **then**   $\rightarrow$ Create a new window.
17:    **let** $windownonce := \nu_1$
18:    **let** $w' := \langle windownonce, \langle \rangle, \bot \rangle$
19:    **let** $s'.\text{windows} := s'.\text{windows} +^{\langle \rangle} w'$
20:   **else**   $\rightarrow$ Use existing top-level window.
21:    **let** $windownonce := s'.\overline{tlw}.\text{nonce}$
22:   **let** $protocol \leftarrow \{\text{P}, \text{S}\}$

```
23:        let host ← Doms
24:        let path ← S
25:        let fragment ← S
26:        let parameters ← [S × S]
27:        let url := ⟨URL, protocol, host, path, parameters, fragment⟩
28:        let req := ⟨HTTPReq, ν₂, GET, host, path, ⟨⟩, parameters, ⟨⟩⟩
29:        call HTTP_SEND(windownonce, req, url, ⊥, ⊥, ⊥, s′)
30:    else if switch ≡ reload then          → Reload some document.
31:        let w̄ ← Subwindows(s′) such that s′.w̄.documents ≠ ⟨⟩
              ↪   if possible; otherwise stop
32:        let url := s′.w̄.activedocument.location
33:        let req := ⟨HTTPReq, ν₂, GET, url.host, url.path, ⟨⟩, url.parameters, ⟨⟩⟩
34:        let referrer := s′.w̄.activedocument.referrer
35:        let s′ := CANCELNAV(s′.w̄.nonce, s′)
36:        call HTTP_SEND(s′.w̄.nonce, req, url, ⊥, referrer, ⊥, s′)
37:    else if switch ≡ forward then
38:        NAVFORWARD(w̄, s′)
39:    else if switch ≡ back then
40:        NAVBACK(w̄, s′)
41: else if m ≡ FULLCORRUPT then          → Request to corrupt browser
42:    let s′.isCorrupted := FULLCORRUPT
43:    stop ⟨⟩, s′
44: else if m ≡ CLOSECORRUPT then          → Close the browser
45:    let s′.secrets := ⟨⟩
46:    let s′.windows := ⟨⟩
47:    let s′.pendingDNS := ⟨⟩
48:    let s′.pendingRequests := ⟨⟩
49:    let s′.sessionStorage := ⟨⟩
50:    let s′.cookies ⊂⟨⟩ Cookies such that
          ↪   (c ∈⟨⟩ s′.cookies) ⟺ (c ∈⟨⟩ s.cookies ∧ c.content.session ≡ ⊥)
51:    let s′.isCorrupted := CLOSECORRUPT
52:    stop ⟨⟩, s′
53: else if ∃⟨reference, request, url, key, f⟩ ∈⟨⟩ s′.pendingRequests
       ↪   such that π₁(dec_s(m, key)) ≡ HTTPResp then      → Encrypted HTTP response
54:    let m′ := dec_s(m, key)
55:    if m′.nonce ≢ request.nonce then
56:        stop
57:    remove ⟨reference, request, url, key, f⟩ from s′.pendingRequests
58:    call PROCESSRESPONSE(m′, reference, request, url, s′)
59: else if π₁(m) ≡ HTTPResp ∧ ∃⟨reference, request, url, ⊥, f⟩ ∈⟨⟩ s′.pendingRequests
       ↪   such that m′.nonce ≡ request.key  then
60:    remove ⟨reference, request, url, ⊥, f⟩ from s′.pendingRequests
61:    call PROCESSRESPONSE(m, reference, request, url, s′)
62: else if m ∈ DNSResponses then          → Successful DNS response
63:    if m.nonce ∉ s.pendingDNS ∨ m.result ∉ IPs ∨ m.domain ≢ π₂(s.pendingDNS).host then
64:        stop
65:    let ⟨reference, message, url, refTB⟩ := s.pendingDNS[m.nonce]
66:    if s′.useTB[url.host] ≡ ⊤ then
67:        let TB_req := ⟨HTTPReq, ν_{n1}, GET, url.host, /OAUTB-prepare, url.parameters, ⟨⟩, ⟨⟩⟩
68:        let s′.tokenBindingRequests := s′.tokenBindingRequests ∪
                           ↪   [ν_{n1}:⟨reference, message, url, refTB⟩]
69:        let s′.pendingRequests := s′.pendingRequests
            ↪   +⟨⟩ ⟨reference, TB_req, url, ν₃, m.result⟩     → TB only with TLS
70:        let message := enc_a(⟨TB_req, ν₃⟩, s′.keyMapping[message.host])
71:    else if url.protocol ≡ S then
72:        let s′.pendingRequests := s′.pendingRequests
            ↪   +⟨⟩ ⟨reference, message, url, ν₃, m.result⟩
73:        let message := enc_a(⟨message, ν₃⟩, s′.keyMapping[message.host])
74:    else
75:        let s′.pendingRequests := s′.pendingRequests
            ↪   +⟨⟩ ⟨reference, message, url, ⊥, m.result⟩
76:    let s′.pendingDNS := s′.pendingDNS − m.nonce
77:    stop ⟨⟨m.result, a, message⟩⟩, s′
78: else          → Some other message
79:    call PROCESS_OTHER(m, a, f, s′)
80: stop
```

The generic HTTPS server is used in the concrete instantiations of clients, authorization servers and resource servers. The placeholder algorithms defined in this section are replaced by the algorithms in the corresponding processes. Here, we use the same model as defined in Appendix E of [44].

In the following, we give the complete definition of the generic HTTPS server, as it is used in the instantiations of the FAPI participants and also referenced within the proof.

*Definition 3 (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: $pendingDNS \in \left[ \mathcal{N} \times \mathcal{T}_{\mathcal{N}} \right]$, $pendingRequests \in \left[ \mathcal{N} \times \mathcal{T}_{\mathcal{N}} \right]$ (both containing arbitrary terms), $DNSaddress \in \mathsf{IPs}$ (containing the IP address of a DNS server), $keyMapping \in \left[ \mathsf{Doms} \times \mathcal{T}_{\mathcal{N}} \right]$ (containing a mapping from domains to public keys), $tlskeys \in \left[ \mathsf{Doms} \times \mathcal{N} \right]$ (containing a mapping from domains to private keys) and $corrupt \in \mathcal{T}_{\mathcal{N}}$ (either $\bot$ if the server is not corrupted, or an arbitrary term otherwise).

Table I shows a list of placeholders for nonces used in these algorithms.

| Placeholder | Usage |
| --- | --- |
| $v_0$ | new nonce for DNS requests |
| $v_1$ | new symmetric key |

Table I: List of placeholders used in the generic HTTPS server algorithm.

We now define the default functions of the generic web server in Algorithms 6–10, and the main relation in Algorithm 11.

---

**Algorithm 6** Generic HTTPS Server Model: Sending a DNS message (in preparation for sending an HTTPS message).

1: **function** HTTPS_SIMPLE_SEND(*reference*, *message*, $s'$)
2:     **let** $s'.\mathrm{pendingDNS}[v_0] := \langle reference, message \rangle$
3:     **stop** $\langle \langle s'.\mathrm{DNSaddress}, a, \langle \mathrm{DNSResolve}, message.\mathrm{host}, v_0 \rangle \rangle \rangle$, $s'$

---

**Algorithm 7** Generic HTTPS Server Model: Default HTTPS response handler.

1: **function** PROCESS_HTTPS_RESPONSE(*m*, *reference*, *request*, *key*, *a*, *f*, $s'$)
2:     **stop**

---

**Algorithm 8** Generic HTTPS Server Model: Default trigger event handler.

1: **function** TRIGGER($s'$)
2:     **stop**

---

**Algorithm 9** Generic HTTPS Server Model: Default HTTPS request handler.

1: **function** PROCESS_HTTPS_REQUEST($m$, $k$, $a$, $f$, $s'$)
2:    **stop**

---

**Algorithm 10** Generic HTTPS Server Model: Default handler for other messages.

1: **function** PROCESS_OTHER($m$, $a$, $f$, $s'$)
2:    **stop**

---

**Algorithm 11** Generic HTTPS Server Model: Main relation of a generic HTTPS server

**Input:** $\langle a, f, m \rangle, s$
1: **let** $s' := s$
2: **if** $s'.\texttt{corrupt} \not\equiv \bot \vee m \equiv \texttt{CORRUPT}$ **then**
3:     **let** $s'.\texttt{corrupt} := \langle \langle a, f, m \rangle, s'.\texttt{corrupt} \rangle$
4:     **let** $m' \leftarrow d_V(s')$
5:     **let** $a' \leftarrow \mathsf{IPs}$
6:     **stop** $\langle \langle a', a, m' \rangle \rangle, s'$
7: **if** $\exists m_{\text{dec}}, k, k', inDomain$ **such that** $\langle m_{\text{dec}}, k \rangle \equiv \mathsf{dec_a}(m, k') \wedge \langle inDomain, k' \rangle \in s.\texttt{tlskeys}$ **then**
8:     **let** $n$, *method*, *path*, *parameters*, *headers*, *body* **such that**
     $\hookrightarrow \langle \texttt{HTTPReq}, n, method, inDomain, path, parameters, headers, body \rangle \equiv m_{\text{dec}}$
     $\hookrightarrow$ **if possible; otherwise stop**
9:     **call** PROCESS_HTTPS_REQUEST($m_{\text{dec}}$, $k$, $a$, $f$, $s'$)
10: **else if** $m \in \mathsf{DNSResponses}$ **then**     $\rightarrow$ Successful DNS response
11:     **if** $m.\texttt{nonce} \notin s.\texttt{pendingDNS} \vee m.\texttt{result} \notin \mathsf{IPs} \vee m.\texttt{domain} \not\equiv s.\texttt{pendingDNS}[m.\texttt{nonce}].2.\texttt{host}$ **then**
12:         **stop**
13:     **let** $\langle reference, request \rangle := s.\texttt{pendingDNS}[m.\texttt{nonce}]$
14:     **let** $s'.\texttt{pendingRequests} := s'.\texttt{pendingRequests}$
         $\hookrightarrow +^{\langle \rangle} \langle reference, request, \nu_1, m.\texttt{result} \rangle$
15:     **let** *message* $:= \mathsf{enc_a}(\langle request, \nu_1 \rangle, s'.\texttt{keyMapping}[request.\texttt{host}])$
16:     **let** $s'.\texttt{pendingDNS} := s'.\texttt{pendingDNS} - m.\texttt{nonce}$
17:     **stop** $\langle \langle m.\texttt{result}, a, message \rangle \rangle, s'$
18: **else if** $\exists \langle reference, request, key, f \rangle \in^{\langle \rangle} s'.\texttt{pendingRequests}$
     $\hookrightarrow$ **such that** $\pi_1(\mathsf{dec_s}(m, key)) \equiv \texttt{HTTPResp}$ **then**     $\rightarrow$ Encrypted HTTP response
19:     **let** $m' := \mathsf{dec_s}(m, key)$
20:     **if** $m'.\texttt{nonce} \not\equiv request.\texttt{nonce}$ **then**
21:         **stop**
22:     **remove** $\langle reference, request, key, f \rangle$ **from** $s'.\texttt{pendingRequests}$
23:     **call** PROCESS_HTTPS_RESPONSE($m'$, *reference*, *request*, *key*, $a$, $f$, $s'$)
24:     **stop**
25: **else if** $m \equiv \texttt{TRIGGER}$ **then**     $\rightarrow$ Process was triggered
26:     **call** PROCESS_TRIGGER($s'$)
27: **stop**

In the following, we will give the formal application-specific model of the participants of the FAPI, i.e., the clients, authorization servers and resource servers. Then, building upon these models, we will give the definition of a FAPI web system with a network attacker.

## A. Clients

Similar to Section H of Appendix F of [44], a client $c \in \mathsf{C}$ is a web server modeled as an atomic DY process $(I^c, Z^c, R^c, s_0^c)$ with the addresses $I^c := \mathsf{addr}(c)$.

*Definition 4.* A *state* $s \in Z^c$ *of client* $c$ is a term of the form $\langle$*DNSaddress, pendingDNS, pendingRequests, corrupt, keyMapping, tlskeys, sessions, issuerCache, oidcConfigCache, jwksCache, clientCredentialsCache, oautbEKM, authReqSigKey, tokenBindings*$\rangle$ with *DNSaddress* $\in \mathsf{IPs}$, *pendingDNS* $\in \left[\mathcal{N} \times \mathcal{T}_\mathcal{N}\right]$, *pendingRequests* $\in \left[\mathcal{N} \times \mathcal{T}_\mathcal{N}\right]$, *corrupt* $\in \mathcal{T}_\mathcal{N}$, *keyMapping* $\in \left[\mathsf{Doms} \times \mathcal{T}_\mathcal{N}\right]$, *tlskeys* $\in \left[\mathsf{Doms} \times K_{\mathrm{TLS}}\right]$ (all former components as in Definition 3), *sessions* $\in \left[\mathcal{N} \times \mathcal{T}_\mathcal{N}\right]$, *issuerCache* $\in \left[\mathcal{T}_\mathcal{N} \times \mathcal{T}_\mathcal{N}\right]$, *oidcConfigCache* $\in \left[\mathcal{T}_\mathcal{N} \times \mathcal{T}_\mathcal{N}\right]$, *jwksCache* $\in \left[\mathcal{T}_\mathcal{N} \times \mathcal{T}_\mathcal{N}\right]$, *clientCredentialsCache* $\in \left[\mathcal{T}_\mathcal{N} \times \mathcal{T}_\mathcal{N}\right]$, *oautbEKM* $\in \mathcal{T}_\mathcal{N}$, *authReqSigKey* $\in \mathcal{N}$ and *tokenBindings* $\in [\mathsf{Doms} \times \mathcal{N}]$.

An *initial state* $s_0^c$ *of* $c$ is a state of $c$ with $s_0^c.\mathtt{pendingDNS} \equiv \langle\rangle$, $s_0^c.\mathtt{pendingRequests} \equiv \langle\rangle$, $s_0^c.\mathtt{corrupt} \equiv \bot$, $s_0^c.\mathtt{keyMapping}$ being the same as the keymapping for browsers, $s_0^c.\mathtt{tlskeys} \equiv tlskeys^c$, $s_0^c.\mathtt{sessions} \equiv \langle\rangle$ and $s_0^c.\mathtt{oautbEKM} \equiv \langle\rangle$.

We require the initial state to contain preconfigured values for $s_0^c.\mathtt{issuerCache}$ (mapping from identities to domains of the corresponding authorization server), $s_0^c.\mathtt{oidcConfigCache}$ (authorization endpoint and token endpoint for each domain of an authorization server) and $s_0^c.\mathtt{jwksCache}$ (mapping from issuers to the public key used to validate their signatures). More precisely, we require that for all authorization servers *as* and all domains $dom_{as} \in \mathsf{dom}(as)$ it holds true that

$$s_0^c.\mathtt{oidcConfigCache}[dom_{as}][\mathtt{token\_ep}] \equiv \langle \mathtt{URL}, \mathsf{S}, dom'_{as}, /\mathtt{token}, \langle\rangle, \langle\rangle\rangle \tag{4}$$

$$s_0^c.\mathtt{oidcConfigCache}[dom_{as}][\mathtt{auth\_ep}] \equiv \langle \mathtt{URL}, \mathsf{S}, dom'_{as}, /\mathtt{auth}, \langle\rangle, \langle\rangle\rangle \tag{5}$$

$$s_0^c.\mathtt{jwksCache}[dom_{as}] \equiv \mathsf{pub}(s_0^{as}.\mathtt{jwk}) \tag{6}$$

$$s_0^c.\mathtt{issuerCache}[id] \in \mathsf{dom}(as) \Leftrightarrow id \in \mathsf{ID}^{as}, \tag{7}$$

for $dom'_{as} \in \mathsf{dom}(as)$. These properties hold true if OIDC Discovery [51] and Dynamic Client Registration [52] are used, as already shown in Lemmas 1 to 4 of [44].

Furthermore, $s_0^c.\mathtt{oidcConfigCache}[issuer][\mathtt{resource\_servers}] = s_0^{as}.\mathtt{resource\_servers}$ shall contain the domains of all resource servers that the authorization server *issuer* supports.

For signing authorization requests, we require that $s_0^c.\mathtt{authReqSigKey}$ contains a key that is only known to $c$.

For each identity $id \in \mathsf{ID}^{as}$ governed by authorization server *as*, the credentials of the client are stored in $s_0^c.\mathtt{clientCredentialsCache}[s_0^c.\mathtt{issuerCache}[id]]$, which shall be equal to $s_0^{as}.\mathtt{clients}[clientId]$, where *clientId* is the client id that was issued to $c$ by *as*.

For OAuth 2.0 Token Binding, we require that $s_0^c.\mathtt{tokenBindings}[d]$ contains a different nonce for each $d \in \mathsf{Doms}$.

The relation $R^c$ is based on the model of generic HTTPS servers (see Section H). The algorithms that differ from or do not exist in the generic server model are defined in Algorithms 12–20. Table II shows a list of all placeholders used in these algorithms.

The scripts that are used by the clients are described in Algorithms 22 and 23. As in [44], the current URL of a document is extracted by using the function GETURL(*tree, docnonce*), which searches for the document with the identifier *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 *docnonce* is found in the tree *tree*, $\diamond$ is returned.

Furthermore, we require that *leak* $\in \mathsf{IPs}$ is an arbitrary IP address. By sending messages to this address, we model the leakage of tokens or messages, as this address can be the address of an attacker.

### Description

In the following, we will describe the basic functionality of the algorithms, focusing on the main differences to the algorithms used in [44].

- Algorithm 12 handles incoming requests. The browser can request a nonce for OAUTB in the `OAUTB-prepare` path. The redirection endpoint requires and checks the corresponding OAUTB header if the client is a web server client using OAUTB. In case of a read-write client, the authorization response contains an id token, which the client checks in Algorithm 19 before continuing the flow.
- Algorithm 13 handles responses. If the client is a read-write client, it checks the second id token contained in the token request before continuing the flow. Furthermore, when using mTLS or OAUTB, the client processes the initial response

| Placeholder | Usage |
|---|---|
| $v_1$ | new login session id |
| $v_2$ | new HTTP request nonce |
| $v_3$ | new HTTP request nonce |
| $v_4$ | new service session id |
| $v_5$ | new HTTP request nonce |
| $v_6$ | new state value |
| $v_7$ | new *nonce* value |
| $v_8$ | new nonce for OAUTB |
| $v_9$ | new PKCE verifier |
| $v_{10}$ | new HTTP request nonce |
| $v_{11}$ | new HTTP request nonce |

Table II: List of placeholders used in the client algorithm.

needed for the actual request that contains either the decrypted mTLS nonce or a signed EKM value. When receiving a resource access nonce, the web server client sends the nonce to the process that sent the request to the redirection endpoint.

- Algorithm 14 starts a new login flow. The client prepares the authorization request, which is redirected by the browser. The request differs depending on the client type and whether the client is an app client or a web server client. In all cases, it creates a signed request JWS for preventing the attack described in Section IV-C. Furthermore, the request leaks at this point, modeling leaking browser logs. The resource server used at the end of the flow is chosen from the set of resource servers provided by the authorization server.
- Algorithms 15 and 16 handle the token request. If the client uses mTLS or OAUTB, it sends an initial request with Algorithm 15. For modeling a misconfigured token endpoint, the read-write client may choose between the preconfigured token endpoint or an arbitrary URL.
- The request to the resource server is made similarly in Algorithms 17 and 18.
- Algorithm 19 checks the first id token and continues the flow, whereas Algorithm 20 uses the id token for logging in the identity.

---

**Algorithm 12** 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.

3:        **let** $headers := [\texttt{ReferrerPolicy:origin}]$    → Set the Referrer Policy for the index page of the client.

4:        **let** $m' := \texttt{enc}_s(\langle \texttt{HTTPResp}, m.\texttt{nonce}, 200, headers, \langle \texttt{script\_client\_index}, \langle\rangle\rangle\rangle, k)$
          → Send *script_client_index* in HTTP response.

5:        **stop** $\langle\langle f, a, m'\rangle\rangle$, $s'$

6:

7:    **else if** $m$.path $\equiv$ /startLogin $\wedge$ $m$.method $\equiv$ POST **then**    → **Serve the request to start a new login.**

8:        **if** $m$.headers$[\texttt{Origin}] \not\equiv \langle m.\texttt{host}, \texttt{S}\rangle$ **then**

9:            **stop**    → Check the Origin header for CSRF protection.

10:        **let** $id := m$.body

11:        **let** $sessionId := v_1$    → Session id is a freshly chosen nonce.

12:        **let** $s'.\texttt{sessions}[sessionId] := [\texttt{startRequest}:[\texttt{message}:m, \texttt{key}:k, \texttt{receiver}:a, \texttt{sender}:f],$
              $\hookrightarrow \texttt{identity}:id]$    → Create new session record.

13:        **call** START_LOGIN_FLOW($sessionId$, $s'$)    → Call the function that starts a login flow.

14:

15:    **else if** $m$.path $\equiv$ /OAUTB-prepare **then**    → For OAUTB between the user-agent and the client

16:        **let** $headers := [\texttt{ReferrerPolicy:origin}]$

17:        **let** $tbNonce := v_8$

18:        **let** $s'.\texttt{oautbEKM} := s'.\texttt{oautbEKM} +^{\langle\rangle} \texttt{hash}(\langle m.\texttt{nonce}, tbNonce, keyMapping[m.\texttt{host}]\rangle)$

19:        **let** $m' := \texttt{enc}_s(\langle \texttt{HTTPResp}, m.\texttt{nonce}, 200, headers, [\texttt{tbNonce}:tbNonce]\rangle, k)$

20:        **stop** $\langle\langle f, a, m'\rangle\rangle$, $s'$

---

| | |
|---|---|
| 21: | **else if** $m$.path $\equiv$ /redirect_ep **then** → **User is being redirected after authentication to the AS.** |
| 22: | **let** $sessionId := m$.headers[Cookie][$\langle$__Secure, sessionId$\rangle$][value] |
| 23: | **if** $sessionId \notin s'$.sessions **then** |
| 24: | **stop** |
| 25: | **let** $session := s'$.sessions[$sessionId$] → Retrieve session data. |
| 26: | **if** $\langle$S, $m$.host, $m$.path, $m$.parameters, $\bot\rangle \not\equiv session$[redirectUri] **then** |
| 27: | **stop** → Check if response was received at the right redirect uri. |
| 28: | **let** $issuer := s'$.issuerCache[$session$[identity]] → Mappings from identites to issuers. |
| 29: | **let** $credentials := s'$.clientCredentialsCache[$issuer$] |
| 30: | **let** $responseType := session$[response_type] |
| | → Determines the flow to use, e.g., code id_token for an OIDC hybrid flow. |
| 31: | **if** $responseType \in \{\langle$code$\rangle, \langle$JARM_code$\rangle\}$ **then** → Authorization code mode: Take data from URL parameters. |
| 32: | **let** $data := m$.parameters |
| 33: | **else** → Hybrid mode: Send script_client_get_fragment to browser to retrieve data from URL fragment |
| 34: | **if** $m$.method $\equiv$ GET **then** |
| 35: | **let** $headers := \langle\langle$ReferrerPolicy, origin$\rangle\rangle$ |
| 36: | **let** $m' := \mathsf{enc_s}(\langle$HTTPResp, $m$.nonce, 200, $headers$, $\langle$script_client_get_fragment, $\bot\rangle\rangle, k)$ |
| 37: | **stop** $\langle\langle f, a, m'\rangle\rangle, s'$ |
| 38: | **else** → POST request: script script_client_get_fragment is sending the data from URL fragment. |
| 39: | **let** $data := m$.body |
| 40: | **if** $credentials$[profile] $\equiv$ rw $\wedge$ $credentials$[client_type] $\equiv$ conf_OAUTB $\wedge$ $credentials$[is_app] $\equiv \bot$ **then** |
| 41: | → Check Provided TB-ID (for PKCE) |
| 42: | **if** Sec-Token-Binding $\notin m$.headers **then** |
| 43: | **stop** |
| 44: | **let** $ekmInfo \leftarrow s'$.oautbEKM |
| 45: | **let** $TB\_Msg\_provided := m$.headers[Sec-Token-Binding][prov] |
| 46: | **let** $TB\_provided\_pub := TB\_Msg\_provided$[id] |
| 47: | **let** $TB\_provided\_sig := TB\_Msg\_provided$[sig] |
| 48: | **if** checksig($TB\_provided\_sig, TB\_provided\_pub) \not\equiv \top$ |
| | $\hookrightarrow \vee$ extractmsg($TB\_provided\_sig) \not\equiv ekmInfo$ **then** |
| 49: | **stop** |
| 50: | **let** $s'$.session[browserTBID] $:= TB\_provided\_pub$ |
| 51: | **let** $s'$.oautbEKM $:= s'$.oautbEKM $- ekmInfo$ |
| 52: | **if** $data$[state] $\not\equiv session$[state] **then** |
| 53: | **stop** → Check $state$ value. |
| 54: | **if** $data$[state] $\equiv \bot$ **then** |
| 55: | **stop** → $state$ value is not valid. |
| 56: | **let** $s'$.sessions[$sessionId$][state] $:= \bot$ → Invalidate state |
| 57: | **let** $s'$.sessions[$sessionId$][redirectEpRequest] $:=$ |
| | $\hookrightarrow$ [message:$m$, key:$k$, receiver:$a$, sender:$f$, data:$data$] → Store incoming request for sending a |
| | resource nonce or the service session |
| | id |
| 58: | **if** $credentials$[profile] $\equiv$ r **then** |
| 59: | **call** PREPARE_TOKEN_REQUEST($sessionId$, $data$[code], $s'$) |
| 60: | **else if** $responseType \equiv \langle$code, id_token$\rangle$ **then** → Check id token |
| 61: | **call** CHECK_FIRST_ID_TOKEN($sessionId$, $data$[id_token], $data$[code], $s'$) |
| 62: | **else** → JARM: Check Response JWS |
| 63: | **call** CHECK_RESPONSE_JWS($sessionId$, $data$[responseJWS], $data$[code], $s'$) |
| 64: | **stop** |

**Algorithm 13** Relation of a Client $R^c$ – Processing HTTPS Responses

1: **function** PROCESS_HTTPS_RESPONSE($m$, *reference*, *request*, *key*, $a$, $f$, $s'$)
2:     **let** *session* := $s'$.sessions[*reference*[session]]
3:     **let** *sessionId* := *reference*[session]
4:     **let** *id* := *session*[identity]
5:     **let** *issuer* := $s'$.issuerCache[*id*]
6:     **let** *profile* := $s'$.clientCredentialsCache[*issuer*][profile]
7:     **let** *isApp* := $s'$.clientCredentialsCache[*issuer*][is_app]
8:     **let** *clientId* := $s'$.clientCredentialsCache[*issuer*][client_id]
9:     **if** *reference*[responseTo] $\equiv$ TOKEN **then**
10:         **if** *session*[scope] $\equiv \langle \rangle \wedge$ *profile* $\equiv$ r **then**
11:             **let** *useAccessTokenNow* := $\top$
12:         **else if** *session*[scope] $\equiv \langle$openid$\rangle \wedge$ *profile* $\equiv$ r **then**
13:             **let** *useAccessTokenNow* $\leftarrow \{\top, \bot\}$
14:         **else if** *profile* $\equiv$ rw $\wedge$ *session*[response_type] $\equiv \langle$code, id_token$\rangle$ **then**     $\rightarrow$ OIDC Hybrid Flow
15:             **let** *firstIdToken* := *session*[redirectEpRequest][data][id_token]
16:             **if** checksig($m$.body[id_token], $s'$.jwksCache[*issuer*]) $\not\equiv \top$ **then**
17:                 **stop**     $\rightarrow$ Check the signature of the id token.
18:             **if** extractmsg($m$.body[id_token])[sub] $\not\equiv$ extractmsg(*firstIdToken*)[sub] **then**
19:                 **stop**     $\rightarrow$ Check if sub is the same as in the first id token ([OIDC], 3.3.3.6).
20:             **if** extractmsg($m$.body[id_token])[iss] $\not\equiv$ extractmsg(*firstIdToken*)[iss] **then**
21:                 **stop**     $\rightarrow$ Check if iss is the same as in the first id token ([OIDC], 3.3.3.6).
22:             **if** extractmsg($m$.body[id_token])[iss] $\not\equiv$ *issuer* **then**
23:                 **stop**     $\rightarrow$ Check the issuer.
24:             **if** extractmsg($m$.body[id_token])[at_hash] $\not\equiv$ hash($m$.body[access_token]) **then**
25:                 **stop**     $\rightarrow$ Check at_hash of second id token (protection against reuse of phised access token).
26:             **if** extractmsg($m$.body[id_token])[aud] $\not\equiv$ *clientId* **then**
27:                 **stop**     $\rightarrow$ Check aud of second id token.
28:             **let** $s'$.sessions[*reference*[session]][idt2_iss] := extractmsg($m$.body[id_token])[iss]
29:             **let** *useAccessTokenNow* $\leftarrow \{\top, \bot\}$
30:         **else if** *profile* $\equiv$ rw $\wedge$ *session*[response_type] $\equiv \langle$JARM_code$\rangle$ **then**     $\rightarrow$ Code flow (JARM)
31:             **let** *requestJWS* := *session*[redirectEpRequest][data][responseJWS]
32:             **if** extractmsg(*requestJWS*)[at_hash] $\not\equiv$ hash($m$.body[access_token]) **then**
33:                 **stop**     $\rightarrow$ Check at_hash of requestJWS (protection against reuse of phised access token).
34:             **if** *session*[scope] $\equiv \langle \rangle$ **then**
35:                 **let** *useAccessTokenNow* := $\top$
36:             **else**
37:                 **let** *useAccessTokenNow* $\leftarrow \{\top, \bot\}$
38:         **if** *useAccessTokenNow* $\equiv \top$ **then**
39:             **call** PREPARE_USE_ACCESS_TOKEN(*reference*[session], $m$.body[access_token], $s'$)
40:         **call** CHECK_ID_TOKEN(*reference*[session], $m$.body[id_token], $s'$)
41:     **else if** *reference*[responseTo] $\equiv$ MTLS_AS **then**
42:         **let** *code* := $s'$.sessions[*sessionId*][code]
43:         **if** $\exists m_{\text{dec}}, k', dom$ **such that** $m_{\text{dec}} \equiv \text{dec}_a(m.\text{body}, k') \wedge \langle dom, k' \rangle \in s'.\text{tlskeys}$ **then**
44:             **let** *mtlsNonce*, *pubKey* **such that** $\langle mtlsNonce, pubKey \rangle \equiv m_{\text{dec}}$ **if possible; otherwise stop**
45:         **if** *pubKey* $\in$ *keyMapping*(*request*.host) **then**     $\rightarrow$ Send nonce only to the process that created it
46:             **call** SEND_TOKEN_REQUEST(*sessionId*, *code*, [type : MTLS, mtls_nonce:*mtlsNonce*], $s'$)
47:         **else**
48:             **stop**
49:     **else if** *reference*[responseTo] $\equiv$ OAUTB_AS **then**
50:         **let** *code* := $s'$.sessions[*sessionId*][code]
51:         **let** *ekm* := hash($\langle m$.nonce, $m$.body[tb_nonce], *keyMapping*[*request*.host]$\rangle$)
52:             $\rightarrow$ Include public key to model Extended Master Secret
53:         **call** SEND_TOKEN_REQUEST(*sessionId*, *code*, [type : OAUTB, ekm:*ekm*], $s'$)
54:     **else if** *reference*[responseTo] $\equiv$ MTLS_RS **then**
55:         **let** *token* := $s'$.sessions[*sessionId*][token]
56:         **if** $\exists m_{\text{dec}}, k', dom$ **such that** $m_{\text{dec}} \equiv \text{dec}_a(m.\text{body}, k') \wedge \langle dom, k' \rangle \in s'.\text{tlskeys}$ **then**
57:             **let** *mtlsNonce*, *pubKey* **such that** $\langle mtlsNonce, pubKey \rangle \equiv m_{\text{dec}}$ **if possible; otherwise stop**
58:         **if** *pubKey* $\in$ *keyMapping*(*request*.host) **then**     $\rightarrow$ Send nonce only to the process that created it
59:             **call** USE_ACCESS_TOKEN(*sessionId*, *token*, [type : MTLS, mtls_nonce:*mtlsNonce*], $s'$)
60:         **else**
61:             **stop**
62:     **else if** *reference*[responseTo] $\equiv$ OAUTB_RS **then**
63:         **let** *token* := $s'$.sessions[*sessionId*][token]
64:         **let** *ekm* := hash($\langle m$.nonce, $m$.body[tb_nonce], *keyMapping*[*request*.host]$\rangle$)
65:         **call** USE_ACCESS_TOKEN(*sessionId*, *token*, [type : OAUTB, ekm:*ekm*], $s'$)

```
66:    else if reference[responseTo] ≡ RESOURCE_USAGE then
67:        let resource := m.body[resource]
68:        let s′.sessions[sessionId][resource] := resource
69:        let request := session[redirectEpRequest]    → Retrieve stored meta data of the request from the browser to
                                                            the redir. endpoint stored in Algorithm 12.
70:        if isApp ≡ ⊥ then    → Send resource access nonce to browser
71:            let headers := [ReferrerPolicy:origin]
72:            let m′ := enc_s(⟨HTTPResp, request[message].nonce, 200, headers, resource⟩, request[key])
73:            stop ⟨⟨request[sender], request[receiver], m′⟩⟩, s′
74:        else    → isApp ≡ ⊤
75:            stop s′
```

---

**Algorithm 14** Relation of a Client $R^c$ – Starting the login flow.

```
1:  function START_LOGIN_FLOW(sessionId, s′)
2:      let session := s′.sessions[sessionId]
3:      let identity := session[identity]
4:      let issuer := s′.issuerCache[identity]
5:      let oidcConfig := s′.oidcConfigCache[issuer]
6:      let authEndpoint := oidcConfig[auth_ep]
7:      let redirectUris := oidcConfig[redirect_uris]    → Set of redirect URIs for the AS.
8:      let resourceServer ← oidcConfig[resource_servers]    → Choose resource server
9:      let s′.sessions[sessionId] := s′.sessions[sessionId] ∪ [RS:resourceServer]
10:     let credentials := s′.clientCredentialsCache[issuer]
11:     let headers := [ReferrerPolicy:origin]
12:     let headers[Set-Cookie] := [sessionId:⟨sessionId, ⊤, ⊤, ⊤⟩]
13:     let profile := credentials[profile]    → either r or rw
14:     let isApp := credentials[is_app]    → either ⊤ or ⊥
15:     let clientType := credentials[client_type]    → pub, conf_JWS, conf_MTLS or conf_OAUTB
16:     if profile ≡ r then
17:         let responseType := ⟨code⟩
18:     else
19:         let responseType ← {⟨code, id_token⟩, ⟨JARM_code⟩}
20:     let redirectUri ← redirectUris    → Auth. response must be received here.
21:     let s′.sessions[sessionId] := s′.sessions[sessionId] ∪ [redirect_uri:redirectUri]
22:     let scope ← {⟨⟩, ⟨openid⟩}
23:     if scope ≡ ⟨openid⟩ ∨ (profile ≡ rw ∧ responseType ≡ ⟨code, id_token⟩) then    → Nonce for obtaining an id token.
24:         let nonce := v_7
25:     else
26:         let nonce := ⟨⟩
27:     if profile ≡ r then
28:         let pkceVerifier := v_9
29:         let pkceChallenge := hash(pkceVerifier)
30:         let s′.sessions[sessionId] := s′.sessions[sessionId] ∪ [pkce_verifier:pkceVerifier]
31:     else if (clientType ≡ pub) ∨ (clientType ≡ conf_OAUTB ∧ isApp ≡ ⊤) then    → OAUTB (app client)
32:         let TB-Id := pub(s′.tokenBindings[authEndpoint.host])
33:         let pkceChallenge := hash(TB-Id)
34:     else if clientType ≡ conf_OAUTB then    → web server client
35:         let pkceChallenge := referred_tb
36:         let headers[Include-Referred-Token-Binding-ID] := ⊤
37:     else    → rw + conf_MTLS
38:         let pkceChallenge := ⟨⟩
39:     let data := [response_type:responseType, redirect_uri:redirectUri, client_id:credentials[client_id],
                        ↪ scope:scope, nonce:nonce, pkce_challenge:pkceChallenge, state:v_6]
40:     let requestJWT := data ∪ [aud:authEndpoint.host]
41:     let requestJWS := sig(requestJWT, s′.authReqSigKey)
42:     let data[request_jws] := requestJWS
43:     let s′.sessions[sessionId] := s′.sessions[sessionId] ∪ data
44:     let authEndpoint.parameters := data
45:     let headers[Location] := authEndpoint
46:     let request := s′.sessions[sessionId][startRequest]
47:     let m′ := enc_s(⟨HTTPResp, request[message].nonce, 303, headers, ⊥⟩, request[key])
48:     let m_leak := ⟨LEAK, authEndpoint⟩
49:     stop ⟨⟨leak, request[receiver], m_leak⟩, ⟨request[sender], request[receiver], m′⟩⟩, s′
```

**Algorithm 15** Relation of a Client $R^c$ – Prepare request to token endpoint.

---

1: **function** PREPARE_TOKEN_REQUEST(*sessionId*, *code*, $s'$)
2:     **let** *session* := $s'$.sessions[*sessionId*]
3:     **let** *identity* := *session*[identity]
4:     **let** *issuer* := $s'$.issuerCache[*identity*]
5:     **let** *credentials* := $s'$.clientCredentialsCache[*issuer*]
6:     **let** *clientId* := *credentials*[client_id]
7:     **let** *clientType* := *credentials*[client_type]
8:     **let** *profile* := *credentials*[profile]
9:     **if** *profile* ≡ rw **then**
10:         **let** *misconfiguredTEp* ← {⊤, ⊥}
11:     **else**
12:         **let** *misconfiguredTEp* := ⊥
13:     **let** $s'$.sessions[*sessionId*][misconfiguredTEp] := *misconfiguredTEp*
14:     **if** *misconfiguredTEp* ≡ ⊤ **then**     → Choose wrong token endpoint.
15:         **let** *host* ← Doms
16:         **let** *path* ← $\mathbb{S}$
17:         **let** *parameters* ← $[\mathbb{S} \times \mathbb{S}]$
18:         **let** *url* := ⟨URL, S, *host*, *path*, *parameters*, ⊥⟩
19:         **let** $s'$.sessions[*sessionId*][token_ep] := *url*
20:     **else**
21:         **let** *url* := $s'$.oidcConfigCache[*issuer*][token_ep]
22:     **let** $s'$.sessions[*sessionId*][code] := *code*
23:     **if** *profile* ≡ r **then**     → Send token request
24:         **if** *clientType* ≡ pub ∨ *clientType* ≡ conf_JWS **then**
25:             **call** SEND_TOKEN_REQUEST(*sessionId*, *code*, ⟨⟩, $s'$)
26:     **if** *clientType* ≡ conf_MTLS **then**     → both profiles
27:         **let** *body* := [client_id:*clientId*]
28:         **let** *message* := ⟨HTTPReq, $\nu_{10}$, GET, *url*.domain, /MTLS-prepare, *url*.parameters, ⟨⟩, *body*⟩
29:         **call** HTTPS_SIMPLE_SEND([responseTo:MTLS_AS, session:*sessionId*], *message*, $s'$)
30:     **else**     → OAUTB
31:         **let** *message* := ⟨HTTPReq, $\nu_{10}$, GET, *url*.domain, /OAUTB-prepare, *url*.parameters, ⟨⟩, ⟨⟩⟩
32:         **call** HTTPS_SIMPLE_SEND([responseTo:OAUTB_AS, session:*sessionId*], *message*, $s'$)

---

**Algorithm 16** Relation of a Client $R^c$ – Request to token endpoint.

---

1: **function** SEND_TOKEN_REQUEST(*sessionId*, *code*, *responseValue*, *s'*)
2:    **let** *session* := *s'*.sessions[*sessionId*]
3:    **let** *identity* := *session*[identity]
4:    **let** *issuer* := *s'*.issuerCache[*identity*]
5:    **if** *session*[misconfiguredTEp] ≡ ⊤ **then**
6:        **let** *url* := *session*[token_ep]
7:    **else**
8:        **let** *url* := *s'*.oidcConfigCache[*issuer*][token_ep]
9:    **let** *credentials* := *s'*.clientCredentialsCache[*issuer*]
10:    **let** *clientId* := *credentials*[client_id]
11:    **let** *clientType* := *credentials*[client_type]
12:    **let** *profile* := *credentials*[profile]
13:    **let** *isApp* := *credentials*[is_app]
14:    **let** *body* := [grant_type:authorization_code, code:*code*,
            ↪    redirect_uri:*session*[redirect_uri], client_id:*clientId*]
15:    **if** *profile* ≡ r **then**
16:        **let** *body*[pkce_verifier] := *session*[pkce_verifier]
17:    **if** *profile* ≡ r ∧ *clientType* ≡ pub **then**
18:        **let** *message* := ⟨HTTPReq, $\nu_2$, POST, *url*.domain, *url*.path, *url*.parameters, ⊥, *body*⟩
19:        **call** HTTPS_SIMPLE_SEND([responseTo:TOKEN, session:*sessionId*], *message*, *s'*)
20:    **else if** *profile* ≡ r ∧ *clientType* ≡ conf_JWS **then**
21:        **let** *clientSecret* := *credentials*[client_secret]
22:        **let** *jwt* := [iss:*clientId*, aud:*url*.domain]
23:        **let** *body*[assertion] := mac(*jwt*, *clientSecret*)
24:        **let** *message* := ⟨HTTPReq, $\nu_2$, POST, *url*.domain, *url*.path, *url*.parameters, ⊥, *body*⟩
25:        **call** HTTPS_SIMPLE_SEND([responseTo:TOKEN, session:*sessionId*], *message*, *s'*)
26:    **else if** *clientType* ≡ conf_MTLS **then**    → both profiles
27:        **if** *responseValue*[type] ≢ MTLS **then**
28:            **stop**
29:        **let** *body*[TLS_AuthN] := *responseValue*[mtls_nonce]
30:        **let** *message* := ⟨HTTPReq, $\nu_2$, POST, *url*.domain, *url*.path, *url*.parameters, ⊥, *body*⟩
31:        **call** HTTPS_SIMPLE_SEND([responseTo:TOKEN, session:*sessionId*], *message*, *s'*)
32:    **else**    → rw with OAUTB
33:        **if** *responseValue*[type] ≢ OAUTB **then**
34:            **stop**
35:        **let** *ekm* := *responseValue*[ekm]
36:        **let** *TB_AS* := *s'*.tokenBindings[*url*.host]    → private key
37:        **let** *TB_RS* := *s'*.tokenBindings[*session*[RS]]    → private key
38:        **let** *TB_Msg_prov* := [id:pub(*TB_AS*), sig:sig(*ekm*, *TB_AS*)]
39:        **let** *TB_Msg_ref* := [id:pub(*TB_RS*), sig:sig(*ekm*, *TB_RS*)]
40:        **let** *headers* := [Sec-Token-Binding:[prov:*TB_Msg_prov*, ref:*TB_Msg_ref*]]
41:        **if** *clientType* ≡ conf_OAUTB **then**    → client authentication
42:            **let** *clientSecret* := *credentials*[client_secret]
43:            **let** *jwt* := [iss:*clientId*, aud:*url*.domain, ]
44:            **let** *body*[assertion] := mac(*jwt*, *clientSecret*)
45:        **if** *isApp* ≡ ⊥ **then**    → web server client: send TBID used by browser
46:            **let** *body*[pkce_verifier] := *session*[browserTBID]
47:        **let** *message* := ⟨HTTPReq, $\nu_2$, POST, *url*.domain, *url*.path, *url*.parameters, *headers*, *body*⟩
48:        **call** HTTPS_SIMPLE_SEND([responseTo:TOKEN, session:*sessionId*], *message*, *s'*)

---

**Algorithm 17** Relation of a Client $R^c$ – Prepare for using the access token.

---

1: **function** PREPARE_USE_ACCESS_TOKEN(*sessionId*, *token*, $s'$)
2:     **let** *session* := $s'$.sessions[*sessionId*]
3:     **let** *identity* := *session*[identity]
4:     **let** *issuer* := $s'$.issuerCache[*identity*]
5:     **let** *credentials* := $s'$.clientCredentialsCache[*issuer*]
6:     **let** *clientType* := *credentials*[client_type]
7:     **let** *profile* := *credentials*[profile]
8:     **let** $s'$.sessions[sessionId][token] := *token*
9:     **let** *rsHost* := *session*[RS]
10:     **if** *profile* $\equiv$ r **then**
11:         **call** USE_ACCESS_TOKEN(*sessionId*, *token*, $\langle\rangle$, $s'$)
12:     **else if** *clientType* $\equiv$ conf_MTLS **then**    $\rightarrow$ rw + mTLS: AT is bound to client
13:         **let** *message* := $\langle$HTTPReq, $\nu_{11}$, GET, *rsHost*, /MTLS-prepare, $\bot$, $\bot$, *body*$\rangle$
14:         **call** HTTPS_SIMPLE_SEND([responseTo:MTLS_RS, session:*sessionId*], *message*, $s'$)
15:     **else**   $\rightarrow$ OAUTB
16:         **let** *message* := $\langle$HTTPReq, $\nu_{11}$, GET, *rsHost*, /OAUTB-prepare, $\bot$, $\bot$, *body*$\rangle$
17:         **call** HTTPS_SIMPLE_SEND([responseTo:OAUTB_RS, session:*sessionId*], *message*, $s'$)

---

**Algorithm 18** Relation of a Client $R^c$ – Using the access token.

---

1: **function** USE_ACCESS_TOKEN(*sessionId*, *token*, *responseValue*, $s'$)
2:     **let** *session* := $s'$.sessions[*sessionId*]
3:     **let** *identity* := *session*[identity]
4:     **let** *issuer* := $s'$.issuerCache[*identity*]
5:     **let** *credentials* := $s'$.clientCredentialsCache[*issuer*]
6:     **let** *clientType* := *credentials*[client_type]
7:     **let** *profile* := *credentials*[profile]
8:     **let** *headers* := [Authorization : $\langle$Bearer, *token*$\rangle$]
9:     **let** *body* := []
10:     **let** *rsHost* := *session*[RS]
11:     **if** *profile* $\equiv$ r **then**
12:         **let** *message* := $\langle$HTTPReq, $\nu_3$, GET, *rsHost*, /resource-r, $\langle\rangle$, *headers*, *body*$\rangle$
13:         **call** HTTPS_SIMPLE_SEND([responseTo:RESOURCE_USAGE, session:*sessionId*], *message*, $s'$)
14:     **else**
15:         **if** *session*[response_type] $\equiv \langle$code, id_token$\rangle$ **then**   $\rightarrow$ OIDC Hybrid Flow
16:             **let** *body*[at_iss] := *session*[idt2_iss]   $\rightarrow$ Send issuer of second id token
17:         **else**   $\rightarrow$ JARM Code Flow
18:             **let** *body*[at_iss] := *session*[JARM_iss]   $\rightarrow$ Send issuer request JWS
19:         **if** *clientType* $\equiv$ conf_MTLS **then**
20:             **if** *responseValue*[type] $\not\equiv$ MTLS **then**
21:                 **stop**
22:             **let** *body*[MTLS_AuthN] := *responseValue*[mtls_nonce]
23:         **else**   $\rightarrow$ OAUTB
24:             **if** *responseValue*[type] $\not\equiv$ OAUTB **then**
25:                 **stop**
26:             **let** *ekm* := *responseValue*[ekm]
27:             **let** *TB_RS* := $s'$.tokenBindings[*rsHost*]
28:             **let** *TB_Msg_prov* := [id:pub(*TB_RS*), sig:sig(*ekm*, *TB_RS*)]
29:             **let** *headers*[Sec-Token-Binding] := [prov:*TB_Msg_prov*]
30:         **let** *message* := $\langle$HTTPReq, $\nu_3$, POST, *rsHost*, /resource-rw, $\langle\rangle$, *headers*, *body*$\rangle$
31:         **call** HTTPS_SIMPLE_SEND([responseTo:RESOURCE_USAGE, session:*sessionId*], *message*, $s'$)

---

**Algorithm 19** Relation of a Client $R^c$ – Check first ID token (without Login).

1: **function** CHECK_FIRST_ID_TOKEN(*sessionId*, *id_token*, *code*, *s′*)    → **Check id token validity.**
2:     **let** *session* := *s′*.sessions[*sessionId*]
3:     **let** *identity* := *session*[identity]
4:     **let** *issuer* := *s′*.issuerCache[*identity*]
5:     **let** *oidcConfig* := *s′*.oidcConfigCache[*issuer*]
6:     **let** *credentials* := *s′*.clientCredentialsCache[*issuer*]
7:     **let** *jwks* := *s′*.jwksCache[*issuer*]
8:     **let** *data* := extractmsg(*id_token*)
9:     **if** *data*[s_hash] $\not\equiv$ hash(*session*[state]) **then**
10:         **stop**    → Check state hash.
11:     **if** *data*[c_hash] $\not\equiv$ hash(*code*) **then**
12:         **stop**    → Check code hash.
13:     **if** *data*[iss] $\not\equiv$ *issuer* $\vee$ *data*[aud] $\not\equiv$ *credentials*[client_id] **then**
14:         **stop**    → Check the issuer and audience.
15:     **if** checksig(*id_token*, *jwks*) $\not\equiv \top$ **then**
16:         **stop**    → Check the signature of the id token.
17:     **if** *data*[nonce] $\not\equiv$ *session*[nonce] **then**
18:         **stop**    → Check nonce.
19:     **call** PREPARE_TOKEN_REQUEST(*sessionId*, *code*, *s′*)

---

**Algorithm 20** Relation of a Client $R^c$ – Check ID token.

1: **function** CHECK_ID_TOKEN(*sessionId*, *id_token*, *s′*)    → **Check id token and create service session.**
2:     **let** *session* := *s′*.sessions[*sessionId*]
3:     **let** *identity* := *session*[identity]
4:     **let** *issuer* := *s′*.issuerCache[*identity*]
5:     **let** *oidcConfig* := *s′*.oidcConfigCache[*issuer*]
6:     **let** *credentials* := *s′*.clientCredentialsCache[*issuer*]
7:     **let** *jwks* := *s′*.jwksCache[*issuer*]
8:     **let** *data* := extractmsg(*id_token*)
9:     **if** *data*[iss] $\not\equiv$ *issuer* $\vee$ *data*[aud] $\not\equiv$ *credentials*[client_id] **then**
10:         **stop**    → Check the issuer and audience.
11:     **if** checksig(*id_token*, *jwks*) $\not\equiv \top$ **then**
12:         **stop**    → Check the signature of the id token.
13:     **if** *data*[nonce] $\not\equiv$ *session*[nonce] **then**
14:         **stop**    → Check nonce.
15:     → User is now logged in. Store user identity and issuer.
16:     **let** *s′*.sessions[*sessionId*][loggedInAs] := $\langle$*issuer*, *data*[sub]$\rangle$
17:     **if** *credentials*[is_app] $\equiv \bot$ **then**    → Send service session id to browser
18:         **let** *s′*.sessions[*sessionId*][serviceSessionId] := $\nu_4$
19:         **let** *request* := *session*[redirectEpRequest]    → Send to sender of request to the redirection endpoint. The request's meta data was stored in PROCESS_HTTPS_REQUEST (Algorithm 12).
20:         **let** *headers* := [ReferrerPolicy:origin]
21:     → Create a cookie containing the service session id.
22:         **let** *headers*[Set-Cookie] := [$\langle$__Secure, serviceSessionId$\rangle$:$\langle \nu_4, \top, \top, \top \rangle$]
23:     → Respond to browser's request to the redirection endpoint.
24:         **let** *m′* := enc$_s$($\langle$HTTPResp, *request*[message].nonce, 200, *headers*, ok$\rangle$, *request*[key])
25:         **stop** $\langle \langle$*request*[sender], *request*[receiver], *m′*$\rangle \rangle$, *s′*
26:     **else**    → app client
27:         **stop** s'

**Algorithm 21** Relation of a Client $R^c$ – Check Response JWS.

1: **function** CHECK_RESPONSE_JWS(*sessionId*, *responseJWS*, *code*, *s'*)    → **Check validity of response JWS.**
2:     **let** *session* := *s'*.sessions[*sessionId*]
3:     **let** *identity* := *session*[identity]
4:     **let** *issuer* := *s'*.issuerCache[*identity*]
5:     **let** *oidcConfig* := *s'*.oidcConfigCache[*issuer*]
6:     **let** *credentials* := *s'*.clientCredentialsCache[*issuer*]
7:     **let** *jwks* := *s'*.jwksCache[*issuer*]
8:     **let** *data* := extractmsg(*responseJWS*)
9:     **if** *data*[state] $\not\equiv$ *session*[state] **then**
10:         **stop**    → Check state hash.
11:     **if** *data*[code] $\not\equiv$ *code* **then**
12:         **stop**    → Check code.
13:     **if** *data*[iss] $\not\equiv$ *issuer* $\vee$ *data*[aud] $\not\equiv$ *credentials*[client_id] **then**
14:         **stop**    → Check the issuer and audience.
15:     **if** checksig(*responseJWS*, *jwks*) $\not\equiv$ $\top$ **then**
16:         **stop**    → Check the signature of the JWS.
17:     **let** *s'*.sessions[*sessionId*][JARM_iss] := *data*[iss]
18:     **call** PREPARE_TOKEN_REQUEST(*sessionId*, *code*, *s'*)

---

**Algorithm 22** Relation of *script_client_index*

**Input:** ⟨*tree*, *docnonce*, *scriptstate*, *scriptinputs*, *cookies*, *localStorage*, *sessionStorage*, *ids*, *secrets*⟩    → **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* $\leftarrow$ {auth, link}    → Non-deterministically decide whether to start a login flow or to follow some link.
2: **if** *switch* $\equiv$ auth **then**    → **Start login flow.**
3:     **let** *url* := GETURL(*tree*, *docnonce*)    → Retrieve own URL.
4:     **let** *id* $\leftarrow$ *ids*    → Retrieve one of user's identities.
5:     **let** *url'* := ⟨URL, S, *url*.host, /startLogin, ⟨⟩, ⟨⟩⟩    → Assemble URL.
6:     **let** *command* := ⟨FORM, *url'*, POST, *id*, $\perp$⟩
            → Post a form including the identity to the client.
7:     **stop** ⟨*s*, *cookies*, *localStorage*, *sessionStorage*, *command*⟩    → Finish script's run and instruct the browser to follow the command (follow post).

8: **else**    → **Follow link.**
9:     **let** *protocol* $\leftarrow$ {P, S}    → Non-deterministically select protocol (HTTP or HTTPS).
10:     **let** *host* $\leftarrow$ Doms    → Non-det. select host.
11:     **let** *path* $\leftarrow$ $\mathbb{S}$    → Non-det. select path.
12:     **let** *fragment* $\leftarrow$ $\mathbb{S}$    → Non-det. select fragment part.
13:     **let** *parameters* $\leftarrow$ [$\mathbb{S} \times \mathbb{S}$]    → Non-det. select parameters.
14:     **let** *url* := ⟨URL, *protocol*, *host*, *path*, *parameters*, *fragment*⟩    → Assemble URL.
15:     **let** *command* := ⟨HREF, *url*, $\perp$, $\perp$⟩    → Follow link to the selected URL.
16:     **stop** ⟨*s*, *cookies*, *localStorage*, *sessionStorage*, *command*⟩    → Finish script's run and instruct the browser to follow the command (follow link).

---

**Algorithm 23** Relation of *script_c_get_fragment*

**Input:** ⟨*tree*, *docnonce*, *scriptstate*, *scriptinputs*, *cookies*, *localStorage*, *sessionStorage*, *ids*, *secrets*⟩
1: **let** *url* := GETURL(*tree*, *docnonce*)
2: **let** *url'* := ⟨URL, S, *url*.host, /redirect_ep, ⟨⟩, ⟨⟩⟩
3: **let** *command* := ⟨FORM, *url'*, POST, *url*.fragment, $\perp$⟩
4: **stop** ⟨*s*, *cookies*, *localStorage*, *sessionStorage*, *command*⟩

## B. Authorization Servers

Similar to Section I of Appendix F of [44], an authorization server $as \in \mathsf{AS}$ is a web server modeled as an atomic process $(I^{as}, Z^{as}, R^{as}, s_0^{as})$ with the addresses $I^{as} := \mathsf{addr}(as)$.

*Definition 5.* A *state* $s \in Z^{as}$ *of an authorization server as* is a term of the form $\langle DNSaddress, pendingDNS, pendingRequests,$ *corrupt, keyMapping, tlskeys, clients* (dict from nonces to terms), *records* (sequence of terms), *jwk* (signing key (only one)), *oautbEKM* (sequence of terms), *accessTokens* (sequence of issued access tokens) $\rangle$ with $DNSaddress \in \mathsf{IPs}$, $pendingDNS \in \left[\mathcal{N} \times \mathcal{T}_{\mathcal{N}}\right]$, $pendingRequests \in \left[\mathcal{N} \times \mathcal{T}_{\mathcal{N}}\right]$, $corrupt \in \mathcal{T}_{\mathcal{N}}$, $keyMapping \in \left[\mathsf{Doms} \times \mathcal{T}_{\mathcal{N}}\right]$, $tlskeys \in \left[\mathsf{Doms} \times K_{\mathrm{TLS}}\right]$ (all former components as in Definition 3), $clients \in \left[\mathcal{N} \times \mathcal{T}_{\mathcal{N}}\right]$, $records \in \mathcal{T}_{\mathcal{N}}$, $jwk \in K_{\mathrm{sign}}$, $oautbEKM \in \mathcal{T}_{\mathcal{N}}$ and $accessTokens \in \mathcal{T}_{\mathcal{N}}$.

An *initial state* $s_0^{as}$ *of as* is a state of *as* with $s_0^{as}.\mathtt{pendingDNS} \equiv \langle\rangle$, $s_0^{as}.\mathtt{pendingRequests} \equiv \langle\rangle$, $s_0^{as}.\mathtt{corrupt} \equiv \bot$, $s_0^{as}.\mathtt{keyMapping}$ being the same as the keymapping for browsers, $s_0^{as}.\mathtt{tlskeys} \equiv tlskeys^{as}$, $s_0^{as}.\mathtt{records} \equiv \langle\rangle$, $s_0^{as}.\mathtt{jwk} \equiv$ $\mathsf{signkey}(as)$, $s_0^c.\mathtt{oautbEKM} \equiv \langle\rangle$ and $s_0^{as}.\mathtt{accessTokens} \equiv \langle\rangle$.

We require $s_0^{as}.\mathtt{clients}[clientId]$ to be preconfigured and to contain information about the client with client id *clientId* regarding the profile (r or rw, accessible via the key $\mathtt{profile}$), the type of the client (pub, conf_JWS, conf_MTLS or conf_OAUTB, accessible via $\mathtt{client\_type}$), the client secret (if the client type is conf_JWS or conf_OAUTB, accessible via $\mathtt{client\_secret}$) and whether the client is an app client or a web server client (either $\top$ or $\bot$, accessible via $\mathtt{is\_app}$).

For checking the signature of signed request JWTs, we require that $s_0^{as}.\mathtt{clients}[clientId][\mathtt{jws\_key}] \equiv \mathsf{pub}(s_0^c.\mathtt{authReqSigKey})$.

The resource servers that are supported by the authorization server shall be contained in $s_0^{as}.\mathtt{resource\_servers} \subset dom(\mathsf{RS})$. For the corresponding resource servers *rs*, we require that $s_0^{rs}.\mathtt{authServ} \in dom(as)$.

The relation $R^{as}$ is based on the model of generic HTTPS servers (see Section H). Algorithm 24 specifies the algorithm responsible for processing HTTPS requests. The script *script_as_form* is defined in Algorithm 25

Furthermore, we require that $leak \in \mathsf{IPs}$ is an arbitrary IP address.

Table III shows a list of all placeholders used in the algorithms of the authorization server.

| Placeholder | Usage |
|---|---|
| $v_1$ | new authorization code |
| $v_2$ | new access token |
| $v_3$ | nonce for mTLS |
| $v_4$ | nonce for OAUTB |

Table III: List of placeholders used in the authorization server algorithm.

**Differences to [44]:** In the /auth2 path, the authorization server requires a signed request JWS in all cases to prevent the attack described in Section IV-C. The authorization server also checks and handles the OAUTB message send by the browser if the client is a confidential client using OAUTB. The authorization response always contains an authorization code. It also contains an id token if the client is a read-write client. The leakage of the authorization response is modeled by sending the response to an arbitrary IP address in case of an app client.

The /token path handles the token request. In case of a read-write client, the access token is bound to the client either via mTLS or OAUTB. To model the leakage of an access token due to phishing, the authorization server sends the access token to an arbitrary IP address in case of a read-write client.

The /MTLS-prepare and /OAUTB-prepare paths send the initial response needed for mTLS and OAUTB.

**Remarks:** To model mTLS, the authorization server sends an encrypted nonce to the client. The client decrypts the message and includes the nonce in the actual request. This way, the client proves possession of the corresponding key.

As we do not model public key infrastructures, the client does not send a certificate to the authorization server (as specified in Section 7.4.6 of [47]). In general, the model uses the function *keyMapping*, which maps domains to the public key that would be contained in the certificate.

We require that the authorization server has a mapping from client ids to the corresponding public keys in $s_0^{as}.\mathtt{clients}[clientId][\mathtt{mtls\_key}]$ (only if the client type is conf_MTLS), with the value $keyMapping[dom_c]$ with $dom_c \in \mathsf{dom}(c)$.

As explained in Section F-C, we do not model access tokens being issued from the authorization endpoint.

**Algorithm 24** Relation of AS $R^{as}$ – Processing HTTPS Requests

1: **function** PROCESS_HTTPS_REQUEST($m$, $k$, $a$, $f$, $s'$)
2:     **if** $m$.path $\equiv$ /auth **then**
3:         **if** $m$.method $\equiv$ GET **then**
4:             **let** $data := m$.parameters
5:         **let** $m' := \mathsf{enc_s}(\langle\mathtt{HTTPResp}, m.\mathtt{nonce}, 200, \langle\langle\mathtt{ReferrerPolicy}, \mathtt{origin}\rangle\rangle, \langle\mathtt{script\_as\_form}, data\rangle\rangle, k)$
6:         **stop** $\langle\langle f, a, m'\rangle\rangle$, $s'$
7:
8:     **else if** $m$.path $\equiv$ /auth2 $\wedge$ $m$.method $\equiv$ POST $\wedge$ $m$.headers[Origin] $\equiv \langle m.\mathtt{host}, \mathtt{S}\rangle$ **then**
9:         **let** $identity := m$.body[identity]
10:        **let** $password := m$.body[password]
11:        **if** $identity$.domain $\notin \mathrm{dom}(as)$ **then**
12:            **stop**
13:        **if** $password \not\equiv \mathsf{secretOfID}(identity)$ **then**
14:            **stop**
15:        **let** $clientId := m$.body[client_id]
16:        **if** $clientId \notin s'$.clients **then**
17:            **stop**
18:        **let** $clientInfo := s'$.clients[$clientId$]
19:        **let** $profile := clientInfo$[profile]
20:        **let** $clientType := clientInfo$[client_type]
21:        **let** $isApp := clientInfo$[is_app]
22:        **if** request_jws $\in m$.body **then**        $\rightarrow$ Request must be as JWS
23:            **let** $requestJWS := m$.body[request_jws]
24:            **if** $\mathsf{checksig}(requestJWS, clientInfo[\mathtt{jws\_key}]) \not\equiv \top$ **then**
25:                **stop**    $\rightarrow$ wrong signature
26:            **let** $requestData := \mathsf{extractmsg}(requestJWS)$
27:            **if** $requestData$[aud] $\not\equiv m$.host **then**
28:                **stop**    $\rightarrow$ wrong audience
29:            **if** $requestData$[client_id] $\not\equiv clientId$ **then**
30:                **stop**    $\rightarrow$ clientId not the same as in body
31:        **else**
32:            **stop**
33:        **if** $profile \equiv \mathtt{rw} \wedge clientType \equiv \mathtt{conf\_OAUTB} \wedge isApp \equiv \bot$ **then**    $\rightarrow$ Check both Token Binding messages
34:            **if** Sec-Token-Binding $\notin m$.headers **then**
35:                **stop**
36:            **let** $ekmInfo \leftarrow s'$.oautbEKM
37:            **let** $TB\_Msg\_provided := m$.headers[Sec-Token-Binding][prov]
38:            **let** $TB\_Msg\_referred := m$.headers[Sec-Token-Binding][ref]
39:            **let** $TB\_provided\_pub := TB\_Msg\_provided$[id]
40:            **let** $TB\_provided\_sig := TB\_Msg\_provided$[sig]
41:            **let** $TB\_referred\_pub := TB\_Msg\_referred$[id]
42:            **let** $TB\_referred\_sig := TB\_Msg\_referred$[sig]
43:            **if** $\mathsf{checksig}(TB\_provided\_sig, TB\_provided\_pub) \not\equiv \top$ **then**
44:                **stop**
45:            **if** $\mathsf{extractmsg}(TB\_provided\_sig) \not\equiv ekmInfo$ **then**
46:                **stop**
47:            **if** $\mathsf{checksig}(TB\_referred\_sig, TB\_referred\_pub) \not\equiv \top$ **then**
48:                **stop**
49:            **if** $\mathsf{extractmsg}(TB\_referred\_sig) \not\equiv ekmInfo$ **then**
50:                **stop**
51:            **let** $s'$.oautbEKM $:= s'$.oautbEKM $- ekmInfo$
52:        **let** $responseType := requestData$[response_type]
53:        **let** $redirectUri := requestData$[redirect_uri]
54:        **let** $state := requestData$[state]
55:        **let** $nonce := requestData$[nonce]
56:        **if** $state \equiv \langle\rangle$ **then**
57:            **stop**    $\rightarrow$ state must be included
58:        **if** $redirectUri \notin^{\langle\rangle} clientInfo$[redirect_uris] **then**
59:            **stop**
60:        **let** $record := [\mathtt{client\_id} : clientId]$    $\rightarrow$ Save data in record
61:        **let** $record$[redirect_uri] $:= redirectUri$
62:        **let** $record$[subject] $:= identity$
63:        **let** $record$[issuer] $:= m$.host

```
64:          let record[nonce] := nonce
65:          let record[scope] := requestData[scope]
66:          let record[response_type] := responseType
67:          let record[code] := ν₁
68:          let record[access_token] := ν₂      → Access token for token request
69:          if profile ≡ r ∨ clientType ≡ pub ∨ (clientType ≡ conf_OAUTB ∧ isApp ≡ ⊤) then      → Save PKCE challenge (unless
                                                                                                      in rw profile with mTLS)
70:              let record[pkce_challenge] := requestData[pkce_challenge]
71:          else if clientType ≡ conf_OAUTB ∧ isApp ≡ ⊥ then
72:              let record[pkce_challenge] := TB_referred_pub
73:          let s′.records := s′.records +⟨⟩ record
74:          let responseData := [code:ν₁]      → Always send code
75:          if (profile ≡ rw ∧ responseType ∉ {⟨code, id_token⟩, ⟨JARM_code⟩}) ∨ (profile ≡ r ∧ responseType ≢ ⟨code⟩) then
76:              stop
77:          if responseType ≡ ⟨JARM_code⟩ then
78:              let responseJWT := [iss : record[issuer], aud : record[client_id], code:record[code],
                        ↪    at_hash:hash(record[access_token]), state:state]
79:              let responseData[responseJWS] := sig(responseJWT, s′.jwk)
80:          if id_token ∈⟨⟩ responseType then
81:              let idTokenBody := [iss : record[issuer], sub : record[subject],
                        ↪    aud : record[client_id], nonce : record[nonce],
                        ↪    c_hash:hash(record[code]), s_hash:hash(state)]
82:              let responseData[id_token] := sig(idTokenBody, s′.jwk)
83:          let responseData[state] := state
84:          if responseType ∈ {⟨code⟩, ⟨JARM_code⟩} then      → Authorization code mode
85:              let redirectUri.parameters := redirectUri.parameters ∪ responseData
86:          else   → Hybrid Mode
87:              let redirectUri.fragment := redirectUri.fragment ∪ responseData
88:          let m′ := encₛ(⟨HTTPResp, m.nonce, 303, ⟨⟨Location, redirectUri⟩⟩, ⟨⟩⟩, k)
89:          if clientInfo[is_app] ≡ ⊤ then      → Leakage of authorization response
90:              stop ⟨⟨leak, a, ⟨LEAK, clientId, ⟨Location, redirectUri⟩⟩⟩, ⟨f, a, m′⟩⟩, s′
91:          else   → is_app ≡ ⊥
92:              stop ⟨⟨f, a, m′⟩⟩, s′
93:      else if m.path ≡ /token ∧ m.method ≡ POST then
94:          let clientId := m.body[client_id]
95:          let code := m.body[code]
96:          let record, ptr such that record ≡ s′.records.ptr ∧ record[code] ≡ code
                        ↪    ∧code ≢ ⊥ if possible; otherwise stop
97:          if record[client_id] ≢ clientId then
98:              stop
99:          let clientInfo := s′.clients[clientId]
100:         let profile := clientInfo[profile]
101:         let clientType := clientInfo[client_type]
102:         if profile ≡ rw ∧ (clientType ≡ pub ∨ clientType ≡ conf_OAUTB) then      → Check   both   Token
                                                                                              Binding messages
103:             let ekmInfo ← s′.oautbEKM
104:             let TB_Msg_provided := m.headers[Sec-Token-Binding][prov]
105:             let TB_provided_pub := TB_Msg_provided[id]
106:             let TB_provided_sig := TB_Msg_provided[sig]
107:             if checksig(TB_provided_sig, TB_provided_pub) ≢ ⊤
                        ↪    ∨extractmsg(TB_provided_sig) ≢ ekmInfo then
108:                 stop      → Wrong signature or ekm value
109:             let TB_Msg_referred := m.headers[Sec-Token-Binding][ref]
110:             let TB_referred_pub := TB_Msg_referred[id]
111:             let TB_referred_sig := TB_Msg_referred[sig]
112:             if checksig(TB_referred_sig, TB_referred_pub) ≢ ⊤
                        ↪    ∨extractmsg(TB_referred_sig) ≢ ekmInfo then
113:                 stop      → Wrong signature or ekm value
114:             let s′.oautbEKM := s′.oautbEKM − ekmInfo
115:         if clientType ≡ conf_JWS ∨ clientType ≡ conf_OAUTB then      → Check JWS
116:             let clientSecret := clientInfo[client_secret]
117:             if checkmac(m.body[assertion], clientSecret) ≢ ⊤ then
118:                 stop      → Invalid MAC
119:             let assertion := extractmsg(m.body[assertion])
120:             if assertion[aud] ≢ m.host ∨ assertion[iss] ≢ clientId then
121:                 stop      → Invalid audience or clientId
```

```
122:      else if clientType ≡ conf_MTLS then
123:            let mtlsInfo ← s'.mtlsRequests[clientId]
124:            if mtlsInfo.1 ≢ m.body[TLS_AuthN] then
125:                 stop
126:            let s'.mtlsRequests[clientId] := s'.mtlsRequests[clientId] − mtlsInfo
127:      if profile ≡ r then
128:            if hash(m.body[pkce_verifier]) ≢ record[pkce_challenge] then
129:                 stop
130:      else if profile ≡ rw ∧ clientType ≡ conf_OAUTB ∧ isApp ≡ ⊥ then
131:            if m.body[pkce_verifier] ≢ record[pkce_challenge] then      → Sec. 5.2 of [OAUTB]
132:                 stop
133:      else if profile ≡ rw ∧ (clientType ≡ pub ∨ (clientType ≡ conf_OAUTB ∧ isApp ≡ ⊤)) then
134:            if hash(TB_provided_pub) ≢ record[pkce_challenge] then      → Sec. 5.1 of [OAUTB]
135:                 stop
136:      if not (record[redirect_uri] ≡ m.body[redirect_uri] ∨
             ↪ (|clientInfo[redirect_uris]| = 1 ∧ redirect_uri ∉ m.body)) then
137:            stop      → If only one redirect URI is registered, it can be omitted.
138:      let s'.records.ptr[code] := ⊥      → Invalidate code
139:      let accessToken := record[access_token]
140:      if profile ≡ rw then      → Create token binding
141:            if clientType ≡ conf_MTLS then
142:                 let s'.accessTokens := s'.accessTokens+⟨⟩
                      ↪ ⟨MTLS, record[subject], clientId, accessToken, mtlsInfo.2, rw⟩
143:            else      → OAUTB
144:                 let s'.accessTokens := s'.accessTokens+⟨⟩
                      ↪ ⟨OAUTB, record[subject], clientId, accessToken, TB_referred_pub, rw⟩
145:      else
146:            let s'.accessTokens := s'.accessTokens +⟨⟩ ⟨record[subject], clientId, accessToken, r⟩
147:      if openid ∈⟨⟩ record[scope] ∨ (record[response_type] ≡ ⟨code, id_token⟩) then      → return id token
148:            let idTokenBody := [iss : record[issuer]]
149:            let idTokenBody[sub] := record[subject]
150:            let idTokenBody[aud] := record[client_id]
151:            let idTokenBody[nonce] := record[nonce]
152:            let idTokenBody[at_hash] := hash(accessToken)      → Mitigate reuse of phished AT
153:            let idToken := sig(idTokenBody, s'.jwk)
154:            let m' := enc_s(⟨HTTPResp, m.nonce, 200, ⟨⟩, [access_token:accessToken, id_token:idToken]⟩, k)
155:      else
156:            let m' := enc_s(⟨HTTPResp, m.nonce, 200, ⟨⟩, [access_token:accessToken]⟩, k)
157:      if profile ≡ rw then
158:            stop ⟨⟨leak, a, ⟨LEAK, clientId, accessToken⟩⟩, ⟨f, a, m'⟩⟩, s'      → Leakage of access token
159:      else      → profile ≡ r
160:            stop ⟨⟨f, a, m'⟩⟩, s'
161:
162:   else if m.path ≡ /MTLS-prepare then
163:      let clientId := m.body[client_id]
164:      let mtlsNonce := ν₃
165:      let clientKey := s'.clients[clientId][mtls_key]
166:      let s'.mtlsRequests[clientId] := s'.mtlsRequests[clientId] +⟨⟩ ⟨mtlsNonce, clientKey⟩
167:      let m' := enc_s(⟨HTTPResp, m.nonce, 200, ⟨⟩, enc_a(⟨mtlsNonce, keyMapping(m.host)⟩, clientKey)⟩, k)
168:      stop ⟨⟨f, a, m'⟩⟩, s'
169:
170:   else if m.path ≡ /OAUTB-prepare then
171:      let tbNonce := ν₄
172:      let s'.oautbEKM := s'.oautbEKM +⟨⟩ hash(⟨m.nonce, tbNonce, keyMapping[m.host]⟩)
173:            → Own public key is needed for modeling Extended Master Secret
174:      let m' := enc_s(⟨HTTPResp, m.nonce, 200, ⟨⟩, [tb_nonce:tbNonce]⟩, k)
175:      stop ⟨⟨f, a, m'⟩⟩, s'
```

---

**Algorithm 25** Relation of *script_as_form*

---

**Input:** ⟨*tree*, *docnonce*, *scriptstate*, *scriptinputs*, *cookies*, *localStorage*, *sessionStorage*, *ids*, *secrets*⟩
1: **let** *url* := GETURL(*tree*, *docnonce*)
2: **let** *url'* := ⟨URL, S, *url*.host, /auth2, ⟨⟩, ⟨⟩⟩
3: **let** *formData* := *scriptstate*
4: **let** *identity* ← *ids*
5: **let** *secret* ← *secrets*
6: **let** *formData*[identity] := *identity*
7: **let** *formData*[password] := *secret*
8: **let** *command* := ⟨FORM, *url'*, POST, *formData*, ⊥⟩
9: **stop** ⟨*s*, *cookies*, *localStorage*, *sessionStorage*, *command*⟩

---

## C. 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 6.* A *state* $s \in Z^{rs}$ *of a resource server rs* is a term of the form $\langle DNSaddress, pendingDNS, pendingRequests, corrupt, keyMapping, tlskeys, mtlsRequests$ (sequence of terms), *oautbEKM* (sequence of terms), *rNonce* (dict from ID to sequence of nonces), *wNonce* (dict from ID to sequence of nonces), *ids* (sequence of ids), *authServ* (domain)$\rangle$ with $DNSaddress \in \text{IPs}$, $pendingDNS \in \left[ \mathcal{N} \times \mathcal{T}_\mathcal{N} \right]$, $pendingRequests \in \left[ \mathcal{N} \times \mathcal{T}_\mathcal{N} \right]$, $corrupt \in \mathcal{T}_\mathcal{N}$, $keyMapping \in \left[ \text{Doms} \times \mathcal{T}_\mathcal{N} \right]$, $tlskeys \in \left[ \text{Doms} \times K_{\text{TLS}} \right]$ (all former components as in Definition 3), $mtlsRequests \in \mathcal{T}_\mathcal{N}$, $oautbEKM \in \mathcal{T}_\mathcal{N}$, $rNonce \in \left[ \text{ID} \times N_r \right]$, $wNonce \in \left[ \text{ID} \times N_w \right]$, $ids \subset \text{ID}$ and $authServ \in \text{Doms}$.

An *initial state* $s_0^{rs}$ *of rs* is a state of $rs$ with $s_0^{rs}.\texttt{pendingDNS} \equiv \langle \rangle$, $s_0^{rs}.\texttt{pendingRequests} \equiv \langle \rangle$, $s_0^{rs}.\texttt{corrupt} \equiv \perp$, $s_0^{rs}.\texttt{keyMapping}$ being the same as the keymapping for browsers, $s_0^{rs}.\texttt{tlskeys} \equiv tlskeys^{rs}$, $s_0^{rs}.\texttt{mtlsRequests} \equiv \langle \rangle$, $s_0^{rs}.\texttt{oautbEKM} \equiv \langle \rangle$, $s_0^{rs}.\texttt{ids}$ being the sequence of identities for which the resource server manages resources, $s_0^{rs}.\texttt{rNonce}$ and $s_0^{rs}.\texttt{wNonce}$ being the set of nonces representing read and write access to resources, where each set contains an infinite sequence of nonces for each $id \in s_0^{rs}.\texttt{ids}$, $s_0^{rs}.\texttt{authServ}$ being the domain of the authorization server that the resource server supports.

The relation $R^{rs}$ is again based on the generic HTTPS server model (see Section H), for which the algorithm used for processing HTTP requests is defined in Algorithm 26.

Table IV shows a list of placeholders used in the resource server algorithm.

| Placeholder | Usage |
|---|---|
| $v_1$ | new nonce for mTLS |
| $v_2$ | new nonce for OAUTB |

Table IV: List of placeholders used in the resource server algorithm.

**Description and Remarks:** A resource server has two paths for requesting access to resources, depending on the profile used in the authorization process. To simplify the protected resource access, we use two disjunct set of nonces, where the set *rNonce* represents read access to a resource, and the set *wNonce* represents write access.

As before, there are paths for requesting nonces required for mTLS and OAUTB, as access tokens are bound to read-write clients.

For checking the token binding of an access token when using mTLS, the nonce chosen by the resource server is encrypted with a public key sent by the client. In contrast to an authorization server, a resource server is not required to check the validity of certificates (e.g., by checking the certificate chain), as this was already done by the authorization server that created the bound access token (Section 4.2 of [32]).

Furthermore, we require that $\forall id \in s_0^{rs}.\texttt{ids} : \text{governor}(id) \equiv s_0^{rs}.\texttt{authServ}$, i.e., the resource server contains only resources of identities that are governed by the authorization server $s_0^{rs}.\texttt{authServ}$.

**Algorithm 26** 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 [32, Section 4.2]
5:         **let** $s'$.mtlsRequests $:= s'$.mtlsRequests $+^{\langle\rangle} \langle mtlsNonce, clientKey \rangle$
6:         **let** $m' := \mathsf{enc_s}(\langle \mathsf{HTTPResp}, m.\mathsf{nonce}, 200, \langle\rangle, \mathsf{enc_a}(\langle mtlsNonce, keyMapping(m.\mathsf{host})\rangle, clientKey)\rangle, k)$
7:         **stop** $\langle\langle f, a, m'\rangle\rangle, s'$
8:
9:     **else if** $m$.path $\equiv$ /OAUTB-prepare **then**
10:         **let** $tbNonce := \nu_2$
11:         **let** $s'$.oautbEKM $:= s'$.oautbEKM $+^{\langle\rangle} \mathsf{hash}(\langle m.\mathsf{nonce}, tbNonce, keyMapping[m.\mathsf{host}]\rangle)$
12:         **let** $m' := \mathsf{enc_s}(\langle \mathsf{HTTPResp}, m.\mathsf{nonce}, 200, \langle\rangle, [\mathtt{tb\_nonce}:tbNonce]\rangle, k)$
13:         **stop** $\langle\langle f, a, m'\rangle\rangle, s'$
14:
15:     **else if** $m$.path $\equiv$ /resource-r **then**
16:         **let** $id \in s'$.ids **such that** check_read_AT$(id, m.header[\mathtt{Authorization}], s'.\mathtt{authServ}) \equiv \top$
           $\hookrightarrow$   **if possible; otherwise stop**
17:         **let** $resource \leftarrow s'.\mathtt{rNonce}[id]$
18:         **let** $m' := \mathsf{enc_s}(\langle \mathsf{HTTPResp}, m.\mathsf{nonce}, 200, \langle\rangle, [\mathtt{resource}:resource]\rangle, k)$
19:         **stop** $\langle\langle f, a, m'\rangle\rangle, s'$
20:
21:     **else if** $m$.path $\equiv$ /resource-rw **then**
22:         **if** at_iss $\notin m$.body **then**
23:             **stop**
24:         **if** $m$.body[at_iss] $\not\equiv s'$.authServ **then**
25:             **stop**
26:         **if** MTLS_AuthN $\in m$.body **then**
27:             **let** $mtlsInfo$ **such that** $mtlsInfo \in^{\langle\rangle} s'$.mtlsRequests
               $\hookrightarrow$   $\land mtlsInfo.1 \equiv m$.body[MTLS_AuthN] **if possible; otherwise stop**
28:             **let** $s'$.mtlsRequests $:= s'$.mtlsRequests $- mtlsInfo$
29:             **let** $id \in s'$.ids **such that** check_mtls_AT$(id, m.header[\mathtt{Authorization}],$
             $\hookrightarrow$  $mtlsInfo.2, s'.\mathtt{authServ}) \equiv \top$ **if possible; otherwise stop**
30:         **else**
31:             **let** $ekmInfo \leftarrow s'$.oautbEKM
32:             **let** $TB\_Msg\_provided := m$.headers[Sec-Token-Binding][prov]
33:             **let** $TB\_provided\_pub := TB\_Msg\_provided[\mathtt{id}]$
34:             **let** $TB\_provided\_sig := TB\_Msg\_provided[\mathtt{sig}]$
35:             **if** checksig$(TB\_provided\_sig, TB\_provided\_pub) \not\equiv \top$ **then**
36:                **stop**
37:             **if** extractmsg$(TB\_provided\_sig) \not\equiv ekmInfo$ **then**
38:                **stop**
39:             **let** $id \in s'$.ids **such that** check_oautb_AT$(id, m.header[\mathtt{Authorization}],$
             $\hookrightarrow$  $TB\_provided\_pub, s'.\mathtt{authServ}) \equiv \top$ **if possible; otherwise stop**
40:             **let** $s'$.oautbEKM $:= s'$.oautbEKM $- ekmInfo$
41:         **let** $read \leftarrow \{\top, \bot\}$
42:         **if** $read \equiv \top$ **then**
43:             **let** $resource \leftarrow s'.\mathtt{rNonce}[id]$
44:         **else**
45:             **let** $resource \leftarrow s'.\mathtt{wNonce}[id]$
46:         **let** $m' := \mathsf{enc_s}(\langle \mathsf{HTTPResp}, m.\mathsf{nonce}, 200, \langle\rangle, [\mathtt{resource}:resource]\rangle, k)$
47:         **stop** $\langle\langle f, a, m'\rangle\rangle, s'$

*D. OpenID FAPI with Network Attacker*

The formal model of the FAPI is based on web system as defined in Definition 27 of [44]. Unless otherwise specified, we adhere to the terms as defined in [44].

A web system $\mathcal{FAPI} = (\mathcal{W}, \mathcal{S}, \text{script}, E^0)$ is called a *FAPI web system with a network attacker*. The components of the web system are defined in the following.

- $\mathcal{W} = \text{Hon} \cup \text{Net}$ 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 $\text{Hon} := \text{B} \cup \text{C} \cup \text{AS} \cup \text{RS}$. DNS servers are subsumed by the network attacker and are therefore not modeled explicitly.
- $\mathcal{S}$ contains the scripts shown in Table V, with string representations defined by the mapping script.
- $E^0$ contains only the trigger events as specified in Definition 27 of [44].

| $s \in \mathcal{S}$ | script$(s)$ |
|---|---|
| $R^{\text{att}}$ | `att_script` |
| *script_c_index* | `script_c_index` |
| *script_c_get_fragment* | `script_get_fragment` |
| *script_as_form* | `script_as_form` |

Table V: List of scripts in $\mathcal{S}$ and their respective string representations.

In addition to the set of nonces defined in [44], we specify an infinite sequence of nonces $N_r$ representing read access to some resources and an infinite sequence of nonces $N_w$ representing write access to some resources. We call these nonces *resource access nonces*.

In the following, we will define terms used within the analysis. All other terms are used as defined in [44] unless stated otherwise.

*Definition 7 (Read Client).* A client $c$ with client id *clientId* issued from authorization server *as* is called *read client* (w.r.t. *as*) if $s_0^{as}.\texttt{clients}[\textit{clientId}][\texttt{profile}] \equiv \texttt{r}$.

*Definition 8 (Read-Write Client).* A client $c$ with client id *clientId* issued from authorization server *as* is called *read-write client* (w.r.t. *as*) if $s_0^{as}.\texttt{clients}[\textit{clientId}][\texttt{profile}] \equiv \texttt{rw}$.

*Definition 9 (Web Server Client).* A client $c$ with client id *clientId* issued from authorization server *as* is called *web server client* (w.r.t. *as*) if $s_0^{as}.\texttt{clients}[\textit{clientId}][\texttt{is\_app}] \equiv \perp$.

A client is an *app client* if it is not a web server client.

*Definition 10 (Client Type).* A client $c$ with client id *clientId* issued from authorization server *as* is of type $t$ (w.r.t. *as*) if $s_0^{as}.\texttt{clients}[\textit{clientId}][\texttt{client\_type}] \equiv t$, for $t \in \{\texttt{pub}, \texttt{conf\_JWS}, \texttt{conf\_MTLS}, \texttt{conf\_OAUTB}\}$.

*Definition 11 (Confidential Client).* A client $c$ with client id *clientId* issued from authorization server *as* is called *confidential* (w.r.t. *as*) if it is of type $\texttt{conf\_JWS}$, $\texttt{conf\_MTLS}$ or $\texttt{conf\_OAUTB}$.

A client is a *public* client if it is not a confidential client.

**Remarks:** As stated in Section V-B, we assume that for all authorization servers, a client is either a read client or a read-write client. This also holds true for the client type and for being a web server client. We also note that these properties do not change, as an honest authorization server *as* never changes the values of $s_0^{as}.\texttt{clients}$.

*Definition 12 (Token Endpoint of Authorization Server).* A message is sent to the token endpoint of the authorization server *as* if it is sent to the URL $\langle \texttt{URL}, \texttt{S}, \textit{dom}_{as}, /\texttt{token}, \textit{param}, \textit{frag} \rangle$ for $\textit{dom}_{as} \in \textit{dom}(as)$ and arbitrary *param* and *frag*.

*Definition 13 (Bound Access Token).* An access token $t$ issued from authorization server *as* is bound to the client with client id *clientId* in the configuration $(S, E, N)$ of a run of a FAPI web system $\mathcal{FAPI}$ with a network attacker, if either $\langle \texttt{MTLS}, \textit{id}, \textit{clientId}, t, \textit{mtlsKey}, \texttt{rw} \rangle \in S(as).\texttt{accessTokens}$ or $\langle \texttt{OAUTB}, \textit{id}, \textit{clientId}, t, \textit{tbKey}, \texttt{rw} \rangle \in S(as).\texttt{accessTokens}$, for some values of $\textit{id}, \textit{mtlsKey}$ and $\textit{tbKey}$.

The access token is bound via mTLS or via OAUTB, depending on the first entry of the sequence.

*Definition 14 (Access Token associated with Client, Authorization Server and Identity).* Let $c$ be a client with client id *clientId* issued to $c$ by the authorization server *as*, and let $\textit{id} \in \texttt{ID}^{as}$. We say that an *access token $t$ is associated with c, as and id* in state $S$ of the configuration $(S, E, N)$ of a run $\rho$ of a FAPI web system, if there is a sequence $s \in S(as).\texttt{accessTokens}$ such that $s \equiv \langle \textit{id}, \textit{clientId}, t, \texttt{r} \rangle$, $s \equiv \langle \texttt{MTLS}, \textit{id}, \textit{clientId}, t, \textit{key}, \texttt{rw} \rangle$ or $s \equiv \langle \texttt{OAUTB}, \textit{id}, \textit{clientId}, t, \textit{key}', \texttt{rw} \rangle$, for some *key* and *key'*.

*Validity of Access Tokens*

For checking the validity of an access token $t$ in state $S$ of a configuration $(S, E, N)$ of a run, we define the following functions, where $\textrm{dom}_{as} \in \textrm{dom}(as)$:

- $\textsf{check\_read\_AT}(\textit{id}, t, \textrm{dom}_{as}) \equiv \top \Leftrightarrow \exists \textit{clientId}$ s.t. $\langle \textit{id}, \textit{clientId}, t, \texttt{r} \rangle \in S(as).\texttt{accessTokens}$.
- $\textsf{check\_mtls\_AT}(\textit{id}, t, \textit{key}, \textrm{dom}_{as}) \equiv \top \Leftrightarrow \exists \textit{clientId}$ s.t. $\langle \texttt{MTLS}, \textit{id}, \textit{clientId}, t, \textit{key}, \texttt{rw} \rangle \in S(as).\texttt{accessTokens}$.
- $\textsf{check\_oautb\_AT}(\textit{id}, t, \textit{key}, \textrm{dom}_{as}) \equiv \top \Leftrightarrow \exists \textit{clientId}$ s.t. $\langle \texttt{OAUTB}, \textit{id}, \textit{clientId}, t, \textit{key}, \texttt{rw} \rangle \in S(as).\texttt{accessTokens}$.

As the profiles of the FAPI are essentially regular OIDC flows secured by additional mechanisms, they follow the same goals. Therefore, the following security definitions are similar to the definitions given in Appendix H of [44]. Notably, the conditions under which these properties hold true are not the same as in the case of regular OIDC, as detailed in Section III-C.

We show that these properties hold true in Theorem 1.

## AUTHORIZATION

Intuitively, authorization means that an attacker should not be able to get read or write access to a resource of an honest identity.

To capture this property, we extend the definition given in [44], where the authorization property states that no access token associated with an honest identity may leak. As we assume that in the Read-Write flow, access tokens may leak to an attacker, our definition states that honest resource servers may not provide access to resources of honest identities to the attacker.

More precisely, we require that if an honest resource server provides access to a resource belonging to an honest user whose identity is governed by an honest authorization server, then this access is not provided to the attacker. This includes the case that the resource is not directly accessed by the attacker, but also that no honest client provides the attacker access to such a resource.

*Definition 15 (Authorization Property).* We say that the FAPI web system with a network attacker $\mathcal{FAPI}$ *is secure w.r.t. authorization* iff for every run $\rho$ of $\mathcal{FAPI}$, every configuration $(S, E, N)$ in $\rho$, every authorization server $as \in$ AS that is honest in $S$ with $s_0^{as}$.resource_servers being domains of honest resource servers, every identity $id \in \mathsf{ID}^{as}$ with $b = \mathsf{ownerOfID}(id)$ being an honest browser in $S$, every client $c \in \mathsf{C}$ that is honest in $S$ with client id *clientId* issued to $c$ by $as$, every resource server $rs \in$ RS that is honest in $S$ such that $id \in s_0^{rs}$.ids, $s_0^{rs}$.authServ $\in \mathsf{dom}(as)$ and with $dom_{rs} \in s_0^{as}$.resource_servers (with $dom_{rs} \in \mathsf{dom}(rs)$), every access token $t$ associated with $c$, $as$ and $id$ and every resource access nonce $r \in s_0^{rs}$.rNonce$[id] \cup s_0^{rs}$.wNonce$[id]$ it holds true that:

If $r$ is contained in a response to a request $m$ sent to $rs$ with $t \equiv m.header[\mathtt{Authorization}]$, then $r$ is not derivable from the attackers knowledge in $S$ (i.e., $r \notin d_\emptyset(S(\mathtt{attacker}))$).

We require that the preconfigured domains of the resource servers of an authorization server are domains of honest resource servers, as otherwise, an access token of a read client can trivially leak to the attacker.

## AUTHENTICATION

Intuitively, an attacker should not be able to log in at an honest client under the identity of an honest user, where the identity is governed by an honest authorization server. All relevant participants are required to be honest, as otherwise, the attacker can trivially log in at a client, for example, if the attacker controls the authorization server that governs the identity.

*Definition 16 (Service Sessions).* We say that there is a *service session identified by a nonce n for an identity id at some client c* in a configuration $(S, E, N)$ of a run $\rho$ of a FAPI web system iff there exists some session id $x$ and a domain $d \in \mathsf{dom}(\mathsf{governor}(id))$ such that $S(r)$.sessions$[x][\mathtt{loggedInAs}] \equiv \langle d, id \rangle$ and $S(r)$.sessions$[x][\mathtt{serviceSessionId}] \equiv n$.

*Definition 17 (Authentication Property).* We say that the FAPI web system with a 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 $c \in \mathsf{C}$ that is honest in $S$, every identity $id \in \mathsf{ID}$ with $as = \mathsf{governor}(id)$ being an honest AS and with $b = \mathsf{ownerOfID}(id)$ being an honest browser in $S$, every service session identified by some nonce $n$ for $id$ at $c$, $n$ is not derivable from the attackers knowledge in $S$ (i.e., $n \notin d_\emptyset(S(\mathtt{attacker}))$).

## SESSION INTEGRITY FOR AUTHENTICATION AND AUTHORIZATION

There are two session integrity properties that capture that an honest user should not be logged in under the identity of the attacker and should not use resources of the attacker.

We first define notations for the processing steps that represent important events during a flow of a FAPI web system, similar to the definitions given in [44].

*Definition 18 (User is logged in).* For a run $\rho$ of a FAPI web system with a network attacker $\mathcal{FAPI}$ we say that a browser $b$ was authenticated to a client $c$ using an authorization server $as$ and an identity $u$ in a login session identified by a nonce *lsid* in processing step $Q$ in $\rho$ with

$$Q = (S, E, N) \xrightarrow[r \to E_{\mathrm{out}}]{} (S', E', N')$$

(for some $S$, $S'$, $E$, $E'$, $N$, $N'$) and some event $\langle y, y', m \rangle \in E_{\mathrm{out}}$ such that $m$ is an HTTPS response matching an HTTPS request sent by $b$ to $c$ and we have that in the headers of $m$ there is a header of the form $\langle \mathtt{Set\text{-}Cookie}, [\langle \mathtt{\_\_Secure}, \mathtt{serviceSessionId} \rangle : \langle ssid, \top, \top, \top \rangle] \rangle$ for some nonce *ssid* such that

$S(c).\texttt{sessions}[lsid][\texttt{serviceSessionId}] \equiv ssid$ and $S(c).\texttt{sessions}[lsid][\texttt{loggedInAs}] \equiv \langle d, u \rangle$ with $d \in \texttt{dom}(as)$. We then write $\texttt{loggedIn}_\rho^Q(b, c, u, as, lsid)$.

*Definition 19 (User started a login flow).* For a run $\rho$ of a FAPI web system with a network attacker $\mathcal{FAPI}$ 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 *script_client_index* in that document, and in that script, executed the Line 7 of Algorithm 22, 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 $\langle \texttt{Set-Cookie}, [\langle \_\_\texttt{Secure}, \texttt{sessionId} \rangle : \langle lsid, \top, \top, \top \rangle] \rangle$. We then write $\texttt{started}_\rho^Q(b, c, lsid)$.

*Definition 20 (User authenticated at an AS).* For a run $\rho$ of a FAPI web system with a network attacker $\mathcal{FAPI}$ we say that the user of the browser $b$ authenticated to an authorization server $as$ using an identity $u$ for a login session identified by a nonce $lsid$ at the client $c$ if there is a processing step $Q$ in $\rho$ with

$$Q = (S, E, N) \rightarrow (S', E', N')$$

(for some $S$, $S'$, $E$, $E'$, $N$, $N'$) in which the browser $b$ was triggered, selected a document loaded from an origin of $as$, executed the script *script_as_form* in that document, and in that script, (1) in Line 4 of Algorithm 25, selected the identity $u$, and (2) we have that the *scriptstate* of that document, when triggered, contains a nonce $s$ such that $scriptstate[\texttt{state}] \equiv s$ and $S(r).\texttt{sessions}[lsid][\texttt{state}] \equiv s$. We then write $\texttt{authenticated}_\rho^Q(b, c, u, as, lsid)$.

*Definition 21 (Resource Access).* For a run $\rho$ of a FAPI web system with a network attacker $\mathcal{FAPI}$ we say that a browser $b$ accesses a resource of identity $u$ stored at resource server $rs$ through the session of client $c$ identified by the nonce $lsid$ in processing step $Q$ in $\rho$ with

$$Q = (S, E, N) \rightarrow (S', E', N')$$

(for some $S$, $S'$, $E$, $E'$, $N$, $N'$) with (1) $\langle \langle \_\_\texttt{Secure}, \texttt{sessionid} \rangle, \langle lsid, y, z, z' \rangle \rangle \in^{\langle \rangle} S(b).\texttt{cookies}[d]$ for $d \in \texttt{dom}(c)$, $y, z, z' \in \mathcal{T}_\mathcal{N}$, (2) $S(c).\texttt{sessions}[lsid][\texttt{resource}] \equiv r$ with $r \in s_0^{rs}.\texttt{rNonce}[u] \cup s_0^{rs}.\texttt{wNonce}[u]$ and (3) $S(c).\texttt{sessions}[lsid][\texttt{resource\_server}] \in \texttt{dom}(rs)$. We then write $\texttt{accessesResource}_\rho^Q(b, r, u, c, rs, lsid)$.

*Session Integrity Property for Authentication for Web Server Clients with OAUTB*

This security property captures that, if the client is a web server client with OAUTB, then (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 22 (Session Integrity for Authentication for Web Server Clients with OAUTB).* Let $\mathcal{FAPI}$ be an FAPI web system with a network attacker. We say that $\mathcal{FAPI}$ *is secure w.r.t. session integrity for authentication* iff for every run $\rho$ of $\mathcal{FAPI}$, every processing step $Q$ in $\rho$ with

$$Q = (S, E, N) \rightarrow (S', E', N')$$

(for some $S$, $S'$, $E$, $E'$, $N$, $N'$), every browser $b$ that is honest in $S$, every $as \in \mathsf{AS}$, every identity $u$, every web server client $c \in \mathsf{C}$ of type $\texttt{conf\_OAUTB}$ that is honest in $S$, every nonce $lsid$, and $\texttt{loggedIn}_\rho^Q(b, c, u, as, lsid)$ we have that (1) there exists a processing step $Q'$ in $\rho$ (before $Q$) such that $\texttt{started}_\rho^{Q'}(b, c, lsid)$, and (2) if $as$ is honest in $S$, then there exists a processing step $Q''$ in $\rho$ (before $Q$) such that $\texttt{authenticated}_\rho^{Q''}(b, c, u, as, lsid)$.

*Session Integrity Property for Authorization for Web Server Clients with OAUTB*

This security property captures that, if the client is a web server client with OAUTB, then (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 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 and receive any resource (for write access).

*Definition 23 (Session Integrity for Authorization for Web Server Clients with OAUTB).* Let $\mathcal{FAPI}$ be a FAPI web system with a network attackers. We say that $\mathcal{FAPI}$ *is secure w.r.t. session integrity for authorization* iff for every run $\rho$ of $\mathcal{FAPI}$, every processing step $Q$ in $\rho$ with

$$Q = (S, E, N) \rightarrow (S', E', N')$$

(for some $S$, $S'$, $E$, $E'$, $N$, $N'$), every browser $b$ that is honest in $S$, every $as \in \mathsf{AS}$, every identity $u$, every web server client $c \in \mathsf{C}$ of type $\texttt{conf\_OAUTB}$ that is honest in $S$, every $rs \in \mathsf{RS}$ that is honest in $S$, every nonce $r$, every nonce $lsid$, we have that if $\texttt{accessesResource}_\rho^Q(b, r, u, c, rs, lsid)$ and $s_0^{rs}.\texttt{authServ} \in \texttt{dom}(as)$, then (1) there exists a processing step $Q'$ in $\rho$ (before $Q$)

53

such that $\text{started}_\rho^{Q'}(b, c, \mathit{lsid})$, and (2) if $\mathit{as}$ is honest in $S$, then there exists a processing step $Q''$ in $\rho$ (before $Q$) such that $\text{authenticated}_\rho^{Q''}(b, c, u, \mathit{as}, \mathit{lsid})$.

*A. General Properties*

*Lemma 1 (Host of HTTP Request).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S,E,N)$ in $\rho$ and every process $p \in \mathsf{C} \cup \mathsf{AS} \cup \mathsf{RS}$ that is honest in $S$ it holds true that if the generic HTTPS server calls PROCESS_HTTPS_REQUEST$(m_{dec},k,a,f,s)$ in Algorithm 11, then $m_{dec}.\mathtt{host} \in \mathrm{dom}(p)$, for all values of $k$, $a$, $f$ and $s$.

PROOF. PROCESS_HTTPS_REQUEST is called only in Line 9 of Algorithm 11. The input message $m$ is encrypted asymmetrically. Intuitively, such a message is only decrypted if the process knows the private TLS key, where this private key is chosen (non-deterministically) according to the host of the decrypted message.

More formally, when PROCESS_HTTPS_REQUEST is called, the stop in Line 8 is not called. Therefore, it holds true that

$$\exists \, inDomain, k' : \langle inDomain, k' \rangle \in S(p).\mathtt{tlskeys} \wedge m_{dec}.\mathtt{host} \equiv inDomain$$
$$\Rightarrow \exists \, inDomain, k' : \langle inDomain, k' \rangle \in \mathtt{tlskeys}^p \wedge m_{dec}.\mathtt{host} \equiv inDomain$$
$$\overset{\mathrm{Def.}}{\Rightarrow} \exists \, inDomain, k' : \langle inDomain, k' \rangle \in \{\langle d, tlskey(d) \rangle | d \in \mathrm{dom}(p)\} \wedge m_{dec}.\mathtt{host} \equiv inDomain$$

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

The first step holds true due to $S(p).\mathtt{tlskeys} \equiv s_0^p.\mathtt{tlskeys} \equiv \mathtt{tlskeys}^p$, as this sequence is never changed by any honest process $p$. ∎

*Lemma 2 (Honest Read Client sends Token Request only to* $\mathrm{dom}(as)$*).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S,E,N)$ in $\rho$, every authorization server $as$ that is honest in $S$, every identity $id \in \mathsf{ID}^{as}$, every read client $c$ that is honest in $S$ and every $sid$, it holds true that if Algorithm 16 (SEND_TOKEN_REQUEST) is called with $sessionId \equiv sid$ and $S(c).\mathtt{sessions}[sid][\mathtt{identity}] \equiv id$, then the messages in SEND_TOKEN_REQUEST are sent only to $d \in \mathrm{dom}(as)$.

PROOF. In Algorithm 16, the client sends messages either in Line 19, Line 25 or Line 31. The HTTPS_SIMPLE_SEND in Line 48 can only be reached by a read-write client.

In all three cases, the message is sent to $url.\mathtt{domain}$, which is equal to $S(c).\mathtt{oidcConfigCache}[issuer][\mathtt{token\_ep}].\mathtt{domain}$ (Line 8), due to $session[\mathtt{misconfiguredTEp}] \equiv \bot$ in Line 5 (for read clients, this is always set to $\bot$ in Line 12 of Algorithm 15).

Let $\overline{dom}$ and $\overline{dom'}$ be from $\mathrm{dom}(as)$.

With this, it holds true that:

$$url.\mathtt{domain}$$
$$= S(c).\mathtt{oidcConfigCache}[issuer][\mathtt{token\_ep}].\mathtt{domain}$$
(Line 4 of Alg. 16)
$$= S(c).\mathtt{oidcConfigCache}[S(c).\mathtt{issuerCache}[identity]][\mathtt{token\_ep}].\mathtt{domain}$$
(Line 3 of Alg. 16)
$$= S(c).\mathtt{oidcConfigCache}[S(c).\mathtt{issuerCache}[S(c).\mathtt{sessions}[sid][\mathtt{identity}]]][\mathtt{token\_ep}].\mathtt{domain}$$
(per assumption)
$$= S(c).\mathtt{oidcConfigCache}[S(c).\mathtt{issuerCache}[id]][\mathtt{token\_ep}].\mathtt{domain}$$
(issuerCache is never modified)
$$= S(c).\mathtt{oidcConfigCache}[s_0^c.\mathtt{issuerCache}[id]][\mathtt{token\_ep}].\mathtt{domain}$$
(per definition)
$$= S(c).\mathtt{oidcConfigCache}[\overline{dom}][\mathtt{token\_ep}].\mathtt{domain}$$
(oidcConfigCache is never modified)
$$= s_0^c.\mathtt{oidcConfigCache}[\overline{dom}][\mathtt{token\_ep}].\mathtt{domain}$$
(per definition)
$$= \overline{dom'}$$

Therefore, the token request is always sent to a domain of $as$. ∎

*Lemma 3 (Code used in Token Request was received at Redirection Endpoint).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every client $c$ that is honest in $S$ it holds true that if Algorithm 13 (PROCESS_HTTPS_RESPONSE) is called with *reference*[responseTo] $\equiv$ TOKEN, then
*request*.body[code] $\equiv S(c)$.sessions[*reference*[session]][redirectEpRequest][data][code], with *request* being an input parameter of PROCESS_HTTPS_RESPONSE.

PROOF. Let *sid* := *reference*[session] be the session id with which PROCESS_HTTPS_RESPONSE is called.

Due to *reference*[responseTo] $\equiv$ TOKEN, the corresponding request was sent in Algorithm 16 (SEND_TOKEN_REQUEST), as this is the only algorithm that uses this reference when sending a message. The code included in the request is always the input parameter of SEND_TOKEN_REQUEST (due to Line 14).

SEND_TOKEN_REQUEST is called in one of the following lines: (1) in Line 46 of Algorithm 13 (PROCESS_HTTPS_RESPONSE), (2) in Line 53 of Algorithm 13 (PROCESS_HTTPS_RESPONSE) or (3) in Line 25 of Algorithm 15 (PREPARE_TOKEN_REQUEST).

*Case 1: Algorithm 13:* The authorization code used for calling SEND_TOKEN_REQUEST is taken from $S(c)$.sessions[*sid*][code]. $S(c)$.sessions[*sid*][code] is only set in Line 22 of Algorithm 15 (PREPARE_TOKEN_REQUEST), where this value is taken from the input of PREPARE_TOKEN_REQUEST.

*Case 2: Algorithm 15:* In this case, the authorization code used for calling SEND_TOKEN_REQUEST is directly taken from the input of PREPARE_TOKEN_REQUEST.

Algorithm 15 (PREPARE_TOKEN_REQUEST) is either directly called at the redirection endpoint (Line 59 of Algorithm 12), called in Line 19 of Algorithm 19 (CHECK_FIRST_ID_TOKEN) or in Line 18 of Algorithm 21 (CHECK_RESPONSE_JWS).

In both functions (Algorithm 19 and Algorithm 21), *code* is an input argument.

As both functions are called only at the redirection endpoint (Line 61 or Line 63 of Algorithm 12), it follows that the authorization code used as an input argument of PREPARE_TOKEN_REQUEST is in all cases the authorization code originally contained in the request that was received at /redirect_ep. In Line 57 of Algorithm 12, the data contained in the request is stored in the session under the key redirectEpRequest. ∎

*Lemma 4 (Receiver of Token Request and Authorization Request for Read Clients).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every authorization server *as* that is honest in $S$, every read client $c$ that is honest in $S$ with client id *clientId* that has been issued to $c$ by *as*, every *sid* being a session identifier for sessions in $S(c)$.sessions it holds true that if a request JWS *reqJWS* is created in Line 41 of Algorithm 14 (called with input argument *sessionId* $\equiv$ *sid*) with extractmsg(*reqJWS*[aud]) $\in$ dom(*as*), then the token request of the session send in Algorithm 16 (when called with *sessionId* $\equiv$ *sid*) is sent only to a domain in dom(*as*).

PROOF. Let extractmsg(*reqJWS*[aud]) $\equiv$ *authEndpoint*.host $\in$ dom(*as*) (Line 40 of Algorithm 14), which means that $s_0^c$.oidcConfigCache[$s_0^c$.issuerCache[$S(c)$.sessions[*sid*][identity]]][auth_ep].host $\in$ dom(*as*)
(Line 2 up to 6; oidcConfigCache and issuerCache are never changed by the client). Per definition of the auth_ep key of oidcConfigCache, it follows that $s_0^c$.issuerCache[$S(c)$.sessions[*sid*][identity]] $\in$ dom(*as*).

Per definition, issuerCache is a mapping from identities to a domain of their governor. This means that $S(c)$.sessions[*sid*][identity] $\in$ ID$^{as}$.

As all conditions of Lemma 2 are fulfilled, the token request sent in Algorithm 16 (when called with the input parameter *sid*) is sent only to a domain of *as*. ∎

*Lemma 5 (EKM signed by Client does not leak).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every process $p \in$ AS $\cup$ RS that is honest in $S$, every domain $dom_p \in$ dom($p$), every key $k \in$ *keyMapping*[$dom_p$], every client $c$ that is honest in $S$, every domain $d \in$ Doms, every key $TB\_key \in s_0^c$.tokenBindings[$d$], every terms $n_1$, $n_2$ and every process $p'$ with $p \neq p' \neq c$ it holds true that sig(hash($\langle n_1, n_2, k \rangle$), $TB\_key$) $\notin d_\emptyset(S(p'))$.

PROOF. Let hash($\langle n_1, n_2, k \rangle$) be signed by $c$.

The honest client signs only in one of the following places: (1) in Line 41 of Algorithm 14 (START_LOGIN_FLOW), (2) in the branch at Line 32 of Algorithm 16 (SEND_TOKEN_REQUEST) or (3) in Line 28 of Algorithm 18 (USE_ACCESS_TOKEN).

Case 1

In Line 41 of Algorithm 14, the client signs the request JWS, which has a different structure than the EKM value.

Case 2

If the client signs a value in the branch at Line 32 of Algorithm 16, it follows that SEND_TOKEN_REQUEST was called in Line 53 of Algorithm 13
(PROCESS_HTTPS_RESPONSE) (due to *responseValue*[type] $\equiv$ OAUTB, L. 33 of Alg. 16).

The client signs *responseValue*[ekm] (Line 35 of Algorithm 16), which is an input argument. This value is created in Line 51 of Algorithm 13, where the key of the sequence is chosen by the client ($k \equiv$ *keyMapping*[*request*.host]).

From this, it follows that *request*.host $\in$ dom($p$). Due to the reference value OAUTB_AS (Line 49) it follows that *request* was sent in Line 32 of Algorithm 15 (PREPARE_TOKEN_REQUEST), as this is the only place where this reference value is used.

As *request*.host $\in$ dom($p$), it follows that the message in PREPARE_TOKEN_REQUEST is sent to $p$ via HTTPS_SIMPLE_SEND. Therefore, the signed EKM value that is sent in SEND_TOKEN_REQUEST is also send to $p$ via HTTPS_SIMPLE_SEND, as the messages are sent to the same domain in both algorithms (the value of *misconfiguredTEp* set in Line 9 ff. of Algorithm 15 is the same in both algorithms, and therefore, the domains to which the messages are sent are the same).

Case 3

For an EKM value signed in Line 28 of Algorithm 18 (USE_ACCESS_TOKEN), the same reasoning as in Case 2 holds true.

The signed EKM value is sent only in the header of a request. As neither an honest authorization server nor an honest resource server sends any message containing a header value of a received request, it follows that $\mathsf{sig}(\mathsf{hash}(\langle n_1, n_2, k \rangle), TB\_key)$ does not leak to any other process. ∎

## B. Client Authentication

In the next lemmas, we prove that only the legitimate client can authenticate itself to the token endpoint. More precisely, we show that if all checks pass at the token endpoint of an authorization server, then the legitimate confidential client send the corresponding token request.

In case of a confidential read client using a JWS for authentication, the message authentication code of the JWS is created with a key that only the client knows (besides the authorization server). We show that the JWS never leaks to an attacker, which implies that if a message is received at the token endpoint containing the JWS, then the message was sent by the honest client (as otherwise, the JWS must have leaked).

In case of a read-write client, we assume that the token endpoint might be misconfigured. In this case, it seems to be possible that the attacker uses the message received from the client to authenticate at the token endpoint of the honest authorization server.

We show that this is not possible for both OAUTB and mTLS clients. In case of a confidential read-write client using OAUTB, the audience value of the assertion has the same value as the host of the token request, which means that the attacker cannot use a received assertion for any other audience.

Similarly, the message that is decrypted by the client for mTLS (to prove possession of the private key) contains the domain of the process that created the message, and therefore, the decrypted nonce is sent to the same process that encrypted the nonce.

**Lemma 6 (JWS Client Assertion created by Client does not leak to Third Party).** For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every authorization server *as* that is honest in $S$, every domain $d \in$ dom($as$), every client $c$ that is honest in $S$ with client id *clientId* and client secret *clientSecret* that has been issued to $c$ by *as*, every client assertion $t \equiv \mathsf{mac}([\mathsf{iss}{:}clientId, \mathsf{aud}{:}d], clientSecret)$ with $clientSecret \equiv S(as).\mathsf{clients}[clientId][\mathtt{client\_secret}]$ and every process $p$ with $c \neq p \neq as$ it holds true that $t \notin d_\emptyset(S(p))$.

PROOF.

Let $t \equiv \mathsf{mac}([\mathsf{iss}{:}clientId, \mathsf{aud}{:}d], clientSecret)$ be the assertion that is created by the client.

Honest read client: sends assertion only to $d$.

In case of read clients, the JWS client assertion is only created within the branch at Line 20 of Algorithm 16 and send in Line 25.

As the host value of the message in Line 24 has the same value as the audience value of the assertion, it follows that the message is sent to $d$.

Honest read-write client: sends assertion only to $d$.

A read-write client using OAUTB creates a client assertion only in Line 44 of Algorithm 16. As the client sets $jwt[\mathsf{aud}] := d$ in Line 43, it follows that $url.\mathsf{domain} \equiv d$. As above, the assertion is sent to $d$ with HTTPS_SIMPLE_SEND.

Honest authorization server: never sends an assertion.

The authorization server receiving the assertion never sends it out, and as it never creates any assertions, we conclude that the assertion never leaks to $p$ due to the authorization server.

Combining both points, we conclude that $p$ never receives an assertion which is created by an honest client and send to an honest authorization server.

As the client secret is unique for each client id, and also neither send out by an honest client nor by the honest authorization server, we conclude that $p$ is never in possession of a client assertion with a valid message authentication code. ∎

*Lemma 7 (mTLS Nonce created by AS does not leak to Third Party).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every authorization server $as$ that is honest in $S$, every client $c$ that is honest in $S$ with client id *clientId* issued by $as$, every *mtlsNonce* created in Line 164 of Algorithm 24 in consequence of a request $m$ received at the /MTLS-prepare path of the authorization server (L. 162 of Alg. 24) with $m.\mathrm{body}[\texttt{client\_id}] \equiv clientId$ and every process $p$ with $as \neq p \neq c$ it holds true that *mtlsNonce* $\notin d_{\emptyset}(S(p))$.

PROOF. The authorization server sends a *mtlsNonce* created in Line 164 of Algorithm 24 only in Line 168, where it is asymmetrically encrypted with the public key

$$
\begin{aligned}
&clientKey \\
&\equiv S(as).\texttt{clients}[clientId][\texttt{mtls\_key}] && \text{(Line 165)} \\
&\equiv s_0^{as}.\texttt{clients}[clientId][\texttt{mtls\_key}] && \text{(value is never changed)} \\
&\equiv keyMapping[dom_c] && \text{(def.)}
\end{aligned}
$$

with $dom_c \in \mathrm{dom}(c)$. The corresponding private key is $tlskeys[dom_c] \in tlskeys^c$, which is only known to $c$. (The *mtlsNonce* saved in $\texttt{mtlsRequests}$ is not sent in any other place).

This implies that the encrypted nonce can only be decrypted by $c$. Such a message is decrypted either in Line 44 or Line 57 of Algorithm 13. (The only other places where a message is decrypted asymmetrically by $c$ is in the generic HTTPS server (Line 8 of Algorithm 11), but this message is not decrypted there 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).

Case 1: Line 57

In this case, it follows that $reference[\texttt{responseTo}] \equiv \texttt{MTLS\_RS}$ (Line 54). The only place where this value is used as a reference is in Line 14 of Algorithm 17 (PREPARE_USE_ACCESS_TOKEN). The corresponding request is sent to

$$
\begin{aligned}
&message.\texttt{host} && \text{(L. 13 of Alg. 17)} \\
&\equiv rsHost && \text{(L. 13)} \\
&\equiv S'(c).\texttt{sessions}[sid][\texttt{RS}] && \text{(L. 9 and L. 2, for some } sid)
\end{aligned}
$$

for a previous state $S'$ within the run.

This value is only set in Line 9 of Algorithm 14 (START_LOGIN_FLOW), where it is chosen from the set of preconfigured resource servers $s_0^c.\texttt{oidcConfigCache}[s_0^c.\texttt{issuerCache}[id]][\texttt{resource\_server}]$, for some $id \in \mathrm{ID}$. (The values of $\texttt{oidcConfigCache}$ and $\texttt{isserCache}$ are not changed by the client and are therefore the same as in the initial state).

These values are required to be from $\mathrm{dom}(rs)$ by definition, for $rs \in \mathrm{RS}$, which means that the initial request with reference being $\texttt{MTLS\_RS}$ is sent to $rsHost \in \mathrm{dom}(rs)$, $rs \in \mathrm{RS}$.

Therefore, the value of $request.\texttt{host}$ in Line 58 of Algorithm 13 is from $\mathrm{dom}(rs)$, and the check in this line fails, which means that the corresponding *stop* is executed and the decrypted nonce is not sent.

Case 2: Line 44

The reference $\texttt{MTLS\_AS}$ is only used in Line 29 of Algorithm 15 (PREPARE_TOKEN_REQUEST). The corresponding request is sent to $url.\texttt{domain}$ (Line 28) via $\texttt{HTTPS\_SIMPLE\_SEND}$. If $url.\texttt{domain} \notin \mathrm{dom}(as)$ for $as \in \mathrm{AS}$, then the check of the public key fails as in Case 1. Otherwise, the initial request is sent to the honest authorization server, which means that the decrypted nonce is sent to the authorization server in SEND_TOKEN_REQUEST (as in both algorithms, the messages are sent to the same domain).

Summing up, the client sends the nonce encrypted by the authorization server only back to the authorization server. As an honest authorization server never sends out such a nonce received in a token request, we conclude that the nonce never leaks to any other process. ∎

*Lemma 8 (Client Authentication).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every authorization server $as$ that is honest in $S$, every domain $d \in \mathrm{dom}(as)$, every confidential client $c$ that is honest in $S$ with client id *clientId* issued by $as$, it holds true that:

If a response is sent in Line 158 or in Line 160 of Algorithm 24 due to a request $m$ received at the token endpoint of $as$ with $m.\text{body}[\texttt{client\_id}] \equiv clientId$, then $m$ was sent by $c$.

PROOF. We assume that $m$ was sent by $p \neq c$ and that the authorization server sends a response, which implies that all (applicable) checks are passed. We also note that the authorization server never sends messages to itself. We distinguish the following cases:

Case 1: $S(as).\texttt{clients}[clientId][\texttt{client\_type}] \equiv \texttt{conf\_JWS}$

In this case, the client id belongs to a read client. As the response in Line 160 is sent, all relevant checks passed successfully, in particular, the checks in Line 117 and Line 120 of Algorithm 24, which means that that $p$ possesses a term $t \equiv \text{mac}([\texttt{iss}:clientId, \texttt{aud}:d], clientSecret)$ with $clientSecret = S(as).\texttt{clients}[clientId][\texttt{client\_secret}]$, contradicting Lemma 6. (As shown in Lemma 1, the host of the message is a domain of the authorization server).

Case 2: $S(as).\texttt{clients}[clientId][\texttt{client\_type}] \equiv \texttt{conf\_OAUTB}$

As in case 1, this is a contradiction to Lemma 6.

Case 3: $S(as).\texttt{clients}[clientId][\texttt{client\_type}] \equiv \texttt{conf\_MTLS}$

Here, the check in Line 124 passes, which means that $p$ knows $mtlsInfo.1$, which is taken from $S(as).\texttt{mtlsRequests}[clientId]$ (Line 123). This sequence was added to $\texttt{mtlsRequests}$ in Line 166, as this is the only place where a term is added to $\texttt{mtlsRequests}$ (in the initial state, $\texttt{mtlsRequests}$ is empty). This implies that the nonce created in Line 164 due to a request $m'$ to $/\texttt{MTLS-prepare}$ with $m'.\text{body}[\texttt{client\_id}] \equiv clientId$ is known to $p$, which is a contradiction to Lemma 7. ∎

## C. PKCE Challenge

*Lemma 9 (PKCE Challenge of Public Client).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every authorization server $as$ that is honest in $S$, every client $c$ of type $\texttt{pub}$ that is honest in $S$ with client id $clientId$ issued to $c$ by $as$ it holds true that:

If $record \in S(as).\texttt{records}$ with $record[\texttt{client\_id}] \equiv clientId$, then $record[\texttt{pkce\_challenge}]$ was created by $c$.

PROOF. Intuitively, this holds true because the authorization request is authenticated, which means that the request and therefore the PKCE challenge were created by the client.

More formally, we note that $S(as).\texttt{records}$ is empty in the initial state. New records are only added in $/\texttt{auth2}$, or more precisely, in Line 73 of Algorithm 24.

Let $m$ be the corresponding request to $/\texttt{auth2}$. Due to $record[\texttt{client\_id}] \equiv clientId$ and Line 15, it follows that $m.\text{body}[\texttt{client\_id}] \equiv clientId$. As the record is added to the state, all (applicable) checks in $/\texttt{auth2}$ have passed.

This means that $\text{checksig}(m.\text{body}[\texttt{request\_jws}], S(as).\texttt{clients}[clientId][\texttt{jws\_key}]) \equiv \top$ (Line 24), which implies that the request JWS was created by $c$ (as the key is preconfigured and not send to any other process).

Due to $record[\texttt{pkce\_challenge}]$ being set from the request JWS (Line 26 and Line 70), it follows that the PKCE code challenge was created by $c$. ∎

## D. Authorization Response

*Lemma 10 (ID Token contained in Authorization Response for Web Server Client does not leak).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every authorization server $as$ that is honest in $S$, every domain $d \in \text{dom}(as)$, every identity $id \in \text{ID}^{as}$ with $b = \text{ownerOfID}(id)$ being an honest browser in $S$, every web server client $c$ that is honest in $S$ with client id $clientId$ that has been issued to $c$ by $as$, every term $n$, every term $h$, every id token $idt \equiv \text{sig}([\texttt{iss}:d, \texttt{sub}:id, \texttt{aud}:clientId, \texttt{nonce}:n, \texttt{s\_hash}:h], k)$ with $k \equiv S(as).\texttt{jwk}$ and every attacker process $a$ it holds true that:

If a request $m$ is sent to the path $/\texttt{auth2}$ of $as$ with $m.\text{body}[\texttt{client\_id}] \equiv clientId$ and $m.\text{body}[\texttt{password}] \equiv \text{secretOfID}(id)$, then for the corresponding response $r$ it holds true that $idt \equiv r.\text{body}[\texttt{id\_token}]$ does not leak to $a$, i.e. , $idt \notin d_\emptyset(S(a))$.

PROOF. We first highlight that $k \equiv s_0^{as}.\texttt{jwk}$, as this value is never changed by the authorization server. Only $as$ knows this value, as it is preconfigured and never transmitted.

Furthermore, we assume that an id token is sent in the response, otherwise, it trivially cannot leak.

Authorization Response does not leak.

Intuitively, the authorization response does not leak because only the honest browser knows the password of the identity, which means that the response is sent back to the honest browser. As the redirection URIs are preregistered, the browser redirects the response directly to the honest client.

This does not hold true in the case of app clients, as we assume that the response can be sent to a wrong app by the operating system.

The formal proof is analogous to the proof of Lemma 6 in [44]. We note that the changes made to the browser algorithms do not change the results of this proof.

ID Token does not leak.

The client never sends an id token. Therefore, we conclude that an id token created at the authorization endpoint does not leak to an attacker.

∎

*Lemma 11 (Code of Read Web Server Clients does not leak).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every authorization server $as$ that is honest in $S$, every domain $d \in \mathrm{dom}(as)$, every identity $id \in \mathrm{ID}^{as}$ with $b = \mathrm{ownerOfID}(id)$ being an honest browser in $S$, every read web server client $c$ that is honest in $S$ with client id *clientId* that has been issued to $c$ by $as$, every code *code* for which there is a record $rec \in^{\langle\rangle} S(as).\mathtt{records}$ with $rec[\mathtt{code}] \equiv code$, $rec[\mathtt{client\_id}] \equiv clientId$, $rec[\mathtt{subject}] \equiv id$, and every attacker process $a$ it holds true that *code* does not leak to $a$, i.e., $code \notin d_\emptyset(S(a))$.

PROOF. The code contained in the authorization response is sent only to the client (analogous to Lemma 10).

Let $m$ be the request received at /auth2 (Line 8 of Algorithm 24) which led to the creation of *rec*.

As the record was added to the state, it follows that the checks in the branch at Line 22 are fulfilled. From this, it follows that the aud value of the request JWS is a domain in $\mathrm{dom}(as)$ (due to Lemma 1).

As shown in Lemma 4, the corresponding token request send by the client is sent only to *as*.

(This does not hold true for read-write clients, as we assume that the token endpoint might be misconfigured, which means that the code might be sent to the attacker.)

As the authorization server does not send the code received at the token endpoint (and only sends newly generated codes at the authorization endpoint), we conclude that *code* does not leak to $a$. ∎

*Lemma 12 (Response JWS created for Web Server Client does not leak).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every authorization server $as$ that is honest in $S$, every domain $d \in \mathrm{dom}(as)$, every identity $id \in \mathrm{ID}^{as}$ with $b = \mathrm{ownerOfID}(id)$ being an honest browser in $S$, every web server client $c$ that is honest in $S$ with client id *clientId* that has been issued to $c$ by $as$, every term $n$ for which a record $rec \in S(as).\mathtt{records}$ exists with $rec[\mathtt{code}] \equiv n$ and $rec[\mathtt{subject}] \equiv id$, every term $s$, every request JWS $t \equiv \mathrm{sig}([\mathtt{iss}{:}d, \mathtt{aud}{:}clientId, \mathtt{code}{:}n, \mathtt{state}{:}s], k)$ with $k \equiv S(as).\mathtt{jwk}$ and every attacker process $a$ it holds true that $t \notin d_\emptyset(S(a))$.

PROOF. Such a term is only created at the authorization endpoint of *as* (in Line 79 of Algorithm 24), where the corresponding authorization request sent to *as* by $b$ (as only $b$ knows the secret of *id*). Therefore, the request JWS follows the same path as an id token created at the authorization endpoint. Due to this, the proof of this lemma is analogous to the proof of Lemma 10.

∎

## E. ID Token

*Lemma 13 (ID Token for Confidential Client created at the Token Endpoint).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every authorization server $as$ that is honest in $S$, every domain $d \in \mathrm{dom}(as)$, every identity $id \in \mathrm{ID}^{as}$, every confidential client $c$ that is honest in $S$ with client id *clientId* that has been issued to $c$ by $as$, every term $n$, every term $h$, every id token $t \equiv \mathrm{sig}([\mathtt{iss}{:}d, \mathtt{sub}{:}id, \mathtt{aud}{:}clientId, \mathtt{nonce}{:}n, \mathtt{at\_hash}{:}h], k)$ with $k \equiv s_0^{as}.\mathtt{jwk}$ and every process $p$ with $c \neq p \neq as$ it holds true that $t \notin d_\emptyset(S(p))$.

PROOF. Such an id token $t$ is only created in Line 153 of Algorithm 24, which is the token endpoint of the authorization server (the only other place where an id token is created by the authorization server is in /auth2 (Line 81), but the id token created there has different attributes, like s_hash).

The id token created at the token endpoint is sent back to the sender of the request in Line 158 or Line 160. As shown in Lemma 8, the request was sent by $c$ (as the client is confidential; more formally, the conditions of the lemma are fulfilled as the request contained $m.\mathrm{body}[\mathtt{client\_id}] \equiv clientId$ due to Lines 94, 97 and 150).

As an honest client never sends out an id token, it follows that $t$ does not leak to $p$. ∎

*Lemma 14 (ID Token created for Web Server Clients).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every authorization server $as$ that is honest in $S$, every identity $id \in \mathrm{ID}^{as}$ with $b = \mathrm{ownerOfID}(id)$ being an honest browser in $S$, every web server client $c$ that is honest in $S$ with client id *clientId* that has been issued to $c$ by $as$, every id token $t$ with $\mathrm{checksig}(t, pub(s_0^{as}.\mathtt{jwk})) \equiv \top$, $\mathrm{extractmsg}(t)[\mathtt{aud}] \equiv clientId$ and $\mathrm{extractmsg}(t)[\mathtt{sub}] \equiv id$ and every attacker process $a$ it holds true that $t \notin d_\emptyset(S(a))$.

PROOF. The private signing key used by the authorization server is preconfigured and never send to other processes. As the signature of the id token is valid, we conclude that it was created by *as* and that its issuer parameter has the value $dom_{as}$, for $dom_{as} \in \mathsf{dom}(as)$. More precisely, this value is equal to $rec[\texttt{issuer}]$ (L. 81 or L. 148 of Alg. 24, where *rec* is a record from records.) This value is only set in Line 63 and is a domain of *as*, as shown in Lemma 1.

If the id token is created at the authorization endpoint (in Line 81 of Algorithm 24), then $t \notin d_\emptyset(S(a))$, as shown in Lemma 10. (The conditions of the lemma are fulfilled as the audience value is *clientId* and taken from the request that was sent to the authorization endpoint and as the request contained the password of *id*).

Otherwise, the id token was created at the token endpoint (in Line 153). As noted in Section V-B, we assume that a web server client is a confidential client. As shown in Lemma 13, it holds true that $t \notin d_\emptyset(S(a))$. ∎

### F. Access Token

*Lemma 15 (Read Access Token does not leak to Attacker).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every authorization server *as* that is honest in *S*, every identity $id \in \mathsf{ID}^{as}$ with $b = \mathsf{ownerOfID}(id)$ being an honest browser in *S*, every read client *c* that is honest in *S* with client id *clientId* that has been issued to *c* by *as*, every access token *t* with $\langle id, clientId, t, r \rangle \in S(as).\texttt{accessTokens}$, every resource server *rs* that is honest in *S* with $dom_{rs} \in s_0^{as}.\texttt{resource\_servers}$ (with $dom_{rs} \in \mathsf{dom}(rs)$) and every attacker process *a* it holds true that if the request to the resource server is sent to $dom_{rs}$ (in Line 13 of Algorithm 18), then $t \notin d_\emptyset(S(a))$.

PROOF. An honest resource server never sends out an access token. The client sends the access token only in Line 13 of Algorithm 18, where it is sent to an honest resource server (per assumption).

In the following, we show that the access token associated with the read client is never sent to the attacker by the authorization server.

Case 1: Confidential Client (Web Server Client or App)

Let *t* be sent in Line 160 of Algorithm 24 (PROCESS_HTTPS_REQUEST) and let *m* be the corresponding request made to the token endpoint of *as*. It holds true that *m* contains the client id *clientId* (due to Line 94 and *clientId* being included in $\langle id, clientId, t, r \rangle$).

From Lemma 8, it follows that *m* was sent by *c*. This means that the token response containing the access token is sent back directly to *c*.

Case 2: Public Client (App)

Let *m* be the message received at the token endpoint of *as* with $m.\mathsf{body}[\texttt{client\_id}] \equiv clientId$ such that the access token *t* is sent back in Line 160 of Algorithm 24.

This implies that the sender of *m* knows a verifier *pkceCV* such that $\mathsf{hash}(pkceCV) \equiv record[\texttt{pkce\_challenge}]$ (due to the check done in Line 128), with $record[\texttt{client\_id}] \equiv clientId$.

As shown in Lemma 9, the corresponding PKCE challenge ($record[\texttt{pkce\_challenge}]$) was created by the honest client *c* (as the challenge was signed by *c* with preconfigured keys). The challenge was created in Line 29 of Algorithm 14 (START_LOGIN_FLOW) as this is the only place where a read client creates a PKCE challenge. Let *reqJWS* be the corresponding request JWS that was created in Line 41.

The authorization server checks the request JWS at the authorization endpoint for each request (Line 22 to Line 30). This check also includes the *aud* value of the request JWS, which is required to be from a domain in $\mathsf{dom}(as)$ (due to Lemma 1).

This implies that $\mathsf{extractmsg}(reqJWS[\texttt{aud}]) \equiv authEndpoint.\mathsf{host} \in \mathsf{dom}(as)$ (Line 40 of Algorithm 14).

As shown in Lemma 4, the corresponding token request of the session is sent only to a domain in $\mathsf{dom}(as)$ (in Algorithm 16, SEND_TOKEN_REQUEST).

The PKCE verifier is a nonce chosen by the client (in Line 28 of Algorithm 14), and send only in the token request. Therefore, only *c* and *as* know the PKCE verifier. As the authorization server never sends messages to itself, it follows that *m* was sent by *c*. Thus, the response (containing the access token) is sent back to *c*. ∎

*Lemma 16 (Access Token bound via mTLS can only be used by Honest Client).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every authorization server *as* that is honest in *S*, every read-write client *c* of type `conf_MTLS` that is honest in *S* with client id *clientId* issued to *c* by *as*, every access token *t* bound to *c* (via mTLS, as defined in Appendix J), every resource server *rs* that is honest in *S* with $dom_{rs} \in s_0^{as}.\texttt{resource\_servers}$ (with $dom_{rs} \in \mathsf{dom}(rs)$) and every message *m* received at an URL $\langle \mathsf{URL}, \mathsf{S}, dom_{rs}, /\texttt{resource-rw}, param, frag \rangle$ with arbitrary *param* and *frag* and with $\mathsf{MTLS\_AuthN} \in m.\mathsf{body}$ and access token $t \equiv m.\mathsf{header}[\texttt{Authorization}]$ it holds true that:

If a response to *m* is sent in Line 47 of Algorithm 26, then the receiver is *c*.

PROOF. Let *m* be a message with $\mathsf{MTLS\_AuthN} \in m.\mathsf{body}$ containing the access token $t \equiv m.\mathsf{header}[\texttt{Authorization}]$ such that a response is sent in Line 47 of Algorithm 26.

61

This implies that the stop in Line 27 and Line 29 are not executed, and therefore, it holds true that

$$m.\text{body}[\texttt{MTLS\_AuthN}] \equiv S(rs).\texttt{mtlsRequests}.ptr.1$$
$$\wedge \; \texttt{check\_mtls\_AT}(id, t, S(rs).\texttt{mtlsRequests}.ptr.2, s_0^{rs}.\texttt{authServ}) \equiv \top$$

for some $ptr$ and $id$, and as $s_0^{rs}.\texttt{authServ}$ is never changed. This is equivalent to

$$m.\text{body}[\texttt{MTLS\_AuthN}] \equiv S(rs).\texttt{mtlsRequests}.ptr.1 \tag{8}$$
$$\wedge \; \langle \texttt{MTLS}, id, clientId, t, S(rs).\texttt{mtlsRequests}.ptr.2, \texttt{rw} \rangle \in S(\text{dom}^{-1}(s_0^{rs}.\texttt{authServ})).\texttt{accessTokens} \tag{9}$$

for the client id $clientId$ of $c$ (as the access token is bound to $c$) and per definition of $\texttt{check\_mtls\_AT}$ in Appendix J.

An entry as in 9 is only created in Line 142 of Algorithm 24 (PROCESS_HTTPS_REQUEST).

As shown in Lemma 8, the request $m'$ that led to the creation of this sequence was sent from the honest client $c$ with client id $clientId$ (intuitively, the client of type $\texttt{conf\_MTLS}$ has authenticated itself at the token endpoint; more formally, the preconditions of the lemma are fulfilled as there is always a response sent by the authorization server after adding an entry to $accessTokens$ and the corresponding request contained $m'.\text{body}[\texttt{client\_id}] \equiv clientId$).

The value $mtlsInfo.2$ in Line 142 of Algorithm 24 is retrieved from $S'(as).\texttt{mtlsRequests}[clientId]$ (in Line 123, for a state $S'$ prior to $S$ within the run). The entries of $\texttt{mtlsRequests}$ are only created in Line 166 of Algorithm 24, where the entry corresponding to the key $clientKey$ is taken from $s_0^{as}.\texttt{clients}[clientId][\texttt{mtls\_key}]$ (due to Line 165 and the entries of $clients$ not being changed by the authorization server). Per definition, this key has the value $keyMapping[dom_c]$ with $dom_c \in \text{dom}(c)$.

This key is equal to $S(rs).\texttt{mtlsRequests}.ptr.2$ (due to 9). This sequence was added to $S(rs).\texttt{mtlsRequests}$ in Line 5 of Algorithm 26, where the key was taken from the corresponding request (Line 4). The nonce chosen in Line 3 is sent as $\texttt{enc}_{\texttt{a}}(\langle mtlsNonce, keyMapping(m.\texttt{host}) \rangle, clientKey)$ (Line 6), where the asymmetric key is the public key of $c$ (as shown above; the key is the same as the key that the authorization server used for encrypting a nonce for $c$). As shown in Lemma 1, $m.\texttt{host}$ is a domain of the resource server.

Analogous to Lemma 7, this nonce is only known to $c$ (and $rs$) and does not leak to any other process. It follows that the request $m$ was sent by $c$, and therefore, the response is sent back to $c$. ∎

*Lemma 17 (Private OAUTB Key does not leak).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every authorization server $as$ that is honest in $S$, every read-write client $c$ of type pub or $\texttt{conf\_OAUTB}$ that is honest in $S$ with client id $clientId$ issued to $c$ by $as$, every $id \in \text{ID}^{as}$, every access token $t$ and every process $p$ it holds true that:

If $\langle \texttt{OAUTB}, id, clientId, t, \texttt{pub}(TB\_ref\_key), \texttt{rw} \rangle \in S(as).\texttt{accessTokens}$ and $TB\_ref\_key \in d_{\emptyset}(S(p))$, then $p = c$.

PROOF. Let $m$ be the token request that led to the creation of the sequence in $S(as).\texttt{accessTokens}$. It holds true that $m.\text{body}[\texttt{client\_id}] \equiv clientId$, as this value is included in the sequence $\langle \texttt{OAUTB}, id, clientId, t, \texttt{pub}(TB\_ref\_key), \texttt{rw} \rangle$ (due to Line 94 of Algorithm 24).

Case 1: Client of Type $\texttt{conf\_OAUTB}$

    After the sequence is added to $S(as).\texttt{accessTokens}$ (in Line 144 of Algorithm 24), there are no further checks that lead to a stop, which means that a response is sent in Line 158. As shown in Lemma 8, $m$ was sent by $c$.

    The only place where an honest client of type $\texttt{conf\_OAUTB}$ sends a message to the token endpoint of an authorization server is in Line 48 of Algorithm 16. This means that $m$ was created by $c$ in Line 47 of Algorithm 16.

    Let $TB\_ref\_pub\_key$ be the key stored in the sequence shown above (in $S(as).\texttt{accessTokens}$).

    It holds true that

$$
\begin{aligned}
&TB\_ref\_pub\_key \\
&\equiv m.\texttt{headers}[\texttt{Sec-Token-Binding}][\texttt{ref}][\texttt{id}] \quad &&\text{(Line 109, Line 110 of Alg. 24)} \\
&\equiv \texttt{pub}(s_0^c.\texttt{tokenBindings}[t]) \quad &&\text{(Line 37, Line 39, Line 40 of Alg. 16)}
\end{aligned}
$$

    for some term $t$ (which is the domain of a resource server, but not relevant at this point). We also note that the value of $\texttt{tokenBindings}$ of the client state is the same as in the initial state as it is not changed by the client.

Case 2: Client of Type pub

    Let $record$ be the record chosen in Line 96 of Algorithm 24 when $m$ is received at the token endpoint and the sequence is added to $S'(as).accessTokens$ (for a state $S'$ prior to $S$ within the run). Due to Lines 94, 97 and 144, it holds true that $record[\texttt{client\_id}] \equiv clientId$.

    As shown in Lemma 9, the PKCE challenge contained in $record[\texttt{pkce\_challenge}]$ was created by $c$. This happens only in Line 33 of Algorithm 14 (START_LOGIN_FLOW).

This implies that $record[\texttt{pkce\_challenge}]$ has the value $\mathsf{hash}(\mathsf{pub}(s_0^c.\texttt{tokenBindings}[t']))$, for some term $t'$. As the private keys in $s_0^c.\texttt{tokenBindings}$ are preconfigured and never send to any process, they are only known to $c$. As $m.\mathsf{body}[\texttt{client\_id}]$ is the client id of a read-write client of type $\texttt{pub}$, $m$ is required to contain an OAUTB provided message $TB\_Msg\_provided \equiv m.\mathsf{headers}[\texttt{Sec-Token-Binding}][\texttt{prov}]$ such that $\mathsf{checksig}(TB\_Msg\_provided[\texttt{sig}], TB\_Msg\_provided[\texttt{id}]) \equiv \top$ (due to Lines 105, 106 and 107).

Due to Line 134, it holds true that $\mathsf{hash}(TB\_Msg\_provided[\texttt{id}]) \equiv record[\texttt{pkce\_challenge}]$, which is equal to $\mathsf{hash}(\mathsf{pub}(s_0^c.\texttt{tokenBindings}[t']))$. Therefore, $TB\_Msg\_provided[\texttt{id}] \equiv \mathsf{pub}(s_0^c.\texttt{tokenBindings}[t'])$. This means that $TB\_Msg\_provided[\texttt{sig}]$ was signed by $c$ (as only $c$ knows the corresponding private key).

Due to Line 107, it follows that the provided Token Binding message is equal to $ekmInfo$, which is taken from $S'(as).\texttt{oautbEKM}$. These values are only created in Line 172 of Algorithm 24. Therefore, $ekmInfo$ is equal to $\mathsf{hash}(\langle n_1, n_2, keyMapping(dom_{as})\rangle)$, for some values $n_1$, $n_2$.

As shown in Lemma 5, the signed EKM value does not leak, and therefore, $m$ was sent by $c$.

To conclude this case, we note that $m$ was created in Algorithm 16, where the key of the referred Token Binding is exactly the same as in Case 1.

In both cases, the private key $TB\_ref\_key$ is taken from $s_0^c.\texttt{tokenBindings}$, which is preconfigured and never sent to any other process. Consequently, only $c$ knows this value. ∎

*Lemma 18 (Access Token bound via OAUTB can only be used by Honest Client).* For any run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S,E,N)$ in $\rho$, every authorization server $as$ that is honest in $S$, every read-write client $c$ of type $\texttt{conf\_OAUTB}$ or $\texttt{pub}$ that is honest in $S$ with client id $clientId$ issued to $c$ by $as$, every access token $t$ bound to $c$ (via OAUTB, as defined in Appendix J), every resource server $rs$ that is honest in $S$ with $dom_{rs} \in s_0^{as}.\texttt{resource\_servers}$ (with $dom_{rs} \in \mathsf{dom}(rs)$) and every message $m$ received at an URL $\langle \texttt{URL}, \mathsf{S}, dom_{rs}, /\texttt{resource-rw}, param, frag\rangle$ with arbitrary $param$ and $frag$ with $\texttt{MTLS\_AuthN} \notin m.\mathsf{body}$ and access token $t \equiv m.\mathsf{header}[\texttt{Authorization}]$ it holds true that:

If a response to $m$ is sent, then the receiver is $c$.

PROOF. Let $m$ be a message with $\texttt{MTLS\_AuthN} \notin m.\mathsf{body}$ and access token $t \equiv m.\mathsf{header}[\texttt{Authorization}]$ such that a response to this message is sent in Line 47 of Algorithm 26.

This implies that all applicable checks until Line 47 have passed successfully. Therefore, it holds true that

$$\mathsf{checksig}(TB\_prov\_sig, TB\_prov\_pub) \equiv \top \qquad \text{(L. 35)} \qquad (10)$$

$$TB\_prov\_msg \equiv S(rs).\texttt{oautbEKM}.ptr \qquad \text{(L. 37 and 31)} \qquad (11)$$

$$\mathsf{check\_oautb\_AT}(id, t, TB\_prov\_pub, s_0^{rs}.\texttt{authServ}) \equiv \top \qquad \text{(L. 39)} \qquad (12)$$

for some $ptr$, $id$ and with $TB\_prov\_sig := m.\mathsf{headers}[\texttt{Sec-Token-Binding}][\texttt{prov}][\texttt{sig}]$, $TB\_prov\_pub := m.\mathsf{headers}[\texttt{Sec-Token-Binding}][\texttt{prov}][\texttt{id}]$ and $TB\_prov\_msg := \mathsf{extractmsg}(TB\_prov\_sig)$. (In this case, $id$ can be an arbitrary identity. Even if this is the identity of an attacker, the token is still bound to the client). We also note that $s_0^{rs}.\texttt{authServ}$ is never changed by the resource server.

Due to 12, it follows that $\langle \texttt{OAUTB}, id, clientId, t, TB\_prov\_pub, \texttt{rw}\rangle \in S(\mathsf{dom}^{-1}(s_0^{rs}.\texttt{authServ})).\texttt{accessTokens}$ (per definition of $\mathsf{check\_oautb\_AT}$), for the client id $clientId$ of $c$ (as the access token is bound to $c$).

Combining this with 10 and Lemma 17, it follows that the provided Token Binding message was created by $c$ (as the signature is valid and the corresponding private key is only known to $c$).

We conclude the proof by showing that such a Token Binding message is directly sent from the client to the resource server, and therefore, cannot leak to and be used by another process.

Due to 11, $c$ has signed $\mathsf{hash}(\langle n_1, n_2, keyMapping(dom_{rs})\rangle)$ for some values $n_1, n_2$ and with $dom_{rs} \in \mathsf{dom}(rs)$ (the only place where these values are created and added to the state of the resource server is in Line 11. As shown in Lemma 1, the host of messages returned by the HTTPS generic server is a domain of the resource server).

As shown in Lemma 5, the signed EKM value does not leak to another process. Therefore, $m$ was sent by $c$, which means that the response containing the resource access nonce is sent back to $c$. ∎

### G. Authorization

*Lemma 19 (Authorization).* For every run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S,E,N)$ in $\rho$, every authorization server $as \in \mathsf{AS}$ that is honest in $S$ with $s_0^{as}.\texttt{resource\_servers}$ being domains of honest resource servers, every identity $id \in \mathsf{ID}^{as}$ with $b = \mathsf{ownerOfID}(id)$ being an honest browser in $S$, every client $c \in \mathsf{C}$ that is honest in $S$ with client id $clientId$ issued to $c$ by $as$, every resource server $rs \in \mathsf{RS}$ that is honest in $S$ such that $id \in s_0^{rs}.\texttt{ids}$,

$s_0^{rs}$.authServ $\in$ dom($as$) and with $dom_{rs} \in s_0^{as}$.resource_servers (with $dom_{rs} \in$ dom($rs$)), every access token $t$ associated with $c$, $as$ and $id$ and every resource access nonce $r \in s_0^{rs}$.rNonce[$id$] $\cup s_0^{rs}$.wNonce[$id$] it holds true that:

If $r$ is contained in a response to a request $m$ sent to $rs$ with $t \equiv m.header$[Authorization], then $r$ is not derivable from the attackers knowledge in $S$ (i.e., $r \notin d_\emptyset(S(\text{attacker}))$).

PROOF. Let $c$, $as$, $rs$, $t$, $r$ and $m$ be given as in the description of the lemma.

Resource Server never sends Resource Access Nonce $r$ to Attacker.

We assume that the resource server sends $r$ to the attacker. Consequently, the attacker sent the message $m$ with the access token $t$ to either the /resource-r or the /resource-rw path of the resource server.

**Case 1:** resource-r

As the attacker receives $r$, we conclude that the check done in Line 16 of Algorithm 26 (PROCESS_HTTPS_REQUEST) passes successfully and the identity chosen is this line is $id$ (as a resource nonce associated with $id$ is returned).

Therefore, it holds true that check_read_AT($id, t, s_0^{rs}$.authServ) $\equiv \top$ ($S(rs)$.authServ is the same as in the initial state, as it is not modified by the resource server), and it follows that $\langle id, clientId, t, r \rangle \in S(as)$.accessTokens (per definition of check_read_AT and as $t$ is associated with $c$).

This sequence is only added to accessTokens by the authorization server if $c$ is a read client (L. 145 of Alg. 24). This means that the attacker is in possession of a valid access token for the read client $c$, which is a contradiction to Lemma 15.

**Case 2:** resource-rw

If MTLS_AuthN $\in m$.body, then the check done in Line 29 passes successfully, which means that $\langle$MTLS, $id, clientId, t, key, $rw$\rangle \in S(as)$.accessTokens, for some $key$. This means that the access token $t$ is bound to $c$ via mTLS and that the client $c$ is a read-write client (due to L. 140 of Alg. 24). Therefore, Lemma 16 holds true, which is a contradiction to the assumption that the response is sent to the attacker.

Otherwise, MTLS_AuthN $\notin m$.body, and the access token $t$ is bound to $c$ via OAUTB (with the same reasoning as in the case of mTLS shown above). Furthermore, $c$ is a read-write client of type conf_OAUTB or pub. Now, Lemma 18 holds true, which is again a contradiction to the assumption that the response containing $r$ is sent to the attacker.

We highlight that if $c$ is a read-write client, the resource server sends $r$ only to $c$, which follows directly from Lemma 16 and Lemma 18.

Client never Sends Resource Access Nonce $r$ to Attacker.

In the following, we will show that the resource nonce $r$ received in Line 67 of Algorithm 13 (PROCESS_HTTPS_RESPONSE) is not sent to the attacker.

**Case 1:** $c$ is an app client

In this case, the resource nonce is not sent at all in Algorithm 13, due to the check in Line 70. (Intuitively, the resource is used directly by the app).

**Case 2:** $c$ is a web server client

We assume that the resource nonce received in Line 67 of Algorithm 13 is sent to the attacker.

The only place where the resource nonce is sent by the client is in Line 73 of Algorithm 13, as this is the only place where the client uses the resource nonce. The nonce saved in the session in Line 68 is not used by the client at any other place. The resource nonce is sent to $request$[sender], where $request$ is retrieved from $S(c)$.sessions[$sid$][redirectEpRequest], for some $sid$.

The only place where $S(c)$.sessions[$sid$][redirectEpRequest] is set by the client is in Line 57 of Algorithm 12 (at the redirection endpoint). Intuitively, this means that the resource access nonce is sent back to the sender of the request to the redirection endpoint.

Let $m'$ be the corresponding request that was received at /redirect_ep (L. 21 of Alg. 12). As we assume that the resource nonce is sent to the attacker, it follows that $m'$ was sent by the attacker. The values of redirectEpRequest are set to the corresponding values of $m'$.

**Subcase 2.1: read client**

If the client is a read client, then the token endpoint is chosen correctly. More precisely, the access token $t$ that is used by the client is associated with $c$, $as$ and $id$. Per definition, it follows that $\langle id, clientId, t, r \rangle \in S(as)$.accessTokens. As shown in Lemma 15, it holds true that the access token $t$ does not leak to the attacker, and therefore, the client sent the token request to $as$ in order to get the access token.

The code included in the token request can only be provided at the redirection endpoint (as shown in Lemma 3), which means that $m'$ contained the code that will be used by the client. This means that there is a record $rec$ at the authorization server that associates the code with the identity and the client id. More precisely, it holds true that $rec \in^{\langle\rangle} S'(as)$.records with $rec$[code] $\equiv code$, $rec$[client_id] $\equiv clientId$ and $rec$[subject] $\equiv id$ (as we assume that

the resource is owned by *id* and due to L. 96, L. 97 and L. 146 of Alg. 24), for a state $S'$ prior to $S$. As the read client is a web server client, all conditions of Lemma 11 are fulfilled, which means that the attacker cannot know such a code.

**Subcase 2.2: read-write client**
As the resource nonce $r$ was sent to $c$ by $rs$, it follows that $m$ was sent by $c$.

**Subcase 2.2.1: OpenID Hybrid Flow:** As Line 47 of Algorithm 26 is executed, it follows that the check in Line 24 passes successfully, and therefore, $m.\text{body}[\texttt{at\_iss}] \equiv s_0^{rs}.\texttt{authServ} \in \text{dom}(as)$ (the value of $\texttt{authServ}$ is not changed by an honest resource server and stays the same as in the initial state).

The client sends messages to the $\texttt{resource-rw}$ path of a resource server only in Line 31 of Algorithm 18 (USE_ACCESS_TOKEN), hence, $m.\text{body}[\texttt{at\_iss}] \equiv S(c).\texttt{sessions}[sid'][\texttt{idt2\_iss}]$, for some $sid'$ (Line 16 of Algorithm 18).

The value of $\texttt{idt2\_iss}$ is only set in Line 28 of Algorithm 13 (PROCESS_HTTPS_RESPONSE), and therefore, the id token received in the token response has an issuer value from $\text{dom}(as)$.

Due to Line 22 of Algorithm 13, it follows that $issuer \in \text{dom}(as)$, and therefore, $s_0^c.\texttt{jwksCache}[issuer] \equiv \text{pub}(s_0^{as}.\texttt{jwk})$ (per definition of $\texttt{jwksCache}$).

This means that the second id token contained in the token response was signed by $as$ (Line 16), as the private key is only known to $as$. Therefore, this id token was created by $as$.

As the id token has the correct hash value for the access token $t$ (Line 24) and as the id token was created by $as$ (in Line 153 of Algorithm 24), it follows that either $\langle \texttt{MTLS}, id', clientId, t, key_1, \texttt{rw} \rangle \in S(as).\texttt{accessTokens}$ or $\langle \texttt{OAUTB}, id', clientId, t, key_2, \texttt{rw} \rangle \in S(as).\texttt{accessTokens}$, for some values $key_1$, $key_2$.

As the access token $t$ is associated with $c$, $as$ and $id$, it follows that $id' \equiv id$, and therefore, the subject attribute of the second id token is equal to $id$ (due to Line 149).

As shown in Lemma 14, such an id token does not leak to an attacker, therefore, we conclude that the token response was sent by $as$.

Let *id_token_authep* be the id token contained in $m'$ (the message that is received at the redirection endpoint).

Due to Line 18 of Algorithm 13, it follows that $id\_token\_authep[\texttt{sub}] \equiv id$. Furthermore, $id\_token\_authep[\texttt{iss}]$ has the same value as in the second id token (due to Line 20), and therefore, *id_token_authep* was signed by $as$ (with the same reasoning as above).

As the client always checks that the audience value of id tokens are equal to its own client id, again Lemma 14 holds true, which is a contradiction to the assumption that the client sends the resource access nonce to the attacker, as otherwise, an id token with a valid signature from $as$ with $id$ being its subject value and the client id of $c$ being its audience value would have leaked to the attacker.

**Subcase 2.2.2: Code Flow with JARM:** As in the case of the Hybrid Flow, the check done in Line 24 of Algorithm 26 passes successfully and it holds true that $m.\text{body}[\texttt{at\_iss}] \equiv s_0^{rs}.\texttt{authServ} \in \text{dom}(as)$.

The client sets this value only in Line 18 of Algorithm 18, where it is set to $S(c).\texttt{sessions}[sid'][\texttt{JARM\_iss}]$ (with $sid'$ being the session identifier of the corresponding session).

This value is only set in Line 17 of Algorithm 21, where it is set to the issuer of the response JWS *respJWS*. As Algorithm 21 is only called at the redirection endpoint of the client (in Line 63 of Algorithm 12), it follows that this response JWS was sent by the attacker.

Therefore, we conclude that $\texttt{extractmsg}(respJWS)[\texttt{iss}] \in \text{dom}(as)$.

Due to the checks done in Lines 7 and 13 of Algorithm 21, it follows that $jwks \equiv s_0^c.\texttt{jwksCache}[dom\_as] \equiv \text{pub}(s_0^{as}.\texttt{jwk})$ (with $dom\_as \in \text{dom}(as)$), which means that the response JWS is signed by $as$.

The request to the resource server was sent by $c$, which means that the checks done by $c$ prior to sending the request passed successfully, i.e., the hash of the access token received in the token response was contained in the response JWS. More precisely, it holds true

$$\texttt{extractmsg}(respJWS)[\texttt{at\_hash}]$$

(Line 57 of Alg. 12)

$$\equiv \texttt{extractmsg}(S''(c).\texttt{sessions}[sid'][\texttt{redirectEpRequest}][\texttt{data}][\texttt{responseJWS}])[\texttt{at\_hash}]$$

(Line 31 and 32 of Alg. 13)

$$\equiv \texttt{extractmsg}(m''.\text{body}[\texttt{access\_token}])$$

for a state $S'$ prior to $S$ and $m''$ being the token response. The token sent to $rs$ in Line 31 of Algorithm 18 is the input argument of the algorithm. In case of the Read-Write profile, Algorithm 18 is only called in Line 59 or Line 65 of Algorithm 13. In both cases, the token is taken from $S'''(c).\texttt{sessions}[sid'][\texttt{token}]$ which is only set in Line 8 of Algorithm 17. Here, the token is the input argument of the algorithm, which is only called in Line 39

of Algorithm 13. This is the access token received in the token response (i.e., equal to $m''$.body[access_token]). Therefore, we conclude that the response JWS contains the hash of the access token that the client sent to the resource server.

As the resource was sent by $rs$ due to an access token $t$ associated with $c$, $as$ and $id$, it follows per definition that either $\langle \texttt{MTLS}, id, clientId, t, key, \texttt{rw} \rangle \in S(as).\texttt{accessTokens}$ or $\langle \texttt{OAUTB}, id, clientId, t, key', \texttt{rw} \rangle \in S(as).\texttt{accessTokens}$, for some values $key$ and $key'$.

The values of the identity and client identifier are both taken from a record $rec \in \bar{S}(as).\texttt{records}$ with $rec[\texttt{access\_token}] \equiv t$ (Lines 97, 139, 142 and 144 of Algorithm 24).

This record is created at the authorization endpoint of the authorization server, directly before the response JWS is created in Line 79 of Algorithm 24.

For creating access tokens, the authorization server chooses fresh nonces for each authorization request (Line 68). Therefore, the hash of the access token included in the JWS is unique for each authorization response Consequently, the response JWS $respJWS$ (from above), contains the code $rec[\texttt{code}]$ (the code that contained in the same record as the access token $t$). However, this record also contains the identity $id$. This is a contradiction to Lemma 12, as the attacker cannot know a response JWS that contains a code associated with $id$, signed by $as$ and with the client id of $c$.

∎

### H. Authentication

*Lemma 20 (Authentication).* For every run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every configuration $(S, E, N)$ in $\rho$, every client $c \in \mathsf{C}$ that is honest in $S$, every identity $id \in \mathsf{ID}$ with $as = \mathsf{governor}(id)$ being an honest authorization server and with $b = \mathsf{ownerOfID}(id)$ being an honest browser in $S$, every service session identified by some nonce $n$ for $id$ at $c$, $n$ is not derivable from the attackers knowledge in $S$ (i.e., $n \notin d_\emptyset(S(\texttt{attacker}))$).

PROOF. Let *clientId* be the client id that has been issued to $c$ by $as$.

If the client is an app client, then no service session id is sent due to the check done in Line 17 of Algorithm 20 (CHECK_ID_TOKEN). In the following, we look at the case of a web server client.

We assume that the service session id is sent to the attacker by the client. This happens only in Line 25 of Algorithm 20.

This message is sent to $request[\texttt{sender}]$, where $request \equiv S(c).\texttt{sessions}[sessionId][\texttt{redirectEpRequest}]$, for some $sessionId$ (Line 19). This value is only set at the redirection endpoint (in Line 57 of Algorithm 12, PROCESS_HTTPS_REQUEST) and is the sender of the message $m$ received at the redirection endpoint of the client. In other words, the service session id is sent back to the sender of the redirection message, and it follows that $m$ was sent by the attacker.

Per Definition 16 (Service Sessions), it holds true that $S(c).\texttt{sessions}[sessionId][\texttt{loggedInAs}] \equiv \langle d, id \rangle$, where $d \in \mathsf{dom}(\mathsf{governor}(id))$. As the identity is governed by $as$, it follows that $d$ is a domain of $as$, therefore, the value for $issuer$ used in Algorithm 20 is a domain of $as$ (Line 16).

CHECK_ID_TOKEN is only called in Line 40 of Algorithm 13

As shown in Lemma 3, the authorization code used for the corresponding token request was included in the request to the redirection endpoint, i.e., the attacker knows this code.

If the client is a read client, then the token request is sent to the authorization server $as$ (as the token endpoint might only be misconfigured in the read-write flow). More precisely, it holds true that $issuer \in \mathsf{dom}(as)$, as shown above. Therefore, $s_0^c.\texttt{issuerCache}[identity] \in \mathsf{dom}(as)$ (Line 4 of Algorithm 20). Per definition of $\texttt{issuerCache}$, it follows that the identity $identity$ chosen in Line 3 is from $\mathsf{ID}^{as}$. In other words, it holds true that $S(c).\texttt{sessions}[sessionId][\texttt{identity}] \in \mathsf{ID}^{as}$. As shown in Lemma 2, the corresponding token request is sent to $d' \in \mathsf{dom}(as)$.

This means that the attacker knows a code such that there is a record $rec \in^{\langle \rangle} S(as).\texttt{records}$ with $rec[\texttt{code}] \equiv code$, $rec[\texttt{client\_id}] \equiv clientId$ and $rec[\texttt{subject}] \equiv id$. This contradicts Lemma 11.

If the client is a read-write client, then it is possible that the code leaks due to a wrongly configured token endpoint. Here, we distinguish between the following cases:

**Case 1: OpenID Connect Hybrid Flow:** We first look at the case that the client uses the OIDC Hybrid flow (for the current flow in which we assume that the client sent the attacker a service session id). Let *id_token_tep* be the id token the client receives in response to the token request. Let *id_token_auth* be the id token received at the redirection endpoint. This implies that the attacker knows *id_token_auth*. Due to the check done in Line 18 of Algorithm 13 (PROCESS_HTTPS_RESPONSE), it holds true that $\mathsf{extractmsg}(id\_token\_auth)[\texttt{sub}] \equiv id$. When the token request is sent in the read-write flow, the first id token is always checked in Algorithm 19 (CHECK_FIRST_ID_TOKEN). Therefore, the audience value of *id_token_auth* is *clientId* (L. 13 of Alg. 19). The signature of the id token is checked with the same key as in Algorithm 20, which means that

*id_token_auth* is signed with the key $s_0^{as}.\mathtt{jwk}$. This contradicts Lemma 14, as the attacker cannot be in possession of such an id token.

**Case 2: Code Flow with JARM:** As noted above, the value of *issuer* in Line 16 of Algorithm 20 is a domain of *as* and the id token received in the token response was signed with the key $s_0^{as}.\mathtt{jwk}$. As the checks done in Algorithm 20 pass successfully, it follows that the $\mathtt{iss}$ value of the id token is a domain of *as* and the *aud* value is the client id of *c* (Line 9). Furthermore, the *sub* value is *id* (as the SSID for this identity is sent to the attacker). As shown in Lemma 14, such an id token does not leak to the attacker, and therefore, we conclude that the token response was sent by *as*. This means that the token request sent by *c* contains a code *code* such that the authorization server *as* creates the id token with the values depicted above. The values for the audience and subject attributes of the id token are taken from a record $rec \in S'(as).\mathtt{records}$ (for a state $S'$ prior to $S$) with $rec[\mathtt{code}]$ being the code received in the token request (due to Lines 95, 96 149 and 150 of Algorithm 24). As the issuer value chosen in Algorithm 20 is a domain of *as*, it follows that $s_0^c.\mathtt{issuerCache}[session[\mathtt{identity}]] \in \mathrm{dom}(as$ (Lines 3 and 4 of Algorithm 20). Therefore, the issuer value of the request JWS received in the authorization response is also a domain of *as* (Lines 3, 4 of Algorithm 21) and the JWS was created and signed by *as* (with the same reasoning as above). We also note that the value of the issuer stored in the session does not change, as the value of the identity is only set in Line 12 of Algorithm 12. However, this contradict Lemma 12, as the attacker sent the message to the redirection endpoint of the client containing a request JWS created by *as* and with a code that is associated with an honest identity. ∎

## I. Session Integrity

In the following, we show that the Read-Write profile of the FAPI, when used with web server clients and OAUTB, provides session integrity for both authentication and authorization. We highlight that this holds true under the assumption that the state value (which is used for preventing CSRF attacks; Section 10.12 of [25]) leaks to the attacker.

*Lemma 21 (Session Integrity Property for Authentication for Web Server Clients with OAUTB).* For every run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every processing step $Q$ in $\rho$ with

$$Q = (S, E, N) \rightarrow (S', E', N')$$

(for some $S$, $S'$, $E$, $E'$, $N$, $N'$), every browser $b$ that is honest in $S$, every $as \in \mathsf{AS}$, every identity $u$, every web server client $c \in \mathsf{C}$ of type $\mathtt{conf\_OAUTB}$ that is honest in $S$, every nonce *lsid*, and $\mathtt{loggedIn}_\rho^Q(b, c, u, as, lsid)$ we have that (1) there exists a processing step $Q'$ in $\rho$ (before $Q$) such that $\mathtt{started}_\rho^{Q'}(b, c, lsid)$, and (2) if *as* is honest in $S$, then there exists a processing step $Q''$ in $\rho$ (before $Q$) such that $\mathtt{authenticated}_\rho^{Q''}(b, c, u, as, lsid)$.

PROOF.

**Part (1):**

This part of the proof is analogous to the proof given in Lemma 10 of [44]. For completeness, we give the full proof for the FAPI model.

Per definition of $\mathtt{loggedIn}_\rho^Q(b, c, u, as, lsid)$ (Definition 18), it holds true that the client *c* sent the service session id to the browser *b*. This happens only in Line 25 of Algorithm 20 (CHECK_ID_TOKEN), where the service session id is sent to $S(c).\mathtt{sessions}[lsid][\mathtt{redirectEpRequest}][\mathtt{sender}]$ (Lines 19 and 25 of Algorithm 20).

This value is only set in Line 57 of Algorithm 12 (at the redirection endpoint of the client), where it is set to sender of the redirection request. In other words, the browser *b* sent the request to the redirection endpoint.

This request contains the nonce *lsid* as a session id (in a cookie; Line 22 of Algorithm 12).

As this cookie contains the secure prefix, it follows that it was set by the client (i.e., it was not set by the network attacker e.g., over a previous HTTP connection).

This means that the client previously sent a response to *b* in Line 49 of Algorithm 14, as this is the only algorithm in which the client sets a cookie containing a login session id (Line 12). This response is sent to $S''(c).\mathtt{sessions}[lsid][\mathtt{startRequest}][\mathtt{sender}]$ (for a state $S''$ prior to $S$ within the same run) (Line 46), which is only set in Line 12 of Algorithm 12, i.e., the browser *b* sent a POST request to the path $/\mathtt{startLogin}$. This request contains a origin header with an origin of the client (checked in Line 8 of Algorithm 12 and due to Lemma 1).

From the two scripts that could send such a request (*script_c_get_fragment* and *script_client_index*), only *script_client_index* (Algorithm 22) sends such a request. Therefore, it holds true that $\mathtt{started}_\rho^{Q'}(b, c, lsid)$ (for a processing step $Q'$ that happens before $Q$).

**Part (2):**

**Login with ID Token from Token Response:** From the definition of $\mathtt{loggedIn}_\rho^Q(b, c, u, as, lsid)$ (Definition 18), it follows that the client *c* sent a response *m* to *b* containing the header $\langle\mathtt{Set\text{-}Cookie}, [\langle\_\_\mathtt{Secure}, \mathtt{serviceSessionId}\rangle : \langle ssid, \top, \top, \top\rangle]\rangle$ for some nonce *ssid*, and it also holds true that

$S(c).\texttt{sessions}[lsid][\texttt{serviceSessionId}] \equiv ssid$ and $S(c).\texttt{sessions}[lsid][\texttt{loggedInAs}] \equiv \langle d, u \rangle$ (with $d \in \mathrm{dom}(as)$).
Let $m_{\mathrm{redir}}^{\mathrm{req}}$ be the request corresponding to the response $m$.

The cookie contains the secure-prefix, which means that it was set in a connection to the client, i.e., it was set by the client. An honest web server client sends such a response only in Line 25 of Algorithm 20 (CHECK_ID_TOKEN). This algorithm is only called in Line 40 of Algorithm 13 (PROCESS_HTTPS_RESPONSE), and therefore, the id token that is used in Algorithm 20 is received in a response with the reference value TOKEN. Let $m_{\mathrm{token}}^{\mathrm{resp}}$ denote this token response.

**Token Response was sent by as:** Due to $S(c).\texttt{sessions}[lsid][\texttt{loggedInAs}] \equiv \langle d, u \rangle$, it follows that the id token received in the token response was signed by $as$: The value of *issuer* chosen in Line 4 of Algorithm 20 is $d$ (as this value is used in $S(c).\texttt{sessions}[lsid][\texttt{loggedInAs}]$ in Line 16). Therefore, the public key used for checking the signature of the id token is $jwks \overset{\mathrm{L.\ 7,\ Alg.\ 20}}{\equiv} S(c).\texttt{jwksCache}[d] \equiv s_0^c.\texttt{jwksCache}[d] \overset{\mathrm{Def.}}{\equiv} \mathrm{pub}(s_0^{as}.\texttt{jwk})$, and the id token was signed with the corresponding private key (as checked in Line 11; we note that the value of jwksCache is never changed by the client, and therefore, is the same as in the initial state). This private key is only known to $as$ and never sent to any other process, which means that the id token was created by $as$.

As shown in Lemma 13, such an id token does not leak to any process other than $c$ and $as$, which means that the token response was sent by $as$. More precisely, all pre-conditions of the lemma are fulfilled, as the authorization server is honest, the client is a web server client (and therefore, it is confidential), the signature of the id token is valid and it has the client id of $c$ as its audience value (which is checked in Line 9 of Algorithm 20). Furthermore, it contains the attribute at_hash (checked in Line 24 of Algorithm 13, when receiving the token response), which means that it was created at the token endpoint of the authorization server.

**Code included in Token Request was provided by b:** As the client sends the service session id to $b$, it follows that $S(c).\texttt{sessions}[lsid][\texttt{redirectEpRequest}][\texttt{sender}]$ is an IP address of $b$ (Lines 2, 19 and 25 of Algorithm 20). Let $m_{\mathrm{token}}^{\mathrm{req}}$ be the token request corresponding to the token response $m_{\mathrm{token}}^{\mathrm{resp}}$. As shown in Lemma 3, it holds true that $m_{\mathrm{token}}^{\mathrm{req}}.\texttt{body}[\texttt{code}] \equiv S''(c).\texttt{sessions}[lsid][\texttt{redirectEpRequest}][\texttt{data}][\texttt{code}]$ (for a state $S''$ prior to $S$).
This value is only set at the redirection endpoint of the client (Line 57 of Algorithm 12), which means that the code used for the token request was sent by $b$. (We note that as the state value is invalidated at the session of the client (Line 56 of Algorithm 12), for each session, only one request to the redirection endpoint is accepted.)

**Together with the code, b also included a TB-ID:** Along with the code, the browser sent a provided Token Binding message with the ID $S(c).\texttt{sessions}[lsid][\texttt{browserTBID}]$ (Line 50 of Algorithm 12). Here, this value is taken from a provided Token Binding message, which the honest browser only sets in Line 6 of Algorithm 4. The private key used by the browser is only used for the client (i.e., for each domain, the browser uses a different key). (Here, we again highlight that the response to the redirection endpoint was sent by the browser, which is honest).

We conclude that the values for the code and the PKCE verifier (which is the Token Binding ID used by the browser for the client; Line 46 of Algorithm 16) included in the token request were both provided by the browser $b$. More precisely, the token request $m_{\mathrm{token}}^{\mathrm{req}}$ (with the reference value TOKEN) is only sent in Line 48 of Algorithm 16 (SEND_TOKEN_REQUEST), as the client is a web server client of type conf_OAUTB (only read-write clients can be of this type).

**The identity u was authenticated by b:** As noted above, the token response was sent by $as$, which means that all checks done by $as$ passed successfully. Therefore, it holds true that $m_{\mathrm{token}}^{\mathrm{req}}.\texttt{body}[\texttt{pkce\_verifier}] \equiv record[\texttt{pkce\_challenge}]$ (Line 131 of Algorithm 24), with $record \in S'''(as).records$ such that $record[\texttt{code}] \equiv m_{\mathrm{token}}^{\mathrm{req}}.\texttt{body}[\texttt{code}]$ (Lines 95 and 96 of Algorithm 24) (for a state $S'''$ prior to $S$). As noted above, $m_{\mathrm{token}}^{\mathrm{req}}.\texttt{body}[\texttt{pkce\_verifier}] \equiv S(c).\texttt{sessions}[lsid][\texttt{browserTBID}]$ (Line 46 of Algorithm 16), which is a Token Binding ID used by $b$.

The value of $record[\texttt{pkce\_challenge}]$ is only set in Line 72 of Algorithm 24 (as the client is a web server client of type conf_OAUTB; we note that this client id is included in the record), where it is set to the value $TB\_referred\_pub \equiv \bar{m}.\texttt{headers}[\texttt{Sec-Token-Binding}][\texttt{ref}][\texttt{id}]$ (Lines 38 and 41 of Algorithm 24), with $\bar{m}$ being the message which the AS receives at the authorization endpoint.

To sum up the previous paragraphs, it holds true that the message $\bar{m}$ contains a valid Token Binding message (i.e., with a valid signature, as this is always checked by the AS) with the Token Binding ID $\bar{m}.\texttt{headers}[\texttt{Sec-Token-Binding}][\texttt{ref}][\texttt{id}] \equiv record[\texttt{pkce\_challenge}] \equiv m_{\mathrm{token}}^{\mathrm{req}}.\texttt{body}[\texttt{pkce\_verifier}] \equiv S(c).\texttt{sessions}[lsid][\texttt{browserTBID}]$, which is a Token Binding ID of the browser $b$. As only $b$ knows the corresponding private key (and does not reveal this key to another process), we conclude that $\bar{m}$ was sent by $b$.

When sending the token response, $as$ does not only check the PKCE verifier, but also retrieves the identity that is then included in the id token from $record$ (Line 149 of Algorithm 24).

The identity is added to the record in the /auth2 path and taken from $\bar{m}$ (Lines 9 and 62 of Algorithm 24), which

means that the identity $u$ was authenticated by $b$, i.e., $b = \text{ownerOfID}(u)$.

**The redirection request was sent from as (to b):**

*Case 1: OIDC Hybrid Flow:* Let $idt_1$ be the id token contained in $m_{\text{redir}}^{\text{req}}$. (As the flow used for the session with the session identifier *lsid* is the OIDC Hybrid flow, such an id token is always required to be included in the request). When receiving the token response, the client checks if the `sub` and `iss` attributes have the same values in both id tokens, and only continues the flow if the values are the same (Lines 18 and 20 of Algorithm 13).

As we know that the second id token (which is used for logging in the end-user) contains the subject $u$ and an issuer being a domain of *as*, it follows that $idt_1$ contains the same values. Furthermore, the id token is signed by *as* (with the same reasoning as above) and contains the client identifier *clientId* as its audience value.

As shown in Lemma 14, such an id token does not leak to the attacker (we note that $b$ is honest and the client is a web server, i.e., all conditions of the lemma are fulfilled).

Analogous to the proof of Lemma 10 of [44], it follows that the request $m_{\text{redir}}^{\text{req}}$ was caused by a redirection from *as*. Instead of the state value, the id token $idt_1$ is a secret value that does not leak to the attacker. In short, the request $m_{\text{redir}}^{\text{req}}$ was not caused by the attacker, as the id token does not leak. The redirect was also not caused by the client $c$, as $c$ does not send messages containing an id token.

As the request $m_{\text{redir}}^{\text{req}}$ contains a state value (which is checked in Line 52 of Algorithm 12), it follows that the request was not created by the scripts *script_client_index* or *script_as_form*, as these scripts do not send messages containing a state parameter.

The script *script_c_get_fragment* sends only data that is contained in the fragment part of its own URI, and only to itself. This means that the script was sent from the client to the browser, which happens only in Line 36 of Algorithm 12, i.e., at the redirection endpoint.

Altogether, we conclude that there was a location redirect which was sent from *as* to $b$ containing the id token.

*Case 2: Authorization Code Flow with JARM:* Let the flow used in the session with the session identifier *lsid* be a Code Flow with JARM. As shown above, it holds true that $b = \text{ownerOfID}(u)$. The code that the client uses at the token endpoint was sent by the browser, and as an id token with the identity of $u$ is returned by the authorization server, it follows that there is a record that contains the code provided by the browser and the identity $u$ (in $S(as).\text{records}$). Furthermore, the code is contained in a response JWS. The issuer value of the JWS is a domain of *as* (as the id token was signed by *as*, it follows that for this particular session, the client is using a domain of *as* as the expected issuer (in *issuerCache*[*session*][identity]), to which the client also compares the issuer of the response JWS. Therefore, its audience value is *clientId*. Now, all conditions of Lemma 12 are fulfilled, which means that such a response JWS does not leak to the attacker. The remaining argumentation is the same as in the first case. ∎

*Lemma 22 (Session Integrity Property for Authorization for Web Server Clients with OAUTB).* For every run $\rho$ of a FAPI web system $\mathcal{FAPI}$ with a network attacker, every processing step $Q$ in $\rho$ with

$$Q = (S, E, N) \rightarrow (S', E', N')$$

(for some $S$, $S'$, $E$, $E'$, $N$, $N'$), every browser $b$ that is honest in $S$, every $as \in \text{AS}$, every identity $u$, every web server client $c \in \text{C}$ of type `conf_OAUTB` 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 $\text{accessesResource}_\rho^Q(b, r, u, c, rs, lsid)$ and $s_0^{rs}.\text{authServ} \in \text{dom}(as)$, 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.

**Part (1):**

Per definition of accessesResource (Definition 21), it holds true that the browser $b$ has a cookie with the session identifier *lsid* for the origin of the client $c$. As this cookie has the secure prefix set, it follows that the cookie was set by $c$, which happens only in Line 49 of Algorithm 14. The remaining reasoning is the same as in the proof of Lemma 21.

**Part (2) (using the OIDC Hybrid Flow)**

Here, we also first prove the property for the OIDC Hybrid Flow and then show the parts that differ when using the Authorization Code Flow in conjunction with JARM.

**Resource was sent from rs:**

Per Definition of accessesResource, it holds true that $c$ saved the resource access nonce $r$ in $S(c).\text{sessions}[lsid][\text{resource}]$. An honest client stores a resource access nonce only in Line 68 of Algorithm 13. Here, $r$ was contained in response to a request $m_{\text{resource}}^{\text{req}}$ with the reference value `RESOURCE_USAGE` (Line 66 of

Algorithm 13). The client sends requests with this reference value only in Line 31 of Algorithm 18 (this holds true as $c$ is a read-write client).

$m_{\text{resource}}^{\text{req}}$ is sent to the */resource-rw* path of $S(c)$.sessions$[lsid][\text{RS}]$ (Lines 2, 10 and 30 of Algorithm 18).

Per Definition of accessesResource, this is a domain of *rs*, which means that $r$ was sent to the client by *rs*. More precisely, the response of the resource server was sent in Line 47 of Algorithm 26 (within the *resource-rw* path).

**Second ID Token was created by as:**

As the checks done in Lines 22 and 24 of Algorithm 26 passed successfully, it follows that $m_{\text{resource}}^{\text{req}}.\text{body}[\texttt{at\_iss}] \in$ dom$(as)$ (per assumption, it holds true that $s_0^{rs}.\texttt{authServ} \in$ dom$(as)$).

This means that the id token contained in the token response was signed by *as*. More precisely, $S''(c)$.sessions$[lsid][\texttt{idt2\_iss}] \in$ dom$(as)$ (Line 16 of Algorithm 18; here, we are only considering the OIDC Hybrid Flow), for some state $S''$ prior to $S$.

This value is only set in Line 28 of Algorithm 13, where it is set to extractmsg$(m_{\text{token}}^{\text{resp}}.\text{body}[\texttt{id\_token}])[\texttt{iss}]$, where $m_{\text{token}}^{\text{resp}}$ is the token response received in Algorithm 13. Therefore, the $\texttt{iss}$ value of this id token is a domain of *as*. Due to Lines 16 and 22 of Algorithm 13, it follows that (with $dom_{as} \in$ dom$(as)$)

$$\text{checksig}(m_{\text{token}}^{\text{resp}}.\text{body}[\texttt{id\_token}], s_0^{as}.\texttt{jwksCache}[dom_{as}]) \equiv \top$$

$$\overset{\text{Def.}}{\Rightarrow} \text{checksig}(m_{\text{token}}^{\text{resp}}.\text{body}[\texttt{id\_token}], \text{pub}(s_0^{as}.\texttt{jwk})) \equiv \top$$

Therefore, we conclude that the id token was signed by *as*. (The values of $\texttt{jwksCache}$ are never changed by the client, which means that they are the same as in the initial state).

**Token Response was sent by as:**

This id token contains a value for the attribute $\texttt{at\_hash}$ (Line 24 of Algorithm 13). As shown in Lemma 13, such an id token does not leak to the attacker (the id token contains the client id of $c$ due to the check done in Line 26 of Algorithm 13). Thus, we conclude that the token response $m_{\text{token}}^{\text{resp}}$ was sent by *as* (clients do not send any messages containing id tokens).

**Access Token used by c:**

The access token $t$ used by the client in the request to *rs* was contained in the token response. More precisely, the message $m_{\text{resource}}^{\text{req}}$ was sent by the client in Algorithm 18 (as noted above). In Line 8 of this algorithm, the client includes the access token in the header of the message. The access token used here is an input parameter of of the function (USE_ACCESS_TOKEN). A read-write client calls this function either in Line 59 or Line 65 of Algorithm 13. In both cases, the access token is taken from $\overline{S}$.sessions$[lsid][\texttt{token}]$ (Line 55 or 63), for some state $\overline{S}$. This value is only set in Line 8 of Algorithm 17, which is only called in Line 39 of Algorithm 13, where the access token is taken from the body of $m_{\text{token}}^{\text{resp}}$. Therefore, the access token was sent to the client by *as*.

**Sequence in** $\texttt{accessTokens}$ **(state of as):**

We note that the token response was sent by *as* to $c$, which means that the corresponding token request was sent by $c$ and contains the client identifier *clientId* (as $c$ is honest).

Before sending the token response in Line 158 of Algorithm 24, the authorization server adds a sequence to $\overline{S}'(as)$.accessTokens (for some state $\overline{S}'$ prior to $S$).

Let *ATSeq* be the sequence added to the state directly before sending the token response $m_{\text{token}}^{\text{resp}}$. As the client identifier received in the token request belongs to a client using OAUTB, this sequence is equal to $\langle \texttt{OAUTB}, u', clientId, t, k, \texttt{rw} \rangle$, for some identity $u'$ and key $k$. The access token $t$ is the same that is included in the token response. In the following, we will show that the $u' \equiv u$.

We first note that for each access token, there is at most one sequence in $\texttt{accessTokens}$ (in the state of the authorization server) containing this access token. This holds true because the authorization server creates fresh authorization codes and access tokens in the */auth2* path (Lines 67 and 68 of Algorithm 24) for each authorization request received at */auth2*. These values are stored in a record in $\texttt{records}$. When the authorization server receives a request to the token endpoint, it chooses the record depending on the authorization code contained in the token request and invalidates the authorization code contained in the record before creating the sequence for $\texttt{accessTokens}$ (Line 138). Therefore, the access token can only be added once to such a sequence.

As the resource server sent a response in Line 47 of Algorithm 26 to $c$, it follows that all (applicable) checks passed successfully. As the request $m_{\text{resource}}^{\text{req}}$ was sent by $c$, it does not contain the key $\texttt{MTLS\_AuthN}$ in its body (Line 19 of Algorithm 18). Therefore, it holds true that check_oautb_AT$(u, t, k, S(rs).\texttt{authServ}) \equiv \top$ (Line 39 of Algorithm 26). (This holds true as the resource server provides access to a resource of the identity $u$). From the definition of check_oautb_AT, it follows that this identity is contained in the sequence *ATSeq*, and therefore, $u' \equiv u$.

**The identity u was authenticated by b:**

The browser $b$ has a cookie with the session identifier *lsid* for the origin of the client $c$, which means that the request to the redirection endpoint of the client was sent by $b$ (as this value is only known to $b$) With the same reasoning as in Lemma 21, it follows that the identity $u$ was authenticated by $b$. We briefly summarize the argumentation: As the request to the redirection endpoint of $c$ was sent by $b$, it follows that the code used in the token request was provided by $b$. Furthermore, $b$ proved possession of a Token Binding ID. As the token response was sent by $as$, it follows that the process that authenticated the identity $u$ proved possession of the Token Binding ID used by the browser when sending the request to the redirection endpoint. As the private key of this Token Binding ID is only known to $b$, it follows that $b$ authenticated $u$.

**ID Token contained in the Redirection Request :**

As we are looking at the Hybrid Flow, the authorization response is required to contain an id token.

Let $idt_1$ be this id token. When receiving this id token in the redirection request, the client stores it in the session (using the key redirectEpRequest; Line 57 of Algorithm 12).

In Lines 18 and 20 of Algorithm 13 (i.e., after receiving the token response), the client only continues the flow if the subject and issuer values of both id tokens have the same value. As shown above, the issuer of the second id token is a domain of $as$, and the signature of $idt_1$ is checked in Line 15 of Algorithm 19. With the same reasoning as above, it can be seen that this id token is signed by $as$ (as the issuer value is a domain of $as$). Furthermore, we note that this function is called in Line 61 of Algorithm 12 (in the Hybrid flow, this function is always called when receiving the redirection request). The client identifier contained in the id token is *clientId* (checked in Line 13 of Algorithm 19, and as the issuer is a domain of $as$).

As showed above, the token response (containing both the id token and the access token) was sent by $as$. Therefore, the id token contained in the token response contains the identity associated with the access token, which is $u$. This means that $idt_1$ also has the subject $u$. As shown in Lemma 14, such an id token does not leak to the attacker.

**Redirection Request was sent to from as (to b):**

Analogous to Lemma 21, it follows that the redirection request was sent to the browser by $as$, as the id token contained in the redirection request does not leak to the attacker.

## Part (2) (using the Code Flow with JARM)

Here, we focus on the parts that are different from the proof for the OIDC Hybrid flow.

**Resource was sent from rs:** This part is the same as in the Hybrid Flow.

**Response JWS was created by as:**

As above, it holds true that $m_{\text{resource}}^{\text{req}}.\text{body}[\texttt{at\_iss}] \in \text{dom}(as)$ (as this is checked at the resource server).

From this, it follows that $S''(c).\texttt{sessions}[lsid][\texttt{JARM\_iss}] \in \text{dom}(as)$, as the value of at_iss of the request to the resource server is set from this value (Line 18 of Algorithm 18; here, we are considering the Code Flow with JARM). This value is only set in Line 17 of Algorithm 21 (CHECK_RESPONSE_JWS), where it is set to $data[\texttt{iss}] \equiv$ extractmsg($respJWS$)[iss] (Line 8 of Alg. 21), with $respJWS$ being the input argument of Algorithm 21.

Here, the first argument of the function (i.e., the session identifier) is *lsid*, as the issuer of the response JWS is saved in the session identified by *lsid* (again Line 18 of Algorithm 18).

This algorithm is only called in Line 63 of Algorithm 12 (at the redirection endpoint), where the $respJWS$ is set to $data[\texttt{responseJWS}]$. Let $m_{\text{redirect}}^{\text{req}}$ be the request which the client received at the redirection endpoint (i.e., at the path /redirect_ep and for the session identifier *lsid*, i.e., the cookie sessionId contained in the request has the value *lsid*).

As we are looking at the Code Flow using JARM, the value of *data* is equal to $m_{\text{redirect}}^{\text{req}}.\texttt{parameters}$ (Line 32), hence, we conclude that the iss value of the JWS $m_{\text{redirect}}^{\text{req}}.\texttt{parameters}[\texttt{responseJWS}]$ is a domain of $as$.

During the checks that is done in CHECK_RESPONSE_JWS, the client also checks the signature of the JWS (Line 15 of Algorithm 21). With the same reasoning as in the case of the Hybrid Flow, it follows that the response JWS contained in $m_{\text{redirect}}^{\text{req}}$ was signed by $as$.

**Token Response was sent by as:** As the client used the access token $t$ it received in the token response at the resource server, it follows that the check of the hash of the access token done in Line 32 of Algorithm 13 passed successfully. Furthermore, the aud value of $respJWS$ is the client identifier *clientId* (checked in Line 13 of Algorithm 21).

In the following, we assume that the token response $m_{\text{resp}}^{\text{token}}$ was sent by the attacker.

As the hash of the access token was included in a response JWS signed by $as$, and as this JWS contains the client identifier *clientId*, we conclude that $as$ created this access token for $c$, i.e., there is a record *rec* within the state of the AS with $rec[\texttt{access\_tokens}] \equiv t$ and $rec[\texttt{aud}] \equiv clientId$.

As we assume that this access token was sent from the attacker, it follows that it previously leaked to the attacker. In order to leak, the access token must first be sent from the AS.

Let $Q''$ be the processing step in which $as$ sent the access token $t$. As the access token is contained in a record which also contains the client identifier *clientId*, and as this is the client identifier of a web server client (which means that

the client is confidential), it follows that the corresponding code *code* was sent by *c* (due to Lemma 8).

However, this means that the client received the code *code* at the redirection endpoint (for this particular flow in which the access token leaks). As *c* is a read-write client, it only accepts signed authorization responses, i.e, the response is either a response JWS or an id token.

We note that this response (either response JWS or id token) was signed by *as*: The client sent (in this particular flow) the code (i.e., the token request) to *as*. The token request is required to contain the redirection uri at which the redirection response was received (checked at the AS in Line 136 of Algorithm 24 (in the FAPI flows, this is always contained in the token request, as the client is required to previously include this in the authorization request). The check of *as* passed, which means that this URI is a redirection URI used by the client for *as*. The FAPI requires these sets of URIs to be disjunct (i.e., for each AS, the client has a different set of redirection URIs). Therefore, this URI belongs to the set of URIs the client uses for *as*. This means that (for this particular session), the issuer (i.e. the authorization server) the client expects is a domain of *as*.

From this, we conclude that the response JWS or id token which the client received in this (previous) flow at the redirection endpoint was signed by *as* (as both id tokens and response JWS always contain an issuer, which is the same that the client stored at the corresponding session. The signature is checked with a key of this issuer).

The request to the redirection endpoint contained a response JWS, as the AS *as* included the access token corresponding to the code in a response JWS. (Otherwise, it would mean that the client would have created an additional id token with the attribute *c_hash* being the hash of the code (due to the check done at the client in Line 11 of Algorithm 19. However, the AS creates either a response JWS or an id token).

The response JWS the client received (in the main flow, i.e., the one in which the client receives the resource) was created by *as* and contains the authorization code *code*. However, the AS creates exactly one response JWS with a particular code (as codes are nonces that are freshly chosen).

This implies that in both flows that we are looking at (i.e., the original flow in which the honest *as* sends out the access token, and the flow for which we assume that the token endpoint is controlled by the attacker), the client received the same response JWS, and in particular, the same state value (which is contained in the response JWS).

However, the state is unique to each session, as it is chosen by the client as a fresh nonce (Line 39 of Algorithm 14). This is a contradiction to the assumption that the check done in Line Line 32 of Algorithm 13 was executed successfully, as this would mean that the client previously accepted the check done in Line 54 of Algorithm 12 (where it is checked if the state was already invalidated) for the response JWS received in the second flow.

Therefore, we conclude that the token endpoint is controlled by the honest *as*, i.e., the token response was sent by *as*.

**Access Token used by c:** As above, the access token that the client uses for the request to the resource server was contained in the token response, i.e., the access token was sent by *as*.

**Sequence in `accessTokens` (state of as):** As in the case of the Hybrid flow, the state of the authorization server *as* contains the sequence $\langle \texttt{OAUTB}, u, clientId, t, k, \texttt{rw} \rangle$.

**The identity *u* was authenticated by *b*:** As above, the identity *u* was authenticated by *b* (due to the check of the Token Binding ID used by the browser for the client, which happens at the AS).

**Response JWS contained in the Redirection Request:** As the identity *u* is governed by an honest browser, and due to Lemma 12, it follows that the response JWS does not leak to the attacker.

The rest of the proof is the same as above, as now, the response JWS is a value that does not leak to the attacker (instead of an id token).

∎

### J. Proof of Theorem

Theorem 1 follows immediately from Lemmas 19, 20, 21, and 22.