

# OPPID: Single Sign-On with Oblivious Pairwise Pseudonyms

Maximilian Kroschewski  
 Hasso Plattner Institute,  
 University of Potsdam  
 maximilian.kroschewski@hpi.de

Anja Lehmann  
 Hasso Plattner Institute,  
 University of Potsdam  
 anja.lehmann@hpi.de

Cavit Özbay  
 Hasso Plattner Institute,  
 University of Potsdam  
 cavit.oezbay@hpi.de

## ABSTRACT

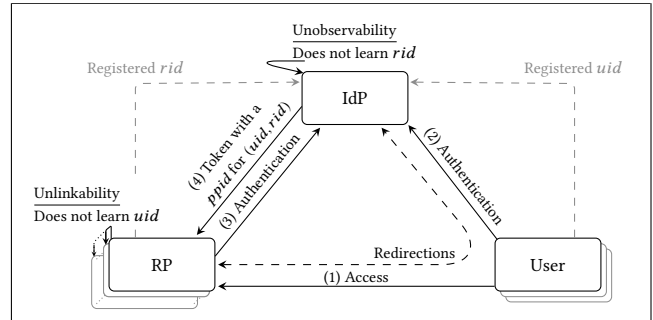
Single Sign-On (SSO) allows users to conveniently authenticate to many Relying Parties (RPs) through a central Identity Provider (IdP). SSO supports unlinkable authentication towards the RPs via *pairwise pseudonyms*, where the IdP assigns the user an RP-specific pseudonym. This feature has been rolled out prominently within Apple’s SSO service. While establishing unlinkable identities provides privacy towards RPs, it actually emphasizes the main privacy problem of SSO: with every authentication request, the IdP learns the RP that the user wants to access. Solutions to overcome this limitation exist, but either assume users to behave honestly or require them to manage long-term cryptographic keys.

In this work, we propose the first SSO system that can provide such pseudonymous authentication in an *unobservable* yet strongly secure and convenient manner. That is, the IdP blindly derives the user’s pairwise pseudonym for the targeted RP without learning the RP’s identity and without requiring key material handled by the user. We formally define the desired security and privacy properties for such unlinkable, unobservable, and strongly secure SSO. In particular, our model includes the often neglected RP authentication: the IdP typically wants to limit its services to registered RPs only and thus must be able to (blindly) verify that it issues the token and pseudonym to such a registered RP. We propose a simple construction that combines signatures with efficient proofs-of-knowledge with a blind, yet verifiable, evaluation of the Hashed-Diffie-Hellman PRF. We prove the security of our construction and demonstrate its efficiency through a prototypical implementation, which requires a running time of 2-20ms per involved party.

## 1 INTRODUCTION

Single Sign-On (SSO) allows users to conveniently authenticate towards multiple online services with the help of a central party, the Identity Provider (IdP). When accessing a service – denoted as a Relying Party (RP) – users are redirected for authentication to the IdP. The IdP then verifies the user and sends a cryptographically signed token attesting the user’s identity *uid* to the RP. The SSO approach frees users from the burden of having to remember a dedicated login credential for each service they want to use, while also providing stronger authentication and simpler deployment for the RPs. Due to these characteristics, SSO has seen widespread adoption in recent years, particularly with major platform providers such as Google, Meta, or Apple serving as IdPs [1, 2, 18].

*Unlinkability via Pseudonyms.* A privacy drawback of SSO systems is that users become linkable across RPs through their identity *uid*, included in each token the IdP signs. Therefore, NIST recommends the use of *Pairwise Pseudonymous Identifiers* [35, §6.2.5] – short *ppid*. The IdP then replaces the user’s identity *uid* in the token with a unique pseudonym *ppid*, which is derived *specifically* for



**Figure 1:** OPPID: Users authenticate to RPs through the IdP. The IdP cannot observe towards which RP the user authenticates, and users are unlinkable via RP-specific pseudonyms *ppids*.

the targeted RP. This protocol feature is supported by the widely-adopted OpenID Connect (OIDC) standard [41, §8], which uses a hash function  $H$  to set  $ppid = H(k, uid, rid)$  where  $k$  is a high-entropy key of the IdP and *uid*, *rid* are the identifiers of the user and RP, respectively. As the IdP assigns deterministic and unique pseudonyms for each user–RP combination, the RP is still ensured that the correct user logs in, and the same user cannot authenticate under multiple pseudonyms, which is known as sybil-resistance. At the same time, the user can engage with different RPs under unlinkable pseudonyms, which has been prominently advertised by Apple in their *Sign in with Privacy* service [5] that uses this feature.

*Main Challenge: Unobservability.* While unlinkable pseudonyms improve user privacy towards RPs, they emphasize another and more fundamental privacy problem of SSO: to derive the user’s RP-specific pseudonym, the IdP must know the RP’s identity *rid* at each user’s login. The pseudonym computation is not the only reason why the *rid* is revealed to the IdP in every authentication request. The most important purpose is to bind the token to the targeted RP for phishing protection, which is done by simply including *rid* in the signed token. Further, the IdP typically wants to restrict its service to *registered* RPs only, which requires some form of authentication from the RP to the IdP too [41].

The lack of unobservability is a significant risk to users’ privacy. The IdP is involved in every online authentication and learns exactly which services and websites users access and when. As SSO is convenient for RPs when only a few IdPs exist, as is currently the case with Google, Meta, and Apple dominating the end-user SSO market, this concentration of information is particularly dangerous.

Thus, an important question is how the convenience and security of SSO can be provided in an *unobservable* manner. This requires that the IdP does not learn the targeted *rid* with every request but can still bind the signed token to the properly authenticated RP and support unlinkability through RP-specific pseudonyms. To maintain convenience for end-users, this should be in the plain SSO

setting, i.e., not relying on any long-term keys or cryptographic credentials managed by the user.

*Partial Solutions Towards Unobservable SSO.* This core challenge of providing SSO in an unobservable manner by the IdP has been addressed in surprisingly few works, and all provide only partial solutions to the problem.

The first work to provide unlinkability and unobservability for users in OIDC was done by Hammann, Sasse, and Basin [26]. Their protocol, denoted as Pairwise POIDC (PPOIDC) [26], lets the IdP *blindly* bind the token to the targeted RP by signing a cryptographic commitment to *rid*. The pseudonym computation (via hash functions) is mostly outsourced to the user and again lets the IdP only *blindly* sign the pseudonym through a commitment. The user must also provide a zero-knowledge proof that the committed pseudonym is indeed derived for *uid*. The protocol does not guarantee that the computed pseudonym was generated for the intended RP, though, allowing corrupt users to generate arbitrary IdP-attested pseudonyms per RP.

The UPPRESSO protocol by Guo et al. [25] also aims at pseudonymous SSO and generates pseudonyms through blind exponentiation of an *rid*-specific group element, enabling the RP to verify that it received a correctly computed pseudonym on its *rid*. The protocol focuses solely on the pseudonym, though, and does not detail how the final token is also strictly bound to the *rid*.

Further, both protocols do not support RP authentication towards the IdP. They only realize a weaker form, where the verification of the RP’s legitimacy is outsourced to the user. Apart from putting more burden on the user, this also implicitly assumes that users must behave honestly. If a user misbehaves, or the user-side verification is not handled properly, the IdP can be tricked into providing its service to malicious and non-registered RPs or sign tokens that assert pseudonymous identities that are incorrect.

The first work to address privacy-preserving RP authentication directly to the IdP was recently done by Kroschewski and Lehmann [31]. Their AIF-ZKP (Authenticated Implicit Flow) protocol ensures that the IdP-issued token is bound to the intended and *authenticated* RP without disclosing *rid* to the IdP. While this approach provides unobservability towards the IdP, the protocol did not provide support for pseudonyms, i.e., it lacks unlinkability.

Thus, there is no protocol – or even security model – for such a fully private yet strongly secure SSO system.

*Concrete Use Case: European Digital Identity Wallet.* Apart from general SSO, there is also a more concrete use case that explicitly demands user authentication with unlinkability, unobservability, and RP authentication: the European Digital Identity Wallet. This Identity Wallet is part of the EU’s eIDAS regulation, which entered into force in May ’24 [17], and aims to establish government-attested and verifiable digital identities with the following requirements:

“Enable privacy-preserving techniques which ensure **unlinkability** [...] [17, §16b] – possibility of users to access services through the use of **pseudonyms** [...] [17, §22] – providers should ensure **unobservability** by not collecting data and not having insight into the transactions of the users [...] [17, §32] – relying parties should provide the information necessary to allow for their identification and **authentication** [...] [17, §17]”

Every EU member state is now tasked with developing such an Identity Wallet for all its citizens and residents, creating an urgent demand for suitable technical solutions.

## 1.1 Our Contributions

In this work, we introduce the first SSO system (OPPID) that combines all properties of unlinkable and unobservable user authentication via a central IdP towards an authenticated RP. More specifically, we propose a protocol where the IdP issues its users strictly RP-bound tokens for a properly authenticated RP and containing RP-specific pseudonyms, yet learns nothing about the RP’s identity. Our protocol achieves its security and privacy properties in a very convenient way, as it still works in the plain SSO setting, i.e., not relying on additional user-managed key material.

*Formal Security Model for OPPID.* The first core challenge is to properly define the security and privacy properties of this 3-party protocol, where each party has complementing views as depicted in Figure 1. In fact, neither of the aforementioned works on pseudonymous SSO provided a formal security model. We formalize *Unlinkability* and *Unobservability* as the two privacy properties, and security is expressed through notions of *Session Binding* and *Request Authentication*. The latter three build upon the model of [31]. The new property of Unlinkability – demanding that two corrupt RPs receiving pseudonymous user authentication cannot decide whether they interact with the same user or not – must carefully exclude trivial wins exploiting the deterministic nature of pseudonyms and their blind computation. Security expressed through Session Binding must hold despite unobservability, in particular guaranteeing that the user can only authenticate under correct and unique pseudonyms  $ppid = F(uid, rid)$  towards an authenticated RP – but where the IdP must not learn anything about *rid*. Our *Session Binding* definition builds upon [31] and discovers and fixes a weakness in their model: to balance security and unobservability, they guarantee Session Binding for *honest* users only, as this allows knowing the RP they intend to authenticate. However, this excludes the most important corruption setting. Thus, beyond extending their Session Binding notion to pseudonymous and unlinkable authentication, we strengthen their model by capturing security for malicious users.

*Provably Secure Protocol  $\pi_{\text{OPPID}}$ .* We propose a protocol that securely realizes all required properties. Our solution builds upon the SSO protocol with privacy-preserving RP authentication from [31] and shows how oblivious – yet strictly binding – pseudonym computation can be added. In a nutshell, [31] uses anonymous credentials for the RP’s authentication towards the IdP and lets the IdP sign a verified commitment on *rid* in its token. To extend this to blindly computed pseudonyms, we rely on a variant of the HashDH (O)PRF [29] to realize  $F(uid, rid)$ . While it is currently not known how such an oblivious PRF can be evaluated on blinded yet *verified inputs* – which would allow ensuring that pseudonyms are computed for the correct *rid* – we circumvent this missing building block: letting the IdP bind non-verified and verified *rid*-derived values in the signed (blinded) token and carefully checking for their consistency in the final token verification, where the *rid* is no longer blind. Thus, we can carry the guarantees from the verified *rid*-bound values over to the ones the IdP had to sign fully

blind and ensure that valid tokens contain properly authenticated pseudonyms  $ppid = F(uid, rid)$ , even with malicious users and RPs.

*Implementation and Evaluation.* To demonstrate the efficiency of our solution, we implemented our protocol using PS signatures [38] and Pedersen commitments [37] for RP authentication, RSA signatures for IdP tokens, and HashDH-style pseudonym computation in the PS signature source group. Our scheme is significantly faster than PPOIDC [26], requiring only 2-17ms per party. We report on the benchmarks of our open-source implementation and compare it in more detail to the closest related works.

## 1.2 Other Related Work

We have already mentioned the related work that is closest to ours: PPOIDC [26], UPPRESSO [25], and AIF-ZKP [31], all of which operate within the plain SSO model. We consider the plain SSO model as one where users do not manage long-term keys or credentials, crucial for convenience and adoption, but this comes with privacy limitations: *colluding* IdP and RPs can trace users. Therefore, we also briefly discuss solutions for pseudonymous user authentication that provide stronger privacy than our work, but at the cost of reduced usability. A summary of such pseudonymous user authentication solutions and a comparison to our work is given in Table 1.

The protocols [13, 21, 24, 42, 43] provide untraceable pseudonymous authentication but either introduce additional parties or rely on user-managed secret keys, thus deviating from the plain SSO setting: PseudoID [13] introduces an additional token service to blindly sign a token that gets bound to a pseudonym and user secret, allowing users to authenticate directly to an RP. Besides the extra party, this approach makes RP authentication towards the IdP impossible due to the token’s independence from the RP’s identity. EL PASSO [43] lets users obtain a short-lived anonymous credential from the IdP, again bound to a user-held key. The user can then locally derive an RP-specific pseudonym and presentation token from that credential and key for each login. This provides untraceable authentication but again detaches the RP authentication from the IdP and requires users to manage a long-term key. PrivSSO [21] requires users to create and manage a dedicated signature key pair for each RP account, which is then bound to a generic IdP token. Using both enables pseudonymous and unobservable authentication towards an RP but relies on even more keys that need to be securely stored and orchestrated by the user. The approaches [24, 42] do not require user keys but leverage secure enclaves on the user side or on an extra party to compute the users’ pseudonyms. The enclave acts as an intermediary between the RPs and the IdP, eliminating the need for the IdP to learn the  $rid$ , while the correctness of the pseudonym is guaranteed through remote attestation.

As an alternative to SSO-(like) solutions, truly user-centric systems [4, 11, 28, 36, 40] exist. They fully remove the role of an on-line IdP and require users to manage their secret keys or anonymous credentials themselves for authentication. While providing the strongest privacy properties, these systems have seen little adoption so far.

## 2 SSO WITH OBLIVIOUS PPIIDS

Before we present our pseudonymous SSO system OPPID, we introduce its entities and detail the properties of pseudonymous user

authentication. Our system builds upon the standard SSO model, where this privacy mechanism is commonly realized via a Pair-wise Pseudonymous Identifier, as outlined by NIST [35, §6.2.5] and further specified by OIDC [41, §8].

### 2.1 Entities & Main Phases

Our OPPID protocol is built for a classic SSO system that encompasses three core entities: Users, Relying Parties (RPs), and a central Identity Provider (IdP):

**Users:** The user is registered with the IdP under a unique username  $uid$ . We assume the IdP handles all user-related registration and authentication but omit those details from our model. For our purposes, the crucial part is that the user is known as  $uid$  to the IdP but has individual pseudonyms  $ppid$  for each Relying Party.

**RPs:** The RP is the service the user wishes to access. The RP is known as  $rid$  to the user and IdP. The RP relies on the IdP for user authentication and for receiving additional user and session information  $ctx$ . The RP knows the user only under their pseudonym  $ppid$ . To use the IdP’s service, the RP must be registered with the IdP.

**IdP:** The IdP is the central authority that RPs and users rely on for authentication. It issues a token  $\tau$ , which asserts to an RP that it is communicating with the user known as  $ppid$ . Apart from the pseudonym, the token is also bound to a particular session referenced by  $sid$ , additional user/session data  $ctx$ , and the targeted RP  $rid$ . While we do not detail how users authenticate to the IdP, our model explicitly covers that only registered and authenticated RPs can use the IdP’s service.

As one of our primary requirements is proper RP authentication, we roughly divide our system into two phases:

*Phase 1: RP Registration.* Before utilizing the IdP’s authentication service, an RP must first register with the IdP. We assume that an RP is uniquely identified through its  $rid$  and denote with  $\mathcal{M}$  the set of registered RPs.

*Phase 2: Authentication.* When users with a unique username  $uid$  want to authenticate towards a specific RP  $rid$ , they initiate the authentication session towards the targeted  $rid$ . Importantly, the user does not reveal her username to the RP. The RP then provides authentication information  $auth$  and a session identifier  $sid$  and sends both – via the user – to the IdP. When forwarding  $sid$ ,  $auth$  to the IdP, the user now reveals the username to the IdP, and we assume that the IdP has the means to check whether the user  $uid$  is correctly authenticated.

While we do not detail how the user authenticates to the IdP, we require that the IdP checks that the request stems from a previously registered RP, i.e.,  $rid \in \mathcal{M}$ . If so, the IdP generates a token  $\tau$  that pseudonymously authenticates the user  $uid$  as  $ppid = F(uid, rid)$  towards  $rid$ , where  $F$  is a pseudonym function we detail next. The final token  $\tau_{fin}$  must be strictly bound to  $rid$ ,  $ppid$ ,  $sid$  and some context  $ctx$ , which stands for additional session/user information vouched for by the IdP.

Approach \ Property	Privacy		Security		Other
	Unobservability	Unlinkability	Req. Auth.	Session Binding	Plain SSO Model
OIDC With Pseudonyms* [41]	○	●	●	●	●
PseudoID [13]	●	●	○	○	○
PPOIDC* [26]	●	●	○	●	●
UPPRESSO* [25] / BISON [27]	●	●	○	●	●
EL PASSO [43]	●	●	○	○	○
AIF-ZKP* [31]	●	○	●	●	●
MISO [42]	●	●	●	●	○
PrivSSO [21]	●	●	○	○	○
Our Work: OPPID	●	●	●	●	●

\*Detailed security comparison given in Sec. 6

**Table 1: Overview of SSO protocols, supporting RP authentication and/or pseudonymous user authentication.**

## 2.2 Pairwise Pseudonymous Identifier

Our system focuses on providing the Pairwise Pseudonymous Identifier (*ppid*) feature of OIDC [41] that hides the user’s *uid* to an RP. The core properties of the *ppid* are:

**Uniqueness:** For every combination of *rid* and *uid*, there exists a unique mapping to a *ppid*. We model this by assuming the *ppid* to be derived through a *deterministic* function  $F$  as

$$ppid = F(uid, rid).$$

**Collision Freeness:** For every *rid* and for all  $uid \neq uid'$ , it must hold that  $F(uid, rid) \neq F(uid', rid)$ , i.e., different users are assigned different pseudonyms towards the same RP *rid*.

**Unlinkable Pseudonyms:** Seeing two pseudonyms for different  $rid_0 \neq rid_1$  with  $ppid_0 = F(uid, rid_0)$  and  $ppid_1 = F(uid', rid_1)$ , it is infeasible to determine whether  $uid = uid'$  or not.

As  $F$  is deterministic and the set of usernames is typically small, the property of unlinkable pseudonyms requires that  $F$  must not be known to the RPs viewing the user’s pseudonyms. This could be achieved by  $F$  being an internal and secret mapping maintained by the IdP or by relying on a keyed function  $F_k$ , where the key  $k$  is only known to the IdP (with the function itself being public). A simple realization for  $F$  is a pseudorandom function.

The challenge we address with our work is to enable the (partially) blind – yet authenticated – *ppid* computation. Specifically, the IdP knows *uid* but not *rid* while ensuring that it computes valid tokens for  $F(uid, rid)$  for the targeted and properly authenticated RP *rid*.

## 2.3 Syntax of OPPID

We present the syntax of OPPID – our *Oblivious Pairwise Pseudonymous Identifier* SSO variant that enables RP authentication towards the IdP and pseudonymous authentication for users towards registered RPs.

*Definition 2.1 (Syntax of OPPID).* In more detail, an OPPID scheme is defined as a tuple of algorithms  $(Setup, KGen_{IdP}, \langle Join_{RP}, Reg_{IdP} \rangle, AInit_U, AReq_{RP}, ARes_{IdP}, AFin_U, Vf_{RP})$ :

$Setup(1^\lambda) \rightarrow pp$  Given the security parameter  $\lambda \in \mathbb{N}$ , returns the public parameters *pp*, which serve as implicit input for all subsequent algorithms.

$KGen_{IdP}(pp) \rightarrow ((isk, \mathcal{M}), ipk)$  Returns the keys for the IdP, where *isk* represents the secret key,  $\mathcal{M}$  the membership state, and *ipk* the public key.

$\langle Join_{RP}(ipk, rid), Reg_{IdP}(isk, rid, \mathcal{M}) \rangle \rightarrow \{(cred, \mathcal{M}'), \perp\}$  An interactive protocol between the RP and IdP. Successful execution results in the RP acquiring a credential *cred*, and the IdP yielding an updated member state  $\mathcal{M}'$ . In case of failure, it returns  $\perp$ .

$AInit_U(ipk, rid) \rightarrow (orid, crid)$  Executed by the user to initialize a token request via an IdP with *ipk* for RP *rid*. It returns a committing value *crid* and an opening *orid*.

$AReq_{RP}(ipk, rid, cred, crid, orid, sid) \rightarrow auth$  Executed by an RP, taking an *rid*, a credential *cred*, user commitment *crid* and opening *orid*, and a random session ID *sid*. It returns the RP authentication *auth*.

$ARes_{IdP}(isk, auth, crid, uid, ctx, sid) \rightarrow \{\tau, \perp\}$  Executed by the IdP using its secret key *isk*, RP authentication data *auth*, user commitment *crid*, context *ctx*, and session identifier *sid*. If the verification of the request fails, it outputs  $\perp$  and a token  $\tau$  otherwise.

$AFin_U(ipk, rid, crid, orid, ctx, sid, \tau) \rightarrow \{(\tau_{fin}, ppid), \perp\}$  Executed by the user to finalize the token  $\tau$ . It takes an RP’s *rid*, user commitment *crid* and opening *orid*, context *ctx*, and *sid*. It outputs  $\perp$  if the inputs are invalid, and the finalized token  $\tau_{fin}$  and her pseudonym *ppid* otherwise.

$Vf_{RP}(ipk, (rid, ppid, ctx, sid), \tau_{fin}) \rightarrow 0/1$  Returns 1 if  $\tau_{fin}$  is valid under *ipk* for  $(rid, ppid, ctx, sid)$  and otherwise 0.

We denote the user, RP, and session space with the sets  $\mathcal{U}$ ,  $\mathcal{R}$ ,  $\mathcal{S}$ , respectively. See Table 2 for an overview of all parameters and App. A for the correctness definition.

*Setup and Registration.* Before offering its authentication service, the IdP generates its key pair  $(isk, ipk)$  based on the public parameters *pp* and initializes its member state  $\mathcal{M}$ . The public key *ipk* is shared with all entities, and tokens issued by the IdP are validated against this key. RPs can then engage in the registration process  $\langle Join_{RP}, Reg_{IdP} \rangle$  with the IdP to obtain their credential *cred*.

*Authentication Flow.* The user authentication (see Figure 3) to an RP *rid* via the IdP with *ipk* involves the following four steps:

Notation	Description
$pp$	Public parameters $pp$ , known by all parties
$isk, ipk, \mathcal{M}$	IdP's secret key $isk$ , public key $ipk$ , RP member state
$rid$	RP's identity $rid \in \mathcal{R}$ registered at the IdP
$cred$	Issued by the IdP to the RP to enable RP authentication
$uid$	User's identity $uid \in \mathcal{U}$ registered at the IdP
$cred, orid$	User commitment $cred$ and opening $orid$ in a session
$auth, sid$	RP authentication for a session referenced by $sid \in \mathcal{S}$
$ctx$	Context that abstracts the user and session information
$\tau$	Authentication token issued by the IdP
$ppid$	RP-specific user pseudonym finalized by the user
$\tau_{fin}$	Token finalized by the user and verified by the RP

Figure 2: Parameters used in an OPPID system.

- (1) The user executes  $\text{AInit}_{\mathcal{U}}$  with  $rid$  to initiate the authentication process, obtaining  $(orid, cred)$ . The user then stores  $orid$  and transmits both values to the RP.
- (2) The RP runs  $\text{AReq}_{\text{RP}}$  to generate  $auth$  using  $cred, orid$ , and  $cred$  to authenticate as a legitimate RP. To ensure freshness, the RP provides a fresh session identifier  $sid$ . The user then forwards  $auth$  to the IdP.
- (3) When the IdP receives a token request from a user  $uid$  for session  $sid$  and implicit authentication  $auth$  for an RP, it executes the algorithm  $\text{ARes}_{\text{IdP}}$ . This results in either a token  $\tau$  or  $\perp$  if the RP authentication fails. The token is now bound to the implicit  $rid$  and explicit  $uid, sid$ , along with additional session information such as timestamps, simplified through context  $ctx$ . We assume that the IdP has properly authenticated  $uid$ , but do not make that explicit here.
- (4) The user runs  $\text{AFin}_{\mathcal{U}}$  to transform the IdP's token  $\tau$  with the committed  $rid$  to verify that the final token corresponds to the initial  $rid$  and to derive an RP-specific pseudonym  $ppid$ . This algorithm takes the user opening  $orid$  and all received information as input to produce the final token  $\tau_{fin}$  and  $ppid$ .

The resulting token  $\tau_{fin}$  is then verified against  $ipk$  to confirm its validity for the tuple  $(rid, ppid, ctx, sid)$ . This explicit verification binds the session information,  $ppid$ , and the RP's  $rid$  together.

### 3 SECURITY MODEL OF OPPID

We now formally define the privacy and security properties expected from an OPPID system, building upon the model of [31]. While we reuse some of their properties (Request Authentication and Unobservability<sup>1</sup>), we also require the additional Unlinkability property and extend their Session Binding model to cover the pseudonymous authentication we aim for. Interestingly, the original model for Session Binding was rather weak, which we strengthen with our work too. We start with a high-level intuition of the desired properties and then present our formal model in the form of game-based security notions.

Note that some requirements are already specified for the pseudonym function  $F$  (see Sec. 2.2): each user must have a single pseudonym per  $rid$  (*Uniqueness*), and distinct users will obtain different pseudonyms for the same  $rid$  (*Collision Freeness*). These two guarantees essentially boil down to requiring that  $F$  is deterministic and injective, so we omit an explicit formalization for these straightforward properties.

<sup>1</sup>These properties were denoted as RP Accountability / RP Hiding in [31]

In addition to these basic pseudonym properties, we require their computation to be done in a blind way by the IdP (*Unobservability*), as well as the unlinkability of the pseudonymous authentication (*Unlinkability*), which also includes the unlinkability of  $F$ .

**Unlinkability:** The user's identity  $uid$  should remain hidden towards RPs – they should only know users under their RP-specific pseudonym  $ppid$ . This implies that when the same user authenticates to two different RPs as  $ppid_0$  and  $ppid_1$ , the two RPs cannot distinguish whether they are communicating with the same user or two different users.

**Unobservability:** The RP's identity  $rid$  should remain hidden towards an IdP during the authentication session. This property assumes that RPs and users are honest and ensures privacy towards a potentially corrupt IdP.

Despite the blind computation of  $ppids$  and the privacy-preserving authentication of RPs, IdP-issued tokens must still be unforgeable and strictly bound to the blindly verified  $rid$  and the correctly computed pseudonym.

**Session Binding:** It is infeasible to create a valid token  $\tau_{fin}$  for a session identified through  $(rid, ppid, ctx, sid)$  that was not properly authenticated or approved by the honest IdP. This notion ensures that the user (and RP) can only generate tokens for the unique and correct  $ppid = F(uid, rid)$  they have jointly authenticated, where RPs must be properly registered with the IdP.

Kroschewski and Lehmann also define the property of *Request Authentication* [31], which complements their Session Binding notion. While Session Binding expresses the security of the final token, this additional property demands that every valid request (including intermediate protocol values) processed by the IdP must originate from a properly registered RP. As this property is not impacted by the blind pseudonym computation focused on in this work, we only restate this notion in our setting and refer to the detailed explanation in App. A.

### 3.1 Oracles

Our definitions are given in a game-based notion, where an adversary  $\mathcal{A}$  runs an experiment with a challenger responsible for managing all honest entities and their private states. These interactions with honest entities are captured through oracles (see Figure 4, right), which we outline before presenting our security games.

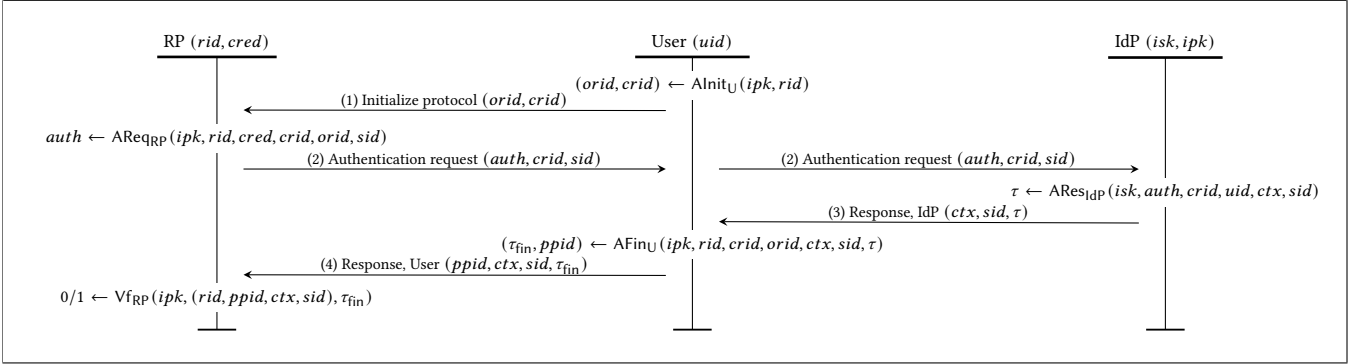
*RP Registration With the IdP.* We give the adversary the ability to register RPs with the IdP, where  $\mathcal{A}$  runs the part of the corrupt party (either RP or IdP) and interacts with the honest counterpart through the oracle.

RegHRP: Runs the registration protocol between an honest RP and the honest IdP. Enables  $\mathcal{A}$  to register an honest RP  $rid$  and later request its authentication data through the  $\text{AReq}_{\text{RP}}$  oracle.

RegCRP: Registers a corrupt RP at the honest IdP.

JoinCIdP: Registers an honest RP with the corrupt IdP. It is only used for our privacy-related Unobservability property, where  $\mathcal{A}$ , the corrupt IdP, aims to break Unobservability.

*Authentication.* We further grant the adversary the capability to intercept, capture, and inject messages between honest parties in an authentication session.



**Figure 3: User authentication in an OPPID scheme to a registered RP, which has previously obtained a credential  $cred$  from the IdP.**

$\text{Alnit}_U$ : Initiates an honest user session with a potentially corrupt  $rid$ , which can later be finalized using the  $\text{AResFin}$  oracle. The main purpose of this oracle is to register the  $rid$  to which an honest user intends to authenticate.

$\text{AReq}_{RP}$ : Returns an honest RP authentication  $auth$  for a potentially adversarially user initiated session, referenced by  $sid$ .

$\text{ARes}_{IdP}$ : Allows  $\mathcal{A}$  to retrieve an honest IdP's response for any token request. This oracle simulates the behavior of an honest IdP and aborts if a token has already been requested for  $sid$ . The inputs, such as the authentication  $auth$ , user commitment  $crid$ , and  $sid$ , could be adversarially generated or partially/fully derived from other oracles.

$\text{AResFin}$ : Enables  $\mathcal{A}$  to obtain a finalized token  $\tau_{fin}$  and  $ppid$  from an honest user session with a (potentially corrupt) RP, initiated using the  $\text{Alnit}_U$  oracle. This oracle simulates secure communication between an honest user and an honest IdP.

$\text{Vf}_{RP}$ : Verifies a session  $(rid, ppid, ctx, sid)$  against a finalized token  $\tau_{fin}$  and the  $ipk$ . It keeps track of tokens presented by the adversary and detects "double-spending", as explained in our Session Binding game.

### 3.2 Unlinkability

Unlinkability captures the core privacy feature of *pseudonymous* authentication, where authentication is done under a pseudonym  $ppid$ 's that hide the user's identity towards malicious RPs. This property requires the unlinkability of pseudonyms produced via  $F$  across RPs, which will be a convenient stepping stone in the formal analysis. Additionally, it ensures that authentication tokens do not leak any information about  $uid$  beyond the requested pseudonym.

We model this property through a classic indistinguishability experiment, which is defined through the game  $\text{Exp}_{\mathcal{A}, \text{OPPID}}^{\text{UNLINK}}$  (see Figure 4). In this game, the IdP is honest, and the adversary  $\mathcal{A}$  controls all RPs. It can register RPs through  $\mathcal{O}.\text{RegCRP}$  and let honest users initiate sessions through  $\mathcal{O}.\text{Alnit}_U$ , and receive tokens and RP-specific  $ppids$  from the honest IdP via  $\mathcal{O}.\text{ARes}_{IdP}$ .

Eventually,  $\mathcal{A}$  outputs two challenge users  $uid_0$  and  $uid_1$  along with common session information  $sid, ctx$  and RP authentication  $auth, crid$ . The game returns the token and  $ppid_b$  for the randomly chosen user  $uid_b$ , requiring the adversary to determine the bit  $b$  better than by guessing.

*Excluding Trivial Wins.* Since  $ppids$  are deterministically (yet blindly) derived for every  $uid, rid$  combination, we must prevent trivial wins exploiting this determinism. Specifically, if the adversary has already learned  $ppid_0$  or  $ppid_1$  through interactions with the oracles, winning this game becomes trivial. Therefore, we ensure that  $\mathcal{A}$  never learns these values through two abort conditions in our game.

Before looking at these conditions, note that Unlinkability is meaningful and defined only for *honest* users. Thus, our challenger computes the pseudonym for an honestly generated  $crid$ , knowing the target RP  $rid$  for which the challenge pseudonym  $ppid_b := F(uid_b, rid)$  is computed.

The first check in our winning condition ensures that  $(uid_d, rid) \notin Q_{ppid}$  for  $d \in \{0, 1\}$ , meaning the adversary never triggered either challenge user  $uid_d$  to initiate an honest session for  $rid$  (via  $\mathcal{O}.\text{Alnit}_U$  and  $\mathcal{O}.\text{ARes}_{IdP}$ ), which would reveal  $ppid_d$ . Here, "honest" implies  $crid_i$  was honestly generated in each request, allowing the challenger to precisely know  $rid$  and the  $ppid_b$  that  $\mathcal{A}$  learned.

When the adversary queries  $\mathcal{O}.\text{ARes}_{IdP}$  with  $crid$  that was not honestly generated, the challenger lacks information about which  $rid$  the adversary requested the token and pseudonym for, requiring stricter abort conditions. This is captured by  $\mathcal{O}.\text{ARes}_{IdP}$  keeping records  $(uid, adv)$  in  $Q_{ppid}$  for such adversarial sessions, and later enforcing that  $(uid_d, adv) \notin Q_{ppid}$  for  $d \in \{0, 1\}$ . Here,  $adv$  denotes that the adversary cannot query  $\mathcal{O}.\text{ARes}_{IdP}$  for *any* adversarial query concerning the challenge users. This is unavoidable as we cannot determine which  $ppid = F(uid, ?)$   $\mathcal{A}$  has obtained.

*Capturing Unlinkability.* Note that  $\mathcal{A}$  can receive pseudonyms  $ppid$  for both challenge users  $uid_0$  and  $uid_1$  for all  $rid' \neq rid$  apart from the target  $rid$  from the challenge query, while using  $\mathcal{O}.\text{Alnit}_U$ . This capability is crucial to capture the desired *unlinkability* of the users' pseudonyms across RPs. Requiring  $\mathcal{A}$  to use  $\mathcal{O}.\text{Alnit}_U$  merely mimics how *honest* users would behave, whom we aim to protect with this property.

**Definition 3.1 (Unlinkability).** An OPPID scheme satisfies Unlinkability if for all PPT adversaries  $\mathcal{A}$ , it holds that

$$\Pr[\text{Exp}_{\mathcal{A}, \text{OPPID}}^{\text{UNLINK}}(\lambda) = 1] \leq 1/2 + \text{negl}(\lambda).$$

*Content of  $ctx$ .* We emphasize that in practice, privacy guaranteed by Unlinkability strongly depends on the information revealed in  $ctx$ . Our model assumes that  $ctx$  is identical for both  $uid_0$  and

<p><i>Unlinkability</i>: <math>\text{Exp}_{\mathcal{A}, \text{OPPID}}^{\text{UNLINK}}(\lambda)</math></p> <hr/> <p><math>pp \leftarrow \text{Setup}(1^\lambda); ((isk, M), ipk) \leftarrow \text{KGen}_{\text{IdP}}(pp); b \leftarrow_{\mathcal{R}} \{0, 1\}</math>  <math>\mathcal{O} := \{\text{RegCRP}, \text{Alnit}_{\text{U}}, \text{ARes}_{\text{IdP}}\}</math>  <math>(uid_0, uid_1, auth, crid, ctx, sid) \leftarrow \mathcal{A}^{\mathcal{O}}(ipk)</math>  Require <math>(rid, crid, orid) \in Q_{\text{rid}}</math>  For <math>d \in \{0, 1\}</math> :  <math>\tau_d \leftarrow \text{ARes}_{\text{IdP}}(isk, auth, crid, uid_d, ctx, sid)</math>  <math>(\tau_{fin_d}, ppid_d) \leftarrow \text{AFin}_{\text{U}}(ipk, rid, crid, orid, ctx, sid, \tau_d)</math>  Require <math>\forall \text{fRP}(ipk, (rid, ppid_d, ctx, sid), \tau_{fin_d}) = 1</math>  <math>b^* \leftarrow \mathcal{A}^{\mathcal{O}}(\tau_{fin_b}, ppid_b)</math>  Abort if for <math>d \in \{0, 1\}</math>: <math>(uid_d, rid) \in Q_{\text{ppid}} \vee (uid_d, adv) \in Q_{\text{ppid}}</math>  Return 1 if <math>b = b^*</math></p> <hr/> <p><i>Unobservability</i>: <math>\text{Exp}_{\mathcal{A}, \text{OPPID}}^{\text{UNOBS}}(\lambda)</math></p> <hr/> <p><math>pp \leftarrow \text{Setup}(1^\lambda); ((isk, M), ipk) \leftarrow \text{KGen}_{\text{IdP}}(pp); b \leftarrow_{\mathcal{R}} \{0, 1\}</math>  <math>\mathcal{O} := \{\text{JoinCIdP}, \text{AReq}_{\text{RP}}\}</math>  <math>(rid_0, rid_1, sid) \leftarrow \mathcal{A}^{\mathcal{O}}((isk, M), ipk)</math>  For <math>d \in \{0, 1\}</math> : Require <math>(rid_d, cred_d) \in \text{HRP}</math>  <math>(orid, crid) \leftarrow \text{Alnit}_{\text{U}}(ipk, rid_b)</math>  <math>auth \leftarrow \text{AReq}_{\text{RP}}(ipk, rid_b, cred_b, crid, orid, sid)</math>  <math>b^* \leftarrow \mathcal{A}^{\mathcal{O}}(auth, crid)</math>  Return 1 if <math>b = b^*</math></p> <hr/> <p><i>Request Authentication</i>: <math>\text{Exp}_{\mathcal{A}, \text{OPPID}}^{\text{REQ-AUTH}}(\lambda)</math></p> <hr/> <p><math>pp \leftarrow \text{Setup}(1^\lambda); ((isk, M), ipk) \leftarrow \text{KGen}_{\text{IdP}}(pp)</math>  <math>\mathcal{O} := \{\text{RegHRP}, \text{AReq}_{\text{RP}}, \text{ARes}_{\text{IdP}}\}</math>  <math>(auth^*, crid^*, uid^*, ctx^*, sid^*) \leftarrow \mathcal{A}^{\mathcal{O}}(ipk)</math>  Return 1 if <math>\text{ARes}_{\text{IdP}}(isk, auth^*, crid^*, uid^*, ctx^*, sid^*) \neq \perp \wedge</math>  <math>(\cdot, crid^*, sid^*) \notin Q_{\text{auth}}</math></p> <hr/> <p><i>Session Binding</i>: <math>\text{Exp}_{\mathcal{A}, \text{OPPID}}^{\text{SES-BIN}}(\lambda)</math></p> <hr/> <p><math>pp \leftarrow \text{Setup}(1^\lambda); ((isk, M), ipk) \leftarrow \text{KGen}_{\text{IdP}}(pp)</math>  <math>\mathcal{O} := \{\text{RegHRP}, \text{RegCRP}, \text{Alnit}_{\text{U}}, \text{AReq}_{\text{RP}}, \text{ARes}_{\text{IdP}}, \text{ARes}_{\text{Fin}}, \text{Vf}_{\text{RP}}\}</math>  <math>(rid^*, ppid^*, ctx^*, sid^*, \tau_{\text{fin}}^*) \leftarrow \mathcal{A}^{\mathcal{O}}(ipk)</math>  Return 1 if <math>\forall \text{fRP}(ipk, (rid^*, ppid^*, ctx^*, sid^*), \tau_{\text{fin}}^*) = 1 \wedge</math>  (1) <math>(\cdot, ctx^*, sid^*) \notin Q_{\tau}</math> // <i>Direct Forgery</i>  (2) <math>(uid, ctx^*, sid^*) \in Q_{\tau}</math> and at least one of the following holds:  (a) <math>ppid^* \neq F(uid, rid^*)</math> // <i>Nym Correctness</i>  (b) <math>(rid, uid, ctx^*, sid^*) \in Q_{\tau_{\text{fin}}} \wedge rid \neq rid^*</math> // <i>RP Binding I</i>  (c) <math>(rid, ppid, ctx^*, sid^*) \in Q_{\text{vf}} \wedge rid \neq rid^*</math> // <i>RP Binding II</i>  (d) <math>rid^* \notin \text{HRP} \cup \text{CRP}</math> // <i>RP Authentication I</i>  (e) <math>rid^* \in \text{HRP} \wedge (rid^*, \cdot, sid^*) \notin Q_{\text{auth}}</math> // <i>RP Authentication II</i></p>	<p><i>Oracle</i>: <math>\text{RegHRP}(rid)</math>      <i>Oracle</i>: <math>\text{JoinCIdP}(ipk, rid)</math></p> <hr/> <p>Require <math>(rid, \cdot) \notin \text{HRP} \cup \text{CRP}</math>      Require <math>(rid, \cdot) \notin \text{HRP} \cup \text{CRP}</math>  // <i>Both RP and IdP are honest</i>      // <i><math>\mathcal{A}</math> being the corrupt IdP</i>  <math>\langle \text{Join}_{\text{RP}}(ipk, rid), \text{Reg}_{\text{IdP}}(isk, rid, M) \rangle</math> Run <math>\text{Join}_{\text{RP}}(ipk, rid)</math> with <math>\mathcal{A}</math>  Upon output <math>(cred, M')</math>      Upon output <math>cred</math>  HRP := <math>\text{HRP} \cup \{(rid, cred)\}</math>      HRP := <math>\text{HRP} \cup \{(rid, cred)\}</math>  Return 1      Return 1</p> <hr/> <p><i>Oracle</i>: <math>\text{RegCRP}(rid)</math>      <i>Oracle</i>: <math>\text{Alnit}_{\text{U}}(rid)</math></p> <p>Require <math>(rid, \cdot) \notin \text{HRP} \cup \text{CRP}</math>      <math>(orid, crid) \leftarrow \text{Alnit}_{\text{U}}(ipk, rid)</math>  // <i><math>\mathcal{A}</math> being the corrupt RP</i>      // <i>Req. for Unlinkability</i>  <math>\text{Reg}_{\text{IdP}}(isk, rid, M)</math> with <math>\mathcal{A}</math>      <math>Q_{\text{rid}} := Q_{\text{rid}} \cup \{(rid, crid, orid)\}</math>  Upon output <math>M'</math>      Return <math>(orid, crid)</math>  CRP := <math>\text{CRP} \cup \{(rid, \cdot)\}</math>  Return 1</p> <hr/> <p><i>Oracle</i>: <math>\text{AReq}_{\text{RP}}(rid, crid, orid, sid)</math></p> <p>Require <math>(rid, cred) \in \text{HRP}</math>  <math>Q_{\text{auth}} := Q_{\text{auth}} \cup \{(rid, crid, sid)\}</math>      // <i>Req. for Session Binding</i>  Return <math>auth \leftarrow \text{AReq}_{\text{RP}}(ipk, rid, cred, crid, orid, sid)</math></p> <hr/> <p><i>Oracle</i>: <math>\text{ARes}_{\text{IdP}}(auth, crid, uid, ctx, sid)</math></p> <p>Require <math>(\cdot, \cdot, sid) \notin Q_{\tau}</math>  <math>\tau \leftarrow \text{ARes}_{\text{IdP}}(isk, auth, crid, uid, ctx, sid)</math>  If <math>(\tau \neq \perp)</math> then  <math>Q_{\tau} := Q_{\tau} \cup (uid, ctx, sid)</math>      // <i>Req. for Session Binding</i>  If <math>(rid, crid, \cdot) \in Q_{\text{rid}}</math> then      // <i>Req. for Unlinkability</i>  <math>Q_{\text{ppid}} := Q_{\text{ppid}} \cup \{(uid, rid)\}</math>  Else <math>Q_{\text{ppid}} := Q_{\text{ppid}} \cup \{(uid, adv)\}</math>  Return <math>\tau</math></p> <hr/> <p><i>Oracle</i>: <math>\text{ARes}_{\text{Fin}}(auth, crid, uid, ctx, sid)</math></p> <p>Require <math>(\cdot, \cdot, sid) \notin Q_{\tau} \wedge (rid, crid, orid) \in Q_{\text{rid}}</math>  <math>\tau \leftarrow \text{ARes}_{\text{IdP}}(isk, auth, crid, uid, ctx, sid)</math>  <math>(\tau_{\text{fin}}, ppid) \leftarrow \text{AFin}_{\text{U}}(ipk, rid, crid, orid, ctx, sid, \tau)</math>  If <math>(\tau \neq \perp)</math>: <math>Q_{\tau} := Q_{\tau} \cup (uid, ctx, sid)</math> // <i>Req. for Session Binding</i>  If <math>(\tau_{\text{fin}}, ppid) \neq \perp</math>: <math>Q_{\tau_{\text{fin}}} := Q_{\tau_{\text{fin}}} \cup \{(rid, uid, ctx, sid)\}</math>  Return <math>(\tau_{\text{fin}}, ppid)</math></p> <hr/> <p><i>Oracle</i>: <math>\text{Vf}_{\text{RP}}((rid, ppid, ctx, sid), \tau_{\text{fin}})</math></p> <p><math>b \leftarrow \text{Vf}_{\text{RP}}(ipk, (rid, ppid, ctx, sid), \tau_{\text{fin}})</math>  If <math>(b = 1)</math>: <math>Q_{\text{vf}} := Q_{\text{vf}} \cup (rid, ppid, ctx, sid)</math> // <i>Req. for Session Binding</i>  Return <math>b</math></p>
--	---

**Figure 4: Our privacy and security definitions.** *Unlinkability* (UNLINK, exp-1), captures the privacy guarantees of *pseudonymous authentication* towards corrupt RPs. *Unobservability* (UNOBS, exp-2) models that an RP’s *rid* is hidden towards a corrupt IdP during an authentication session. *Request Authentication* (REQ-AUTH, exp-3) captures the security of an authentication request sent to an IdP, ensuring it originates from a registered RP. *Session Binding* (SES-BIN, exp-4) defines the security of the authentication session. Oracles are defined on the right side. All sets are initially empty.

$uid_1$ , as revealing different  $ctx$  would render distinguishing the pseudonym and associated token trivial again. Thus, any implementation of our protocol must ensure that  $ctx$  does not disclose information that could link/identify users behind their pseudonyms.

### 3.3 Unobservability

This property captures that a malicious IdP does not learn anything about the RP’s identity  $rid$  in an authentication request, meaning it cannot observe where the user wants to authenticate to. Specifically, the IdP should not be able to distinguish whether a user repeatedly authenticates to the same RP or different ones. This

property was formally introduced in [31], and we simply adapt this to our notation. The game is represented as  $\text{Exp}_{\mathcal{A}, \text{OPPID}}^{\text{UNOBS}}$  in Figure 4. Unobservability is defined through an indistinguishability game, where the adversary, acting as a corrupt IdP, can setup RPs and obtain their authentication data via the corresponding oracles. Eventually the adversary chooses two RPs  $rid_0$  and  $rid_1$  and receives the authenticated request  $auth_b, crid_b$  of either of them. The adversary wins if it can determine  $b$  better than by guessing.

*Definition 3.2 (Unobservability).* An OPPID scheme satisfies Unobservability if for all PPT adversaries  $\mathcal{A}$ , it holds that

$$\Pr[\text{Exp}_{\mathcal{A}, \text{OPPID}}^{\text{UNOBS}}(\lambda) = 1] \leq 1/2 + \text{negl}(\lambda).$$

### 3.4 Session Binding

The Session Binding property ensures that despite the privacy-preserving computation of an authentication token and pseudonym, the content  $(rid, ppid, ctx, sid)$  in the token is strictly unforgeable and pseudonyms are correctly formed. This comprises the classic unforgeability for all inputs directly seen and vouched for by the IdP – which are  $(ctx, sid)$  for a *Direct Forgery* – but also all blindly signed information. The blindly signed information is  $(rid, ppid)$ , which the IdP vouches for in a session  $sid$  for user  $uid$  with context  $ctx$ . If the IdP indeed created a token for  $(ctx, sid)$ , the blindly signed information must be consistent with its view and the intentions of all honest users and RPs. More precisely, the following must hold:

- If a user  $uid$  requested a token in session  $sid$ , it is infeasible to create a valid token for  $rid, sid$  and  $ppid \neq F(uid, rid)$  – *Nym Correctness*.
- If an honest user  $uid$  intended to authenticate to an RP  $rid$  in session  $sid$ , it must be infeasible to create a valid token for  $sid$  and another  $rid' \neq rid$  – *RP Binding I*.
- If a corrupt user  $uid$  authenticated to an RP  $rid$  in a session  $sid$ , it is infeasible to generate valid authentication tokens for  $sid$  and more than *one* RP – *RP Binding II*.
- If an RP  $rid$  is not properly registered, it is infeasible to generate a valid token for  $rid$  – *RP Authentication I*.
- If an honest RP  $rid$  never authenticated for session  $sid$ , it is infeasible to create a valid token for  $rid$  – *RP Authentication II*.

We model the aforementioned properties through the game  $\text{Exp}_{\mathcal{A}, \text{OPPID}}^{\text{SES-BIN}}$  (see Figure 4), following the classic unforgeability setting. In this game, the IdP is honest, and the adversary can register both corrupt and honest RPs using oracles  $\mathcal{O}.\text{RegHRP}$  and  $\mathcal{O}.\text{RegCRP}$ , respectively, storing their registrations in HRP and CRP. The adversary initiates sessions for honest users via  $\mathcal{O}.\text{Anit}_U$ , obtains authentications from honest RPs via  $\mathcal{O}.\text{AReq}_{RP}$ , and acquires tokens from the IdP using  $\mathcal{O}.\text{ARes}_{IdP}$  (for corrupt users) and  $\mathcal{O}.\text{ARes}_{Fin}$  (for honest users). Additionally, we utilize a verification oracle  $\mathcal{O}.\text{Vf}_{RP}$  to detect if the adversary attempts to reuse the same token across multiple (possibly corrupt) RPs.

The adversary can interact arbitrarily with these oracles and must output a forgery consisting of a valid token  $\tau_{fin}^*$  for session  $(rid^*, ppid^*, ctx^*, sid^*)$  that verifies under  $ipk$ . The adversary wins the game if this forgery is non-trivial, meaning it breaks any of the guarantees listed above, which are captured through dedicated winning conditions.

*Direct and Indirect Forgeries.* First, note that whenever the honest IdP creates a token for a session identified through  $(uid, ctx, sid)$ , these values are stored in  $Q_\tau$ . Thus, in the game, we check if the forgery is for  $(\cdot, ctx^*, sid^*) \notin Q_\tau$ . If this occurs,  $\mathcal{A}$  has produced a token for a session never attested by the honest IdP and wins under condition 1 (*Direct Forgery*).

The second category, an *Indirect Forgery*, means that the IdP did sign  $(ctx^*, sid^*)$ , but the additionally blindly signed or derived values  $rid^*, ppid^*$  are inconsistent with the behavior of the other (honest) parties. This is captured under winning condition 2 and branches according to the properties we discussed earlier. In the following, we focus on Nym Correctness and refer to App. A for a

detailed explanation of the RP Binding and Authentication properties.

*Nym Correctness.* Despite the blind pseudonym computation, a corrupt user  $uid$  must not be able to derive a token for any pseudonym other than the one uniquely defined through  $F(uid, rid)$ , where  $rid$  is the RP specified in the token. This is captured in Condition (a), where  $\mathcal{A}$  wins if  $ppid^* \neq F(uid, rid^*)$ . This condition leverages the fact that the IdP receives  $uid$  as input, and we store  $sid, uid$  in  $Q_\tau$  whenever a token is generated. Therefore, when the adversary outputs  $(rid^*, ppid^*, ctx^*, sid^*)$ , we can look up  $uid$  in  $Q_\tau$  for  $sid^*$  and verify the correctness of the pseudonym for  $uid, rid^*$ . Recall that we already required  $F$  to produce unique pseudonyms, so this precisely defines the one pseudonym that is valid here, and  $\mathcal{A}$  wins if it can produce a valid token for any other pseudonym value.

Note that this property, together with the uniqueness requirement of  $F$ , ensures *sybil-resistance*. This prevents malicious users from exploiting pseudonymous authentication to create several identities towards a single RP, which was not guaranteed in [26].

*RP Binding & Authentication.* The winning condition (b) for RP Binding I exploits that we know the intended  $rid$  when a session  $sid^*$  is started by an honest user through  $\mathcal{O}.\text{Anit}_U$ . Thus, if the adversary outputs a token for any  $rid^* \neq rid$  for such a session, it wins the game. For a session  $sid^*$  initiated by a corrupt user towards a corrupt RP, we never know the exact RP the user wants to authenticate to:  $\mathcal{A}$  invokes  $\mathcal{O}.\text{ARes}_{IdP}$  with adversarially chosen inputs  $auth$  and  $crid$ , both of which hide  $rid$ . Thus, our guarantees are weaker here and follow the spirit of one-more unforgeability: the adversary wins if it has previously "presented" a valid token for some  $rid$  to the  $\text{Vf}_{RP}$  oracle, yet later outputs a token for the same  $sid^*$  but with  $rid^* \neq rid$  as a forgery. Catching such a "double spending" attack is the reason why we have the  $\text{Vf}_{RP}$  oracle: it runs purely on public values, but essentially asks the adversary to commit to one view, and later output a contradicting one as it's forgery (see App. A for further explanation).

The preceding two properties ensure that the IdP-generated token is bound to the blindly received  $rid$ . Additionally, conditions (d) and (e) further ensure that only legitimate RPs can request such tokens. In condition (d), the adversary wins if it manages to produce a valid token for some  $rid^*$  that has never been registered, meaning  $rid^* \notin \text{HRP} \cup \text{CRP}$ . If the  $rid^*$  in the forgery belongs to an honest RP, we further let  $\mathcal{A}$  win if  $(rid^*, sid^*) \notin Q_{auth}$ , indicating that the honest RP had never authenticated for that particular session  $sid^*$ .

*Definition 3.3 (Session Binding).* An OPPID scheme satisfies Session Binding if for all PPT adversaries  $\mathcal{A}$ , it holds that

$$\Pr[\text{Exp}_{\mathcal{A}, \text{OPPID}}^{\text{SES-BIN}}(\lambda) = 1] \leq \text{negl}(\lambda).$$

*Uniqueness of sid.* We remark that our security notion relies on the fact that an honest IdP issues a single token per session  $sid$ . This ensures the desired freshness guarantees and uniquely identifies the session context that the IdP attests. As is typical in such protocols, we therefore assume that  $sid$  is unique per IdP and do not rely on the cryptographic protocol to enforce this. Thus, an implementation of OPPID must implement measures to ensure the freshness of  $sid$  at the application layer.



*Weaknesses of [31].* Our notion builds upon the RP Session Binding model of [31]. Apart of adding Nym Correctness – which is the core functional extension needed for our work – our model significantly enhances the overall security guarantees provided by their notion. We address the following two weaknesses:

Firstly, the original model [31] only ensures security for sessions involving *honest* users. This restriction excludes scenarios involving "corrupt" user sessions, which are critical in real-world applications. The justification for this limitation lies in the blindness of *rid* towards the IdP: it is argued that the game requires the view of honest users to determine their targeted RP. However, this dependency is necessary only for security properties that depend on *rid*. Our security model carefully separates these dependencies into several sub-cases, where only RP Binding I necessitates the restriction to honest users. We demonstrate that even for corrupt users, a weaker form of *rid*-binding should be realized, as expressed in our RP Binding II condition.

Secondly, the original model captures security only against *corrupt* RPs, arguing that honest RPs do not give the adversary any advantage. However, higher security guarantees should ideally apply to sessions involving honest RPs as well. Specifically, the adversary should not be able to create any valid *rid*-bound token for sessions that the RP never authenticated, which we formalize in our RP Authentication II property. We stress that both are oversights in their security model only, as the protocol from [31] also satisfies our stronger security notion. However, related works such as PPOIDC [26] and UPPRESSO [25] do become insecure when users are corrupt, highlighting the need of a security model that properly captures malicious behavior. A more detailed comparison with the security model of [31] is provided in App. A.

### 3.5 Privacy Limitation: No Untraceability

While OPPID significantly enhances privacy in SSO, its guarantees are notably weaker compared to "full-fledged" privacy-preserving authentication systems. The primary limitation lies in the lack of *untraceability*, meaning OPPID does not provide privacy protection when the IdP and RP collude. We discuss these limitations here and argue that they are inherent in any SSO-like system.

When the IdP and (some) RPs collude, they can trace users through several means. First, through the deterministic pseudonyms, which is inherent in any pseudonymous SSO system where the only secret input to the pseudonym computation,  $ppid = F_k(uid, rid)$ , is controlled by the IdP. As *uid* and *rid* are public information and typically stem from small "brute-forceable" sets, a colluding RP *rid* and IdP can determine the user behind an RP-specific pseudonyms through re-computation of the pseudonyms of all users and comparison against the *ppid* they want to identify. This is independent of *how* these pseudonyms are computed, and merely exploits their determinism. Possible means to mitigate that would be to distribute the IdP's key *k* among several IdPs or requiring some cryptographic input from the user. Both would deviate from core principles of SSO though, which relies on a single entity and does not assume the users to manage keys or credentials.

Second, a user's token request and the finalized token can be linked to each other, not only through the pseudonym, but also through the *sid* values known to both the IdP and RP, and through

the timing between sessions that are handled simultaneously by both the IdP and RP. One can design a protocol where *sid* values are not revealed to the IdP in the clear, but the impact would be limited as the sessions can still be linked through the timing information. To avoid that linkage, one would have to break the immediate connection between the IdP and RP, e.g., by letting the IdP issue (somewhat) long-term credentials to users, and rely on techniques such as anonymous credentials for untraceable authentication from the user to the RP. In fact, this approach has already been proposed by EL PASSO [43], but gives up on the convenience advantage of plain SSO as it requires users to manage a long-term key.

## 4 BUILDING BLOCKS

This section introduces the necessary building blocks. The security parameter is denoted as  $\lambda \in \mathbb{N}$ , and the symbol  $\perp$  represents failure. Note that all algorithms may use global parameters  $pp$ , such as shared groups, instead of  $1^\tau$ , and may also provide additional public parameters. For simplicity, we omit explicit mention of these public parameters or the algorithms used to generate them.

*Commitment Scheme.* A commitment scheme  $\text{COM} = (\text{Com}, \text{Open})$  produces a commitment *com* and its corresponding opening *o* using the algorithm  $\text{Com}(m)$ . The algorithm  $\text{Open}(m, \text{com}, o)$  outputs 1 if *com* is a valid commitment for *m*, and 0 otherwise. The commitment scheme must satisfy hiding and binding properties.

*Non-interactive Zero-Knowledge Proofs.* In a *non-interactive* zero-knowledge proof system [7, 20], the prover and verifier possess the statement *s* and some public context *x*. The prover generates a proof  $\pi \leftarrow \text{NIZK}\{(w) : s(w)\}(x)$  that convinces the verifier that  $s(w) = 1$ , without revealing *w* to the verifier and ensuring that  $\pi$  is bound to *x*. We require the proof system to be zero-knowledge and simulation-sound [23]. Our NIZK instantiation uses generalized Schnorr-proofs [9], made non-interactive through the Fiat-Shamir heuristic [20], including *x* in the challenge hash.

*Signature Scheme.* A signature scheme is a tuple of algorithms  $S_1 = (\text{KGen}, \text{Sign}, \text{Vf})$ , with key generation  $(pk, sk) \leftarrow \text{KGen}(1^\lambda)$ , signing  $\sigma \leftarrow \text{Sign}(sk, m)$ , and verification as  $0/1 \leftarrow \text{Vf}(pk, m, \sigma)$ . We need  $S_1$  to be Existentially Unforgeable under a Chosen Message Attack (EUF-CMA) [22]. In our implementation, we use RSA signatures for compatibility with existing standards.

We also require a signature scheme  $S_2 = (\text{KGen}, \text{Sign}, \text{Vf})$  that supports the creation of efficient NIZKs. The NIZK should prove knowledge of a valid signature  $\sigma$  on a message *m* under *pk* without revealing the message or signature. In our construction, we combine  $S_2$  signatures with commitments using a NIZK proof that demonstrates knowledge of a signature on a committed message as

$$\text{NIZK}\{(m, \sigma, o) : \text{Vf}(pk, m, \sigma) = 1 \wedge \text{Open}(m, \text{com}, o) = 1\}(\text{com}).$$

This proof discloses only the commitment *com* while verifying the possession of a valid signature  $\sigma$  under *pk* on *m* and an opening *o* to the commitment *com* for the signed message. We instantiate this scheme with PS signatures [38], which provide all these features.

*Pseudorandom Functions.* We require a pseudorandom function  $y \leftarrow \text{PRF}(k, x)$  that produces output indistinguishable from random, towards an adversary not knowing the key *k*. We need two

different pseudorandom functions, one that produces pseudorandom values in  $\mathbb{Z}_q$  and can simply be HMAC with a proper output mapping; and a second function that allows for the partially-blind evaluation needed in our protocol. For the latter, we use the DL-based  $\text{PRF}(k, x) := H(x)^k$  [34] in a group  $\mathbb{G}$  of prime order  $q$ .

## 5 OUR OPPID CONSTRUCTION

This section introduces our OPPID protocol  $\pi_{\text{OPPID}}$ , which combines oblivious  $ppid$  generation with a recent privacy-preserving RP authentication approach for the OIDC Implicit Flow [31]. We first outline the adapted authentication process from [31], present the construction of our pseudonym function and its semi-blind evaluation in the context of joint SSO authentication, and then proceed with the security analysis of our protocol. The detailed protocol, including oblivious  $ppid$  generation, is given in Figure 5.

### 5.1 Privacy-Preserving RP Authentication

Our protocol  $\pi_{\text{OPPID}}$  builds upon [31], which enables privacy-preserving RP authentication in SSO. The core idea therein is that an RP obtains obtaining a privacy-preserving credential from the IdP that includes the RP's identifier  $rid$ . This credential uses a signature scheme with efficient proofs, enabling the RP to authenticate to the IdP in a blind yet verifiable manner by sending a commitment to  $rid$  and proving ownership of a valid signature on  $rid$ . The IdP then verifies the proof and signs the commitment as part of the authentication token. Both the user and RP know the opening to the commitment and can verify that the token is indeed issued for the intended RP. The authentication token in [31] always contains the username  $uid$  that is vouched for by the IdP. In our protocol, the key modification is replacing  $uid$  with the pseudonym  $ppid$ , computed in a blind yet controlled way, detailed in Sec. 5.2.

The protocol in [31] also captures revocation by making the RP's credentials short-lived and encoding an epoch that must be revealed in every authentication. On the protocol level, adding epochs to credentials is very simple, but it makes the security model and analysis significantly more complex. Thus, we only use the core idea in our protocol and use it as a basis for integrating our privacy-preserving  $ppid$  generation. The main steps we use from [31] are related to setup, registration, and basic authentication.

*Setup & Registration.* The IdP generates key pairs for two signature schemes:  $(sk_1, pk_1) \leftarrow_{\mathbb{R}} S_1.\text{KGen}(1^\lambda)$ , a standard scheme used for signing the authentication token, and  $(sk_2, pk_2) \leftarrow_{\mathbb{R}} S_2.\text{KGen}(1^\lambda)$ , which supports efficient proofs for RP authentication. During the registration of an RP with identifier  $rid$ , the IdP issues a credential  $cred := \sigma_{rid}$  for  $rid$ .

The public parameters  $pp \leftarrow \text{OPPID.Setup}(1^\lambda)$  serve as implicit inputs for all algorithms. They include the descriptions of the underlying groups and the potential public parameters of the commitment scheme  $\text{Com}$ , signature schemes  $S_2$  and  $S_1$ , and the zero-knowledge proof system  $\text{NIZK}$ .

*Basic Authentication – From [31] Without Pseudonyms.* For authentication of a user  $uid$  to an RP  $rid$  in a session identified through  $sid$ , the user, RP, and IdP proceed as follows:

$\text{KGen}_{\text{IdP}}(pp) \rightarrow ((isk, \mathcal{M}), \text{ipk})$
$(sk_1, pk_1) \leftarrow_{\mathbb{R}} S_1.\text{KGen}(1^\lambda)$ ; $(sk_2, pk_2) \leftarrow_{\mathbb{R}} S_2.\text{KGen}(1^\lambda)$ ; $k \leftarrow_{\mathbb{R}} \{0, 1\}^\lambda$ Return $((sk_1, sk_2, k), \emptyset), (pk_1, pk_2)$
$\langle \text{Join}_{\text{RP}}(\text{ipk}, \text{rid}), \text{Reg}_{\text{IdP}}(\text{isk}, \text{rid}, \mathcal{M}) \rangle \rightarrow \{(cred, \mathcal{M}'), \perp\}$
RP : Initiate registration for $rid$ IdP : Parse $isk$ as $(\cdot, sk_2, \cdot)$ ; Require $rid \notin \mathcal{M}$ $\sigma_{rid} \leftarrow S_2.\text{Sign}(sk_2, \text{rid})$ ; $\mathcal{M}' \leftarrow \mathcal{M} \cup \{rid\}$ RP : Return $\sigma_{rid}$ ; IdP : Return $\mathcal{M}'$
$\text{Alnit}_{\text{U}}(\text{ipk}, \text{rid}) \rightarrow (\text{orid}, \text{crid})$
$(com, o) \leftarrow_{\mathbb{R}} \text{Com}(rid)$ ; $r \leftarrow_{\mathbb{R}} \mathbb{Z}_q^*$ ; $\bar{x} \leftarrow H(rid)^r$ Return $((r, o), (\bar{x}, com))$
$\text{AReq}_{\text{RP}}(\text{ipk}, \text{rid}, cred, \text{crid}, \text{orid}, \text{sid}) \rightarrow \text{auth}$
Parse $cred$ as $\sigma_{rid}$ , $ipk$ as $(\cdot, pk_2)$ , $crid$ as $(\bar{x}, com)$ , $orid$ as $(r, o)$ Require $H(rid)^r = \bar{x} \wedge \text{Open}(rid, com, o) = 1$ Return $\pi \leftarrow \text{NIZK}\{(rid, o, \sigma_{rid}) : S_2.\text{Vf}(pk_2, rid, \sigma_{rid}) = 1 \wedge \text{Open}(rid, com, o) = 1\}(com, \bar{x}, sid)$
$\text{ARes}_{\text{IdP}}(\text{isk}, \text{auth}, \text{crid}, \text{uid}, \text{ctx}, \text{sid}) \rightarrow \{\tau, \perp\}$
Parse $ipk$ as $(\cdot, pk_2)$ , $isk$ as $(sk_1, \cdot, k)$ , $crid$ as $(\bar{x}, com)$ , $auth$ as $\pi$ Require that $\pi$ verifies w.r.t. $(pk_2, com, \bar{x}, sid)$ and $\bar{x} \in \mathbb{G}$ $uk \leftarrow \text{PRF}(k, \text{uid})$ ; $\bar{y} \leftarrow \bar{x}^{uk}$ $\sigma_\tau \leftarrow S_1.\text{Sign}(sk_1, (com    \bar{x}    \bar{y}    \text{ctx}    \text{sid}))$ ; Return $(\sigma_\tau, \bar{y})$
$\text{AFin}_{\text{U}}(\text{ipk}, \text{rid}, \text{crid}, \text{orid}, \text{ctx}, \text{sid}, \tau) \rightarrow \{(\tau_{\text{fin}}, ppid), \perp\}$
Parse $ipk$ as $(pk_1, \cdot)$ , $crid$ as $(\cdot, com)$ , $orid$ as $(r, o)$ , $\tau$ as $(\sigma_\tau, \bar{y})$ $\bar{x} \leftarrow H(rid)^r$ ; $y \leftarrow \bar{y}^{-r}$ Require $\text{Open}(rid, com, o) = S_1.\text{Vf}(pk_1, (com    \bar{x}    \bar{y}    \text{ctx}    \text{sid}), \sigma_\tau) = 1$ Return $((com, o, r, \bar{y}, \sigma_\tau), y)$
$\text{Vf}_{\text{RP}}(\text{ipk}, (rid, ppid, \text{ctx}, \text{sid}), \tau_{\text{fin}}) \rightarrow 0/1$
Parse $ipk$ as $(pk_1, \cdot)$ , $\tau_{\text{fin}}$ as $(com, o, r, \bar{y}, \sigma_\tau)$ $\bar{x} \leftarrow H(rid)^r$ ; $y \leftarrow \bar{y}^{-r}$ Return 1 if $\text{Open}(rid, com, o) = S_1.\text{Vf}(pk_1, (com    \bar{x}    \bar{y}    \text{ctx}    \text{sid}), \sigma_\tau) = 1 \wedge ppid = y$

Figure 5:  $\pi_{\text{OPPID}}$  protocol construction of our OPPID system.

- (1) Initialization: The user creates a commitment  $com$  on the intended RP's  $rid$ , sends  $com$  and the opening  $o$  to the RP, and keeps all values in its state for later finalization.
- (2) RP Authentication: When the RP receives a well-formed  $com$  and  $o$  for its  $rid$ , it generates its authentication  $auth$  by proving possession of the IdP's  $S_2$  signature  $\sigma_{rid}$  on  $rid$  and proving that  $rid$  is also contained in  $com$ . It sends the proof  $\pi$  to the user, who forwards it along with the commitment  $com$  to the IdP.
- (3) Token Generation: When the IdP receives the commitment  $com$  and the proof  $\pi$  from the user  $uid$ , it verifies the proof, ensuring that authentication is requested by a registered RP. If verification succeeds, the IdP generates a standard  $S_1$  signature  $\sigma_\tau$  on the commitment  $com$ , the session  $sid$ , and the context  $ctx$ .
- (4) Finalization & Verification: To finalize the token  $\tau$ , the user opens the commitment  $com$  with  $o$  and the intended  $rid$  and verifies the IdP's signature. If the verification is successful, the opening  $o$  is added to  $\tau$ , creating a verifiable binding of the IdP's signature to a specific RP.

### 5.2 Oblivious PPID Generation

The protocol outlined above would not be very useful yet, as it does not include a user identifier, which was simply  $uid$  in [31]. We now want to include a pseudonym  $ppid = F(uid, rid)$  in every token,

where  $uid$  is the known user that the IdP has authenticated (outside of our protocol), and  $rid$  is the RP the user wants to authenticate to, but which must not be revealed to the IdP. We first describe the core function  $F$  and then explain how to compute it in a semi-blind and controlled manner as our model requires.

*Pseudonym Function  $F$ .* Our pseudonym function  $F$  is a keyed deterministic function, which combines the standard DL-based pseudorandom function  $F_{DL}(k, x) := H(x)^k$  [34], operating in a group  $\mathbb{G}$  of prime order  $q$ , and a hash function  $H : \mathcal{R} \mapsto \mathbb{G} \setminus \{1\}$ , with a standard PRF  $: \{0, 1\}^\lambda \times \mathcal{U} \mapsto \mathbb{Z}_q$ . The sets  $\mathcal{U}$  and  $\mathcal{R}$  represent the user and RP space, respectively. For compactness, we will sometimes write  $F_k(uid, rid)$  to refer to this keyed function:

$$F_{DL+PRF} = F_k(uid, rid) := H(rid)^{\text{PRF}(k, uid)}.$$

The desired pseudonym unlinkability is directly ensured as  $F_{DL+PRF}$  is a secure PRF. In fact,  $F_{DL}$  is often used to derive so-called scope-exclusive pseudonyms in the context of anonymous credentials, group signatures and DAA [8, 14, 33]. In these works, the exponent is a user-managed secret key, while we rely on the IdP to maintain them. Note that UPPRESSO [25] employs a similar keyed pseudonym function, and we refer for the comparison to Sec. 6.

*Partially-Blind Evaluation of  $F$ .* The construction  $F_{DL+PRF}$  has also been used as a *partially-blind OPRF* in prior works [12, 30, 32]. This gives us the needed capability for blind evaluation of the function on a hidden  $rid$  and a revealed  $uid$ : the user blinds the inner hash as  $\bar{x} \leftarrow H(rid)^r$  for a random  $r$  and sends this blinded value along with  $uid$  to the IdP. The IdP responds with  $\bar{y} \leftarrow \bar{x}^{\text{PRF}(k, uid)}$ , where the exponent depends on the revealed  $uid$ . The user can then unblind the response to  $ppid \leftarrow \bar{y}^{-r}$ , which yields the expected pseudonym:

$$ppid = \bar{y}^{-r} = \left( (H(rid)^r)^{\text{PRF}(k, uid)} \right)^{-r} = H(rid)^{\text{PRF}(k, uid)}.$$

*The Need for Verifiability.* Finally, another essential feature required for our function is verifiability: the RP must be assured that a received  $ppid$  was computed for the correct  $rid$ . While the IdP is trusted here, one approach could be to delegate this verification task to the IdP by employing a partially-blind OPRF with *committed and verifiable inputs*. This approach would enable the IdP to verify that it evaluates the blind function on the same  $rid$  that is authenticated through *auth* (and contained in *com*). Such a function would serve as a suitable building block, but it is not known whether such an OPRF is feasible [10]. Existing constructions like the Dodis-Yampolskiy (O)PRF [15], which operate on homomorphically encrypted inputs, allow for proofs of well-formedness but do not extend to the partially-blind setting needed here. Similarly, constructions like (2)HashDH, which enable partial blindness, lack efficient and composable proofs of correct inputs as their input is a perfectly blinded hash value that destroys all algebraic structure.

*Adding Partial Verifiability to Our Protocol.* Interestingly, we can work around this non-existent building block by incorporating several straightforward steps, building upon UPPRESSO [25] and the base protocol from [31]. The resulting protocol is detailed in Figure 5.

First, in addition to the commitment and opening  $com, o$  for  $rid$ , we let the user compute  $\bar{x} = H(rid)^r$  and send all values, including

$r$  and  $o$ , to the RP. The RP verifies that both  $\bar{x}$  and  $com$  open to its  $rid$ . Only upon successful verification does the RP provide its authentication *auth*, proving ownership of a valid credential for the committed  $rid$  in *com*. The RP also binds its NIZK proof  $\pi$  to  $\bar{x}$ . This ensures  $\bar{x}$  correctness when either the RP or user is honest.

Second, upon receiving a verified authentication request, the IdP includes both the blinded input  $\bar{x}$  and the blinded output  $\bar{y}$  in its token, as done in [25], signing  $\sigma_\tau \leftarrow S_1.\text{Sign}(sk_1, (com || \bar{x} || \bar{y} || ctx || sid))$ . Crucially, the signature binds the blinded and non-verified values  $\bar{x}$  and  $\bar{y}$  used for  $ppid$  to the commitment  $com$  on  $rid$ , for which the RP provided a valid NIZK proof.

Third, the user incorporates the blinding value  $r$  used to hide  $rid$  in  $\bar{x}$  as part of the final token  $\tau_{fin}$ . Thus,  $\tau_{fin} = (com, o, r, \bar{y}, \sigma_\tau)$ , and the verification function  $Vf_{RP}(ipk, (rid, ppid, ctx, sid), \tau_{fin})$  performs the following crucial checks:

- Verify that the IdP's signature  $\sigma_\tau$  is valid for the recomputed  $\bar{x} = H(rid)^r$ , where  $rid$  is the one provided in the verification.
- Ensure that the pseudonym  $ppid$  satisfies  $ppid = \bar{y}^{-r}$ , where  $\bar{y}$  is signed by the IdP and  $r$  is the blinding value leading to the correct  $\bar{x}$ .
- Confirm that the  $rid$  contained in  $\bar{x}$  matches the one in *com*.

These checks extend the guarantees from *com* to *ppid*, leveraging our three-party setting where both the user and RP are aware of  $rid$ , and both verify that  $\bar{x}$  and  $com$  are valid for the intended  $rid$ . Thus, as long as either the RP or the user in the session remains honest, the derived pseudonym  $ppid = F(uid, rid)$  is ensured to be correct.

*Invalid Pseudonyms – If Both RP and User Are Corrupt.* If both the RP and user are corrupt, they have some leeway, but none that is harmful. A malicious RP  $rid$  and a malicious user  $uid$  can request and obtain pseudonyms  $ppid = F(uid, rid')$  for arbitrary  $rid' \neq rid$  by sending a blinded  $\bar{x}$  to the IdP that does not contain the correct  $rid$ . The IdP will compute the pseudonym based on the incorrect  $rid'$  but will bind it to the verified commitment  $com$ , which can only be opened to  $rid'$ . In the final token, the  $rid$  is no longer blinded, and verification includes a check whether the commitment contains the same  $rid$  as the blinded  $\bar{x}$  used in the  $ppid$  computation. This check will fail, rendering the entire token and pseudonym invalid.

Furthermore, note that this "attack" is only feasible for malicious users  $uid$ , as  $uid$  is revealed to the IdP and used to compute  $ppid$ . Thus, a malicious RP and user  $uid$  cannot trick the IdP into computing pseudonyms for any other user  $uid' \neq uid$ .

### 5.3 Security Analysis

We have already informally sketched how the different security properties are guaranteed in our protocol description. Now, we formally prove that our protocol  $\pi_{\text{OPPID}}$  satisfies all security and privacy properties defined in Sec. 3.

*Properties of  $F_{DL+PRF}$ .* Let us first analyze the properties of our pseudonym function  $F_{DL+PRF} = F_k(uid, rid) := H(rid)^{\text{PRF}(k, uid)}$ . Recall that our OPPID model requires this function to provide unique, collision-free, and unlinkable pseudonyms (see Sec. 2.2). It is easy to see that all properties are satisfied due to  $F_{DL+PRF}$  being deterministic, injective, and a secure PRF. These properties are further elaborated in App. B.

*Analysis of  $\pi_{\text{OPP}\text{ID}}$ .* We now turn to the proofs of our three core properties, ensuring the correct yet privacy-preserving computation of pseudonymous authentication tokens with respect to our pseudonym function  $F_{\text{DL}+\text{PRF}}$ .

**THEOREM 5.1 (UNLINKABILITY).**  *$\pi_{\text{OPP}\text{ID}}$  satisfies Unlinkability if  $H$  is a random oracle, PRF is a secure pseudorandom function and the DDH assumption holds in  $\mathbb{G}$ .*

This proof relies on the pseudorandomness of  $F_{\text{DL}+\text{PRF}}$ , as shown in App. B.2, under the assumptions that  $H$  is a random oracle, PRF is a secure pseudorandom function, and the DDH assumption holds in  $\mathbb{G}$ . We now provide a proof sketch and refer to App. B.1 for the full proof, where we also discuss why a one-more-type assumption, often required for OPRFs, is not required.

**PROOF SKETCH.** In the Unlinkability game, the adversary's objective is to determine the user  $uid_b$  behind a pseudonym  $ppid_b$  and token  $\tau_{\text{fin}b}$ , generated for  $rid$  and either  $uid_0$  or  $uid_1$ . The adversary has oracle access to the honest IdP and can learn the pseudonyms of  $uid_0$  and  $uid_1$  for all RPs except  $rid$  (as otherwise winning is trivial).

We already know that  $F_{\text{DL}+\text{PRF}}$  is a secure PRF, meaning the  $ppids$  themselves do not leak any information about the contained  $uid$ , except what is deterministically derived. The only part in the IdP's response  $\tau_{\text{fin}}$  that depends on  $uid$  is the PRF output  $\bar{y} = \bar{x}^{\text{PRF}(k,uid)}$ , where  $\bar{x}$  is the value  $F_k(uid, rid) = H(rid)^{\text{PRF}(k,uid)}$  blinded with a random  $r$ . Note that Unlinkability holds for honest users only, ensuring that  $\bar{x} = H(rid)^r$  in the challenge query is a valid input. Thus, the token does not provide the adversary with any information beyond  $ppid$ .

What remains to be shown is that a malicious RP, possibly colluding with a malicious user  $uid^*$ , cannot exploit the partially blind evaluation of  $F_{\text{DL}+\text{PRF}}$  to obtain dedicated  $ppids$  of either of the honest challenge users  $uid_0$  or  $uid_1$  illegitimately. Specifically, they cannot obtain their pseudonyms through oracle queries *not* intended for either  $uid_0$  or  $uid_1$  (as for the challenge users, the oracles enforce honest user behavior and honestly generated inputs  $\bar{x}$ ).

It is easy to see that this scenario is infeasible because the  $ppid$  depends on the  $uid$ , which the IdP learns in clear and uses in its computation. Therefore, there is no opportunity to manipulate the  $uid$  and its impact on the  $ppid$  computation. While a malicious RP and user could potentially trick the IdP into computing a pseudonym for an arbitrary  $rid$  that does not match the one authenticated via *auth*, they can only do so for a malicious  $uid^* \neq uid_0, uid_1$ , which does not provide any advantage in winning the Unlinkability game.  $\square$

**THEOREM 5.2 (UNOBSERVABILITY).**  *$\pi_{\text{OPP}\text{ID}}$  satisfies Unobservability if COM is hiding, and the NIZK is zero-knowledge.*

This proof is essentially the same as in [31]. It follows from the fact that the IdP receives the  $rid$  in a commitment and within a zero-knowledge proof. The only difference here is that  $\mathcal{A}$  also receives  $\bar{x} = H(rid)^r$ , which is the blinded hash of  $rid$ . As the blinding is information-theoretic, no additional assumptions are needed. We provide a simple proof in App. B.2 for completeness.

**THEOREM 5.3 (SESSION BINDING).**  *$\pi_{\text{OPP}\text{ID}}$  satisfies Session Binding if the  $S_1$  and  $S_2$  scheme are EUF-CMA secure, COM is binding, and the NIZK is zero-knowledge and simulation extractable.*

We give a proof sketch below and refer for the full proof to App. B.3.

**PROOF SKETCH.** The proof branches along the winning conditions in the Session Binding game. Direct forgeries, condition (1), require the adversary to output a  $S_1$  signature on a fresh  $(\cdot, ctx^*, sid^*)$ , which leads to a  $S_1$  forgery.

For Nym Correctness, condition (2a), it must be infeasible to output  $(rid^*, ppid^*, ctx^*, sid^*, \tau_{\text{fin}}^*)$ , where  $\tau_{\text{fin}}^*$  is a valid token for a pseudonym  $ppid^* \neq F(uid, rid^*)$ . That is, the pseudonym in the forgery does not match the expected and unique pseudonym  $F(uid, rid^*)$  for the user  $uid$  that the honest IdP had seen in the session  $sid^*$ , and the  $rid^*$  the token is targeted for.

Recall that the final token  $\tau_{\text{fin}}^*$  contains  $r^*$ ,  $\bar{y}^*$ , and  $\sigma_{\tau}^*$ , and being valid implies that  $ppid^* = \bar{x}^* = \bar{y}^{*-r^*}$  and  $\sigma_{\tau}^*$  is a valid signature on  $\bar{x}^*$ ,  $\bar{y}^*$ ,  $ctx^*$ , and  $sid^*$  vouched by the IdP. If  $\bar{x}^*$ ,  $\bar{y}^*$  (in  $\tau_{\text{fin}}^*$ ) are the values indeed seen by the IdP in session  $sid^*$  and the verification checks that  $\bar{x}^* = H(rid^*)^{r^*}$  and  $ppid^* = \bar{y}^{*-r^*}$ , they uniquely determine the correct pseudonym value for which verification succeeds.

RP Binding 1, condition (2b), implies that the adversary was able to present an honest user's token for a different  $rid$  than it was intended for by the honest user. The adversary can succeed in this case only by finding an opening collision for the commitment  $com^*$  in the finalized token of the honest user, which occurs negligibly due to the binding property of COM.

RP Binding 2, condition (2c), occurs when the adversary can create two valid tokens with the same  $(ctx^*, sid^*)$ , but distinct  $rid$  and  $rid^*$  values. As the IdP issues a token for  $sid$  values only once, the adversary needs to either forge a  $S_1$  signature on  $sid^*$  with a new commitment to  $rid^*$  or find an opening collision for  $com^*$  for distinct  $rid$  and  $rid^*$  values, breaking the binding property of COM.

RP Authentication 1, condition (2d), occurs when the adversary forges a valid token for a non-registered  $rid$  value. This requires forging a  $S_2$  signature for the non-registered  $rid^*$ , finding a COM collision and using the credential on another  $rid$  value than  $rid^*$  to create the NIZK proof, or creating a NIZK proof for an invalid statement, i.e., without knowing a valid credential on  $rid^*$  at all.

RP Authentication 2, condition (2e), is similar to condition (2d) but considers an adversary who forges a valid token for honest RPs by using corrupted RPs. This case closely follows condition (2d), except for a technical detail about the required extractability notion from the NIZK, which is explained in detail in App. B.3.  $\square$

**THEOREM 5.4 (REQUEST AUTHENTICATION).**  *$\pi_{\text{OPP}\text{ID}}$  satisfies Request Authentication if the  $S_2$  scheme is EUF-CMA secure, and the NIZK is zero-knowledge and simulation extractable.*

The proof of Request Authentication essentially follows [31] and is given in App. B.4 for completeness. This property ensures that it is hard to forge an authentication request  $auth^*$  for a fresh tuple  $(crid^*, sid^*)$ . In our construction, we set  $crid = (\bar{x}, com)$  to the blinded pseudonym input and commitment on  $rid$ , and  $auth$  is a NIZK  $\pi$  of a valid credential  $cred$  for  $rid$ . The proof  $\pi$  is also bound to  $(crid, sid)$ . It easily follows that without having a valid credential, which is a  $S_2$  signature on the honest RP's  $rid$ , an adversary cannot compute such a valid proof  $auth^*$  for a fresh  $(crid^*, sid^*)$ .

## 6 EVALUATION & DISCUSSION

In this section, we present our prototypical implementation of  $\pi_{\text{OPPID}}$  and compare it to the most related work, evaluating both efficiency and security.

### 6.1 Security Comparison with Related Protocols

The works most closely related to ours are standard OIDC with pseudonyms [41], AIF-ZKP [31], PPOIDC [26], and UPPRESSO [25]. These protocols have been selected for a detailed comparison as they also operate within the plain-SSO model, meaning they do not require the user to manage any long-term keys or credentials, nor do they rely on additional parties or dedicated hardware modules.

*OIDC* [41]. OIDC is the most widely deployed SSO protocol for user authentication, supporting both RP authentication and RP-specific pseudonyms. According to its specification [41, §8.1], the IdP creates a pseudonym  $ppid$  as  $H(uid||rid||k)$ , where  $H$  is a cryptographic hash function and  $k$  is a secret, high-entropy random string held by the IdP. Upon receiving an authenticated request from the RP  $rid$  via the user  $uid$ , the IdP computes  $ppid$  and signs  $rid$ ,  $ppid$ , and the session data. While this provides Request Authentication, Session Binding, and Unlinkability, OIDC does not achieve Unobservability.

*AIF-ZKP* [31]. This work, outlined in Sec. 5.1, provides the foundation of our protocol to achieve Session Binding, including strong RP authentication and Unobservability simultaneously. Our Session Binding model is stronger than that in [31], and our analysis shows that the original protocol already satisfied this stronger notion as well. The protocol reveals  $uid$  in clear text to RPs and does not support pseudonyms, thus failing to provide Unlinkability.

*PPOIDC* [26]. The PPOIDC protocol aims to turn OIDC into an unobservable protocol. To blindly bind the IdP’s token to a particular  $rid$  and compute the  $ppid$  it mostly relies on hash functions, serving as commitments. More precisely, when the user wants to authenticate to  $rid$ , it first computes  $com_{rid} := H(rid||r)$  for a random  $r \leftarrow_{\mathbb{R}} \{0, 1\}^\lambda$ . The pseudonym computation uses a hash function again and also makes the non-standard assumption that  $uid$  is a high-entropy value that the user retrieves from the IdP at every login. The user computes  $ppid := H(uid||rid)$  and the commitment  $com_{ppid} := H(ppid||r')$  with  $r' \leftarrow_{\mathbb{R}} \{0, 1\}^\lambda$ . The user also generates a zero-knowledge proof  $\pi$  that  $com_{ppid}$  is derived for her  $uid$  and sends  $com_{rid}$ ,  $com_{ppid}$  and  $\pi$  to the IdP.

The IdP verifies that the proof is valid for the authenticated  $uid$ , and then signs  $com_{rid}$  and  $com_{ppid}$  in its token. The user forwards the IdP’s signature, randomness  $r, r'$  and her  $ppid$  to the RP, which checks that the IdP-signed commitment  $com_{rid}$  correctly opens to its own  $rid$  and  $com_{ppid}$  opens to  $ppid$ .

The protocol does not consider RP authentication towards the IdP and thus cannot satisfy these parts of the Session Binding property or Request Authentication. RP authentication is not entirely missing either though, as their protocol issues certificates to the RP upon registration and relies on the user to verify them when they start a session. This only provides security if all users are honest though, and we discuss the difference to our IdP-centric authentication at the end of this section. The PPOIDC still partially

satisfies Session Binding, as the IdP blindly signs the commitment  $com_{rid}$  of an user-verified  $rid$ , which ensures RP Binding I and II.

However, PPOIDC does not achieve *Nym Correctness*: The proof  $\pi$  only ensures that  $ppid$  is computed on the correct  $uid$ , but does not guarantee that the committed  $rid$  in  $com_{rid}$  matches the one used to compute  $ppid$ , allowing corrupt users to obtain arbitrary IdP-certified pseudonyms.

*UPPRESSO* [25]. The protocol focuses solely on the blindly computed pseudonyms in SSO, which are computed as  $ppid := rid^{k_{uid}}$  where  $k_{uid}$  is a user-specific secret key in  $\mathbb{Z}_q$  maintained by the IdP. Instead of hashing  $rid$  to the group as in  $\pi_{\text{OPPID}}$ , their protocol relies on  $rid$  already being a proper and random group element. This is done by letting the IdP issue a certificate on a randomly chosen group element  $rid \in \mathbb{G}$  (where  $\mathbb{G}$  is a cyclic group of order  $q$ ) to the RP when it registers.

When a user wants to authenticate to an RP, it receives and verifies the certified  $rid$  from the RP and blinds it as  $\overline{rid} := rid^r$  using a random  $r$ . The IdP then receives  $\overline{rid}$  from the user  $uid$ , computes  $\overline{ppid} := \overline{rid}^{k_{uid}}$ , and signs both  $\overline{rid}$  and  $\overline{ppid}$  in its token. The RP then receives all signed values and  $r$  and unblinds them to  $ppid := rid^{k_{uid}}$ .

In terms of security, UPPRESSO achieves both privacy-related properties due to perfect blinding of  $rid$  and pseudonyms computed via a classic DL-based PRF. The protocol also achieves some form of RP-Binding, as the signed  $\overline{rid}$  could serve as a commitment to  $rid$  which can be verified using  $r$  (this is not made publicly verifiable in their protocol though). In contrast to our scheme, this would additionally require to also verify that  $rid$  is the correct group element – which was simply computing  $H(rid)$  in our scheme. However, UPPRESSO does not achieve RP authentication as part of Session Binding or Request Authentication, as the RP does not authenticate to the IdP.

Comparing UPPRESSO to our protocol, we made three key changes to the pseudonym computation: First, we set  $k_{uid} := \text{PRF}(k, uid)$ , where  $k$  is a secret key held by the IdP. This eliminates the need for the IdP to manage a secret key table that grows linearly with the number of users. Second, we compute the pseudonym on  $H(rid)$  instead of a *certified*  $rid \in \mathbb{G}$  directly, removing the need for users to verify a certificate on  $rid$  to check its validity. Using a malformed  $rid$  will allow malicious RPs to link users, which is prevented through our hash computation. Third, our protocol ensures that any valid and publicly verifiable IdP token contains the correct and RP-authenticated pseudonym through consistency checks between the verified commitment and blinded input. In UPPRESSO, the final token can include malformed pseudonyms when corrupt users and RPs collude, allowing the RP to falsely claim an inflated user base asserted by the IdP. Other changes primarily focus on providing IdP-side RP authentication and ensuring the final token’s public verifiability for the targeted  $rid$ .

We note that the recently proposed BISON protocol [27] for blindly computed SSO pseudonyms follows essentially the same approach as UPPRESSO as well. Although the authors claim a more generic approach, their construction boils down to the same protocol already taken by UPPRESSO, and our analysis and comparison transfers to BISON too.

*User-Side RP Authentication.* In PPOIDC and UPPRESSO, RP authentication has been shifted from an IdP-verified setting to a user-verified one. In this approach, the user receives an RP’s certificate to verify that it is properly registered with the IdP and provides the correctly certified values for the cryptographic protocol. UPPRESSO requires two RP and IdP scripts that the user receives from each party, which handle certificate transfer and verification within the session.

Interestingly, both protocols suggest sending the plain certificate from the RP to the user, i.e., without binding the certificate to a key and session nonce. This approach would immediately allow phishing attacks if no additional cross-verification of consistent certificates from the authenticated TLS session and the verified SSO-protocol values is performed.

Furthermore, relying on proper RP authentication by the user compromises full IdP control over the token it issues. As discussed above, the IdP can be tricked into signing tokens for malformed pseudonyms or invalid *rids*, violating the correctness and non-repudiation guarantees typically expected from such an IdP. Overall, handling RP authentication on the user side is rather fragile and requires trust in the honest execution on the user’s device. Therefore, our OPPID system aims for IdP-side RP authentication.

*RP Revocation.* Achieving Unobservability, which means hiding the *rid* from the IdP during RP authentication, rules out classic revocation strategies based on blacklisting the revoked *rids*. This is a well-understood challenge and has been addressed by [31] through an epoch-based strategy. In this approach, the RP credentials used in anonymous authentication requests are short-lived, lasting, e.g., for a week or a month. Thus, the RP must regularly re-obtain its membership credential from the IdP, which will refuse to do so if the RP has been revoked. The short-lived credentials can be realized by having the IdP also sign the current epoch along with the RP’s *rid* in the membership credential. During authentication, the *rid* remains hidden, but the epoch must be revealed and valid.

Since  $\pi_{\text{OPPID}}$  is built upon [31], integrating this revocation mechanism only requires adding an epoch to the issued RP credential. This is straightforward and has no impact on how the *ppid* is generated or verified. We omitted revocation to simplify the security model and focus on oblivious *ppid* computation.

## 6.2 Implementation and Evaluation

We now report on the implementation of our protocol and compare its efficiency to related works.

*Instantiation of Building Blocks.* We instantiated  $\pi_{\text{OPPID}}$  as follows: for the IdP’s standard signature  $S_1$ , we used RSA-SHA256 2048-bit to comply with current industry standards. Building upon [31], we used PS signatures [38] on curve BLS12-381 [6] for  $S_2$ , which includes the bilinear group description  $(q, \mathbb{G}_1, \mathbb{G}_2, \mathbb{G}_T, g_1, g_2, e)$  in the public parameters  $pp$ . The commitment is instantiated with Pedersen commitments [37] in  $\mathbb{G}_1$ , and the pseudonym function is executed in  $\mathbb{G}_1$  as well, with hashing to the curve [19] and the “exponent” PRF realized with HMAC-SHA256 mapping elements to  $\mathbb{Z}_q$  [19, §5.3].

For comparison with OIDC, PPOIDC, and UPPRESSO, we used the same RSA-signature for IdP-signed tokens in all schemes. For

Entity Protocol \ Ops.	RP		User		IdP
	<i>Vf</i>	<i>Req</i>	<i>Fin</i>	<i>Init</i>	<i>Res</i>
OIDC [41]	0.06	n/a	n/a	n/a	1.41
AIF-ZKP [31]	1.79	8.84	1.31	1.05	16.23
PPOIDC [26]	0.11	n/a	n/a	2647.89	6.98
UPPRESSO [25]	1.37	0.85	n/a	0.70	2.08
Our Work: $\pi_{\text{OPPID}}$	2.13	9.42	2.03	1.68	17.33

**Table 2: The execution mean in milliseconds.**

OIDC, we followed the specification [41, §8.1] and used SHA256 for the *ppid* computation. For PPOIDC, no implementation was available, so we implemented their scheme from scratch. We chose SHA256 for  $H$  to instantiate  $F(\text{uid}, \text{rid})$  and the commitment scheme, as originally proposed in [26]. For the ZKP  $\pi$  of a pre-image of  $H$ , we chose ZoKrates [16] – the zkSNARK is created by the user and verified by the IdP. For UPPRESSO, we implemented its core cryptographic operations: IdP-issued RP certificates are realized via an RSA-SHA256 signature on *rid* only, and pseudonyms are computed in  $\mathbb{G}_1$ .

*Evaluation Results.* The benchmark results of all schemes are summarized in Table 2, with all operations performed on an Apple M1 CPU (8-core, 2020, 3.2 GHz). Our implementation and benchmarks are available at [3] for reproducibility.

Our protocol  $\pi_{\text{OPPID}}$ , achieving all desired security and privacy properties, is highly efficient. User operations and token verification each take only 2ms, proof verification at the IdP requires only 17ms after a 9ms generation time at the RP. We now discuss these results in relation to the closest variants, PPOIDC and UPPRESSO.

Our scheme is significantly faster than PPOIDC, which requires an expensive proof generation of 2.5s by the user, which came for a fast IdP verification of only 7ms. One might be able to speed up the proof generation of a hash preimage by using alternative zkSNARK setups, but that would increase verification or communication costs [39]. UPPRESSO has almost an identical pseudonym computation as in  $\pi_{\text{OPPID}}$  and does not involve the costs for RP authentication, thus it is slightly faster than ours – but also provides less security.

Regarding communication costs, an element in  $(\mathbb{Z}_q, \mathbb{G}_1, \mathbb{G}_2, \mathbb{G}_T)$  requires (32, 48, 92, 576) bytes respectively. Therefore,  $\pi_{\text{OPPID}}$  requires only 864 bytes  $(3\mathbb{Z}_q + 4\mathbb{G}_1 + 1\mathbb{G}_T)$  for the authentication proof and the blinded and committed *rid* sent to the IdP. The user-generated ZKP in PPOIDC is relatively small at 997 bytes, but the associated costs are substantial due to the large pre-compiled circuit (96435 constraints, 231 MB) and a 37 MB proving key. Fetching these files from the IdP for each login is clearly impractical. The overhead introduced by user-side RP certificates in PPOIDC/UPPRESSO is relatively minor at 32 bytes for the signature in addition to the signed 4-byte *rid*. Thus, in scenarios where RP authentication would *not* be needed, our protocol remains more efficient than PPOIDC.

## ACKNOWLEDGMENTS

This research was partially funded by the HPI Research School on Data Science and Engineering. It was also supported by the German Federal Ministry of Education and Research (BMBF) through funding of the ATLAS project under reference number 16KISA037.

## REFERENCES

- [1] Google 2023. *OpenID Connect*. Google. <https://developers.google.com/identity/protocols/oauth2/openid-connect>
- [2] Apple Inc. 2023. *Sign in with Apple*. Apple Inc. [https://developer.apple.com/documentation/sign\\_in\\_with\\_apple](https://developer.apple.com/documentation/sign_in_with_apple)
- [3] 2024. *OPPID*. <https://github.com/jmakr0/OPPID>
- [4] Gergely Alpár, Fabian van den Broek, Brinda Hampiholi, Bart Jacobs, Wouter Lueks, and Sietse Ringers. 2017. IRMA: practical, decentralized and privacy-friendly identity management using smartphones. *HotPETS 2017* (2017).
- [5] Apple. 2023. *Sign in with Apple & Privacy*. <https://www.apple.com/legal/privacy/data/en/sign-in-with-apple/>
- [6] Paulo S. L. M. Barreto, Ben Lynn, and Michael Scott. 2003. Constructing Elliptic Curves with Prescribed Embedding Degrees. In *SCN 02 (LNCS, Vol. 2576)*, Stelvio Cimato, Clemente Galdi, and Giuseppe Persiano (Eds.). Springer, Heidelberg, Germany, Amalfi, Italy, 257–267. [https://doi.org/10.1007/3-540-36413-7\\_19](https://doi.org/10.1007/3-540-36413-7_19)
- [7] Manuel Blum, Paul Feldman, and Silvio Micali. 1988. Non-Interactive Zero-Knowledge and Its Applications (Extended Abstract). In *20th ACM STOC*. ACM Press, Chicago, IL, USA, 103–112. <https://doi.org/10.1145/62212.62222>
- [8] Jan Camenisch, Manu Drijvers, and Anja Lehmann. 2016. Anonymous Attestation Using the Strong Diffie Hellman Assumption Revisited. In *Trust and Trustworthy Computing*, Michael Franz and Panos Papadimitratos (Eds.). Springer International Publishing, Cham, 1–20.
- [9] Jan Camenisch, Aggelos Kiayias, and Moti Yung. 2009. On the Portability of Generalized Schnorr Proofs. In *EUROCRYPT 2009 (LNCS, Vol. 5479)*, Antoine Joux (Ed.). Springer, Heidelberg, Germany, Cologne, Germany, 425–442. [https://doi.org/10.1007/978-3-642-01001-9\\_25](https://doi.org/10.1007/978-3-642-01001-9_25)
- [10] Silvia Casacuberta, Julia Hesse, and Anja Lehmann. 2022. SoK: Oblivious Pseudorandom Functions. Cryptology ePrint Archive, Report 2022/302. <https://eprint.iacr.org/2022/302>.
- [11] James Conners, Corey Devenport, Stephen Derbidge, Natalie Farnsworth, Kyler Gates, Stephen Lambert, Christopher McClain, Parker Nichols, and Daniel Zapala. 2022. Let’s authenticate: Automated certificates for user authentication. In *Network and Distributed Systems Security Symposium (NDSS)*.
- [12] Poulami Das, Julia Hesse, and Anja Lehmann. 2022. DPaSE: Distributed Password-Authenticated Symmetric-Key Encryption, or How to Get Many Keys from One Password. In *ASIACCS 22*, Yuji Suga, Kouichi Sakurai, Xuhua Ding, and Kazuo Sako (Eds.). ACM Press, Nagasaki, Japan, 682–696. <https://doi.org/10.1145/3488932.3517389>
- [13] Arkajit Dey and Stephen Weis. 2010. PseudoID: Enhancing Privacy in Federated Login. In *Hot Topics in Privacy Enhancing Technologies*. Sciencdo, Berlin, Germany, 95–107.
- [14] Jesus Diaz and Anja Lehmann. 2021. Group Signatures with User-Controlled and Sequential Linkability. In *PKC 2021, Part I (LNCS, Vol. 12710)*, Juan Garay (Ed.). Springer, Heidelberg, Germany, Virtual Event, 360–388. [https://doi.org/10.1007/978-3-030-75245-3\\_14](https://doi.org/10.1007/978-3-030-75245-3_14)
- [15] Yevgeniy Dodis and Aleksandr Yampolskiy. 2005. A Verifiable Random Function with Short Proofs and Keys. In *PKC 2005 (LNCS, Vol. 3386)*, Serge Vaudenay (Ed.). Springer, Heidelberg, Germany, Les Diablerets, Switzerland, 416–431. [https://doi.org/10.1007/978-3-540-30580-4\\_28](https://doi.org/10.1007/978-3-540-30580-4_28)
- [16] Jacob Eberhardt and Stefan Tai. 2018. ZoKrates - Scalable Privacy-Preserving Off-Chain Computations. In *2018 IEEE International Conference on Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom) and IEEE Smart Data (SmartData)*. 1084–1091. [https://doi.org/10.1109/Cybermatics\\_2018.2018.00199](https://doi.org/10.1109/Cybermatics_2018.2018.00199)
- [17] EU. 2024. *Regulation 2024/1183 of the european parliament and of the council of 11 april 2024 amending regulation no 910/2014 as regards establishing the european digital identity framework*. <https://eur-lex.europa.eu/eli/reg/2024/1183/oj>
- [18] Facebook. 2021. *OpenID Connect*. <https://developers.facebook.com/docs/facebook-login>
- [19] A. Faz-Hernandez, S. Scott, N. Sullivan, R. S. Wahby, and C. A. Wood. 2023. *Hashing to Elliptic Curves*. RFC 9380. RFC Editor.
- [20] Amos Fiat and Adi Shamir. 1987. How to Prove Yourself: Practical Solutions to Identification and Signature Problems. In *CRYPTO’86 (LNCS, Vol. 263)*, Andrew M. Odlyzko (Ed.). Springer, Heidelberg, Germany, Santa Barbara, CA, USA, 186–194. [https://doi.org/10.1007/3-540-47721-7\\_12](https://doi.org/10.1007/3-540-47721-7_12)
- [21] Ge Gao, Yuan Zhang, Yaqing Song, and Shiyu Li. 2024. PrivSSO: Practical Single-sign-on Authentication against Subscription/Access Pattern Leakage. *IEEE Transactions on Information Forensics and Security* (2024).
- [22] Shafi Goldwasser, Silvio Micali, and Ronald L. Rivest. 1988. A Digital Signature Scheme Secure Against Adaptive Chosen-message Attacks. *SIAM J. Comput.* 17, 2 (April 1988), 281–308.
- [23] Jens Groth. 2006. Simulation-Sound NIZK Proofs for a Practical Language and Constant Size Group Signatures. In *ASIACRYPT 2006 (LNCS, Vol. 4284)*, Xuejia Lai and Kefei Chen (Eds.). Springer, Heidelberg, Germany, Shanghai, China, 444–459. [https://doi.org/10.1007/11935230\\_29](https://doi.org/10.1007/11935230_29)
- [24] Chengqian Guo, Fan Lang, Qiongqiao Wang, and Jingqiang Lin. 2022. UP-SSO: Enhancing the User Privacy of SSO by Integrating PPID and SGX. In *2021 International Conference on Advanced Computing and Endogenous Security*. IEEE, 01–05.
- [25] Chengqian Guo, Jingqiang Lin, Quanwei Cai, Wei Wang, Fengjun Li, Qiongqiao Wang, Jiwu Jing, and Bin Zhao. 2021. Uppresso: Untraceable and unlinkable privacy-preserving single sign-on services. *arXiv preprint arXiv:2110.10396* (2021).
- [26] Sven Hammann, Ralf Sasse, and David A. Basin. 2020. Privacy-Preserving OpenID Connect. In *ASIACCS 20*, Hung-Min Sun, Shih-Pyng Shieh, Guofei Gu, and Giuseppe Ateniese (Eds.). ACM Press, Taipei, Taiwan, 277–289. <https://doi.org/10.1145/3320269.3384724>
- [27] Jakob Heher, Lena Heimberger, and Stefan More. 2024. BISON: Blind Identification through Stateless scOpe-specific derivation. *arXiv preprint arXiv:2406.01518* (2024).
- [28] Marios Isaakidis, Harry Halpin, and George Danezis. 2016. UnlimitedID: Privacy-preserving federated identity management using algebraic MACs. In *Proceedings of the 2016 ACM on Workshop on Privacy in the Electronic Society*. 139–142.
- [29] Stanislaw Jarecki, Aggelos Kiayias, and Hugo Krawczyk. 2014. Round-Optimal Password-Protected Secret Sharing and T-PAKE in the Password-Only Model. In *ASIACRYPT 2014, Part II (LNCS, Vol. 8874)*, Palash Sarkar and Tetsu Iwata (Eds.). Springer, Heidelberg, Germany, Kaoshiung, Taiwan, R.O.C., 233–253. [https://doi.org/10.1007/978-3-662-45608-8\\_13](https://doi.org/10.1007/978-3-662-45608-8_13)
- [30] Stanislaw Jarecki, Hugo Krawczyk, and Jason Resch. 2018. Threshold Partially-Oblivious PRFs with Applications to Key Management. Cryptology ePrint Archive, Report 2018/733. <https://eprint.iacr.org/2018/733>.
- [31] Maximilian Kroschewski and Anja Lehmann. 2023. Save The Implicit Flow? Enabling Privacy-Preserving RP Authentication in OpenID Connect. *Proceedings on Privacy Enhancing Technologies* 4 (2023), 96–116.
- [32] Anja Lehmann. 2019. ScrambleDB: Oblivious (Chameleon) Pseudonymization-as-a-Service. *PoPETS 2019*, 3 (July 2019), 289–309. <https://doi.org/10.2478/popets-2019-0048>
- [33] Anna Lysyanskaya, Ronald L. Rivest, Amit Sahai, and Stefan Wolf. 1999. Pseudonym Systems. In *SAC 1999 (LNCS, Vol. 1758)*, Howard M. Heys and Carlisle M. Adams (Eds.). Springer, Heidelberg, Germany, Kingston, Ontario, Canada, 184–199. [https://doi.org/10.1007/3-540-46513-8\\_14](https://doi.org/10.1007/3-540-46513-8_14)
- [34] Moni Naor, Benny Pinkas, and Omer Reingold. 1999. Distributed Pseudo-random Functions and KDCs. In *EUROCRYPT’99 (LNCS, Vol. 1592)*, Jacques Stern (Ed.). Springer, Heidelberg, Germany, Prague, Czech Republic, 327–346. [https://doi.org/10.1007/3-540-48910-X\\_23](https://doi.org/10.1007/3-540-48910-X_23)
- [35] NIST. 2023. *SP 800-63 Digital Identity Guidelines*. <https://pages.nist.gov/800-63-4/sp800-63c.html>
- [36] Christian Paquin and Greg Zaverucha. 2013. U-prove cryptographic specification v1. 1. *Technical Report, Microsoft Corporation* (2013).
- [37] Torben P. Pedersen. 1992. Non-Interactive and Information-Theoretic Secure Verifiable Secret Sharing. In *CRYPTO’91 (LNCS, Vol. 576)*, Joan Feigenbaum (Ed.). Springer, Heidelberg, Germany, Santa Barbara, CA, USA, 129–140. [https://doi.org/10.1007/3-540-46766-1\\_9](https://doi.org/10.1007/3-540-46766-1_9)
- [38] David Pointcheval and Olivier Sanders. 2016. Short Randomizable Signatures. In *CT-RSA 2016 (LNCS, Vol. 9610)*, Kazuo Sako (Ed.). Springer, Heidelberg, Germany, San Francisco, CA, USA, 111–126. [https://doi.org/10.1007/978-3-319-29485-8\\_7](https://doi.org/10.1007/978-3-319-29485-8_7)
- [39] Ethereum Research. 2023. *Benchmarking ZKP Development Frameworks: the Pantheon of ZKP*. <https://ethresear.ch/t/benchmarking-zkp-development-frameworks-the-pantheon-of-zkp/14943>
- [40] Ahmad Sabouri and Kai Rannenberg. 2015. ABC4Trust: protecting privacy in identity management by bringing privacy-ABCs into real-life. In *Privacy and Identity Management for the Future Internet in the Age of Globalisation: 9th IFIP WG 9.2, 9.5, 9.6/11.7, 11.4, 11.6/SIG 9.2. 2 International Summer School, Patras, Greece, September 7-12, 2014, Revised Selected Papers* 9. Springer, 3–16.
- [41] N. Sakimura, J. Bradley, M. Jones, B. de Medeiros, and C. Mortimore. 2014. *OpenID Connect Core 1.0*. [https://openid.net/specs/openid-connect-core-1\\_0.html](https://openid.net/specs/openid-connect-core-1_0.html)
- [42] Rongwu Xu, Sen Yang, Fan Zhang, and Zhixuan Fang. 2023. MISO: Legacy-compatible Privacy-preserving Single Sign-on using Trusted Execution Environments. In *2023 IEEE 8th European Symposium on Security and Privacy (EuroS&P)*. IEEE, 352–372.
- [43] Zhiyi Zhang, Michal Król, Alberto Sonnino, Lixia Zhang, and Etienne Riviere. 2021. EL PASSO: Efficient and Lightweight Privacy-preserving Single Sign On. *PoPETS 2021*, 2 (April 2021), 70–87. <https://doi.org/10.2478/popets-2021-0018>

## A OMITTED MODEL PARTS

Here, we provide the omitted parts and further explanations of our security model as well as compare it to the Session Binding model from [31].

## A.1 Correctness

Recall that we denote with  $\mathcal{R}$ ,  $\mathcal{S}$ , and  $\mathcal{U}$  the RP, session, and user spaces. An OPPID scheme – as defined in Sec. 2.3 – is correct if for all  $\lambda \in \mathbb{N}$ , setup and RP registrations with  $rid \in \mathcal{R}$

$$\begin{aligned} pp &\leftarrow \text{Setup}(1^\lambda), ((isk, \mathcal{M}), ipk) \leftarrow \text{KGen}_{\text{IdP}}(pp), \\ (cred, \mathcal{M}') &\leftarrow \langle \text{Join}_{\text{RP}}(ipk, rid), \text{Reg}_{\text{IdP}}(isk, rid, \mathcal{M}) \rangle, \end{aligned}$$

all authentication sessions with  $sid \in \mathcal{S}$  of user  $uid \in \mathcal{U}$  to RP  $rid$

$$\begin{aligned} (orid, crid) &\leftarrow \text{Alnit}_{\text{U}}(ipk, rid) \\ auth &\leftarrow \text{AReq}_{\text{RP}}(ipk, rid, cred, crid, orid, sid) \\ \tau &\leftarrow \text{ARes}_{\text{IdP}}(isk, auth, crid, uid, ctx, sid) \\ (\tau_{\text{fin}}, ppid) &\leftarrow \text{AFin}_{\text{U}}(ipk, rid, crid, orid, ctx, sid, \tau), \end{aligned}$$

result in

$$\text{Vf}_{\text{RP}}(ipk, (rid, ppid, ctx, sid), \tau_{\text{fin}}) = 1.$$

## A.2 Session Binding: RP Binding & RP Auth.

Here, we provide more intuition on how the properties of RP Binding and RP Authentication are captured in our Session Binding game  $\text{Exp}_{\mathcal{A}, \text{OPPID}}^{\text{SES-BIN}}$  in Figure 4. Note that RP Binding II and RP Authentication II define security aspects that were not covered by the original model in [31]. We give a more detailed comparison in the following section.

*RP Binding.* If an honest user  $uid$  wants to authenticate to a certain RP  $rid$  in a session  $sid$ , the adversary wins if it can produce a token for the same session, but which is valid for a different  $rid^*$ . This models phishing attacks from malicious RPs. To capture this property, we need to know what the intent of the honest user was, which is done through  $\mathcal{O}.\text{Alnit}_{\text{U}}$  and  $\mathcal{O}.\text{AResFin}$ . Using the bookkeeping in these oracles, we can define an honest user’s intent as  $(rid, uid, ctx, sid)$  in  $\mathcal{Q}_{\tau_{\text{fin}}}$  when the finalized token is computed. The adversary wins if it can come up with a token for that session but with  $rid^* \neq rid$ , as defined in condition (b) for *RP Binding I*.

For a session initiated by a corrupt user towards a corrupt RP, we never know the exact RP the user wants to authenticate to:  $\mathcal{A}$  invokes  $\mathcal{O}.\text{ARes}_{\text{IdP}}$  on adversarial chosen inputs  $auth$  and  $crid$ , which both hide the  $rid$ . So our guarantees are naturally weaker here than for honest users. What we do guarantee, is that the adversary cannot re-use the same token it receives from the IdP, and claim it to be valid towards multiple corrupt RPs. This is the reason we have  $\mathcal{O}.\text{Vf}_{\text{RP}}$ . It might be surprising that we provide an oracle that runs purely on public values, and would mimic something that is typically run internally by the adversary. What we want to model here is that if the adversary presents a token for  $(rid, ppid, ctx, sid)$  somehow *publicly*, e.g., towards a judge or some external *honest* entity, and later presents a token for the same context  $ctx, sid$  but for a different  $rid^* \neq rid$  (which is its forgery). Thus, despite the honest IdP never learning the exact RP the corrupt user wanted to authenticate, we guarantee that this user can authenticate to at most *one* RP. This *RP Binding II* is captured in winning condition (c), and can be seen as a notion similar to the one-more unforgeability property in blind signatures.

*RP Authentication.* The previous two properties ensured that the IdP-generated token is bound to the blindly received  $rid$ . We

further want to guarantee that only legitimate RPs can request such tokens, which is captured in conditions (d) and (e). In condition (d), the adversary wins if it manages to produce a valid token for some  $rid^*$ , yet this RP has never registered, i.e.,  $rid^* \notin \text{HRP} \cup \text{CRP}$  where HRP and CRP are the sets of all honest and corrupt RPs, the IdP has registered through  $\mathcal{O}.\text{RegHRP}$  and  $\mathcal{O}.\text{RegCRP}$ .

If the  $rid^*$  in the forgery belongs to an honest RP, we want even stronger guarantees and ensure that it must be infeasible to win if the honest RP had never authenticated for that particular session  $sid^*$ . As all authentication requests for honest RPs are handled through  $\mathcal{O}.\text{AReq}_{\text{RP}}$ , where we store each query  $(rid, sid)$  in  $\mathcal{Q}_{\text{auth}}$ , this translates to simply checking that  $(rid^*, sid^*) \notin \mathcal{Q}_{\text{auth}}$ .

## A.3 Comparison to Session Binding from [31]

Before we compare the security guarantees of our Session Binding model to the one from [31], we outline two fundamental *functional* differences between both models.

*Functional Differences.* The focus of our work is *pseudonymous* user authentication towards RPs, where the pseudonyms are blindly computed by the IdP. This was the main challenge and manifests as Nym Correctness in our model. In the work of Kroschewski and Lehmann [31], all authentication tokens contained the globally unique user identifier  $uid$  that was directly vouched for by the IdP.

The focus of [31] was on RP authentication, and their work explicitly models epoch-based credentials to allow adaptive RP corruption and enable their revocation. Our model omits this approach for simplicity. We stress that epoch-based credentials are straightforward to add at the construction level but significantly complicate the security model. In fact, epoch-based renewal is an orthogonal aspect to our focus on pseudonymous identifiers, and we can consider our setting as a single-epoch – and pseudonym-extended – version of [31].

*Original Model Only Considers Honest Users + Corrupt RPs.* For better comparison, we state the original Session Binding property translated to our single-epoch view and syntax on the left in Fig. 6. This makes it easy to see that, ignoring the obvious differences due to the different functional properties, our Session Binding model is significantly stronger than [31], as shown on the right in Fig. 6.

The original definition has the weakness that it only considers the security of sessions between honest users and corrupt RPs. This was justified by the argument that (i) honest users are necessary in order to know the  $rid$  a user wanted to authenticate to in a particular session and compare it to the adversary’s forgery, and (ii) honest RPs do not give the adversary an advantage. We will now explain why this restriction to honest users and corrupt RPs led to a security model that does not capture all desirable properties, and how we incorporated them into our definition.

*Adding Security for Sessions of Corrupt User.* Regarding (i): While we indeed do not know the intended  $rid$  when an IdP issues a token towards a corrupt user and corrupt RP, the model should not abandon security in such scenarios. What we still care about — and in fact, this could be seen as the most crucial security property of SSO — is that the IdP’s signature cannot be used out of context. Simply removing the restriction of honest users in the winning condition already enhances security. This is what our game achieves



$\text{Exp}_{\mathcal{A}, \text{OPP}\text{ID}}^{\text{RP Session Binding}}(\lambda)$	This Work	Model [31]
$pp \leftarrow \text{Setup}(1^\lambda); ((isk, M), ipk) \leftarrow \text{KGen}_{\text{IdP}}(pp)$ $\mathcal{O} := \{\text{RegCRP}, \text{Alnit}_U, \text{AReqRP}, \text{ARes}_{\text{IdP}}, \text{AResFin}\}$ $(rid^*, ppid^*, ctx^*, sid^*, \tau_{\text{fin}}^*) \leftarrow \mathcal{A}^{\mathcal{O}}(ipk)$ Return 1 if $\forall \text{RP}(ipk, (rid^*, uid^*, ctx^*, sid^*), \tau_{\text{fin}}^*) = 1 \wedge "uid^*" \text{ is honest}$ and at least one of the following holds: (a) $(rid^*, uid^*, ctx^*, sid^*) \notin \mathcal{Q}_{\tau_{\text{fin}}}$ (b) $(rid^*, uid^*, ctx^*, sid^*) \in \mathcal{Q}_{\tau_{\text{fin}}} \wedge rid^* \notin \text{CRP}$	(1) Direct Forgery (2) Indirect Forgery: (a) Nym Correctness (b) RP Binding I (c) RP Binding II (d) RP Authentication I (e) RP Authentication II	(a), but for honest users only  (a), via $uid$ correctness, but for honest users only (a) —, missing (b), but for honest users only —, missing, partially covered via RP Accountability

**Figure 6:** Left: RP Session Binding from [31] translated to our syntax, where we abuse notation and consider  $F(uid, rid) = uid$ . The honest user requirement was enforced in [31] by checking that the  $uid^*$  of the forgery was never used in a "malicious" query via  $\mathcal{O}.\text{ARes}_{\text{IdP}}$  but always via  $\mathcal{O}.\text{AResFin}$  – which allows the challenger to know the intended  $rid$  for every session. Right: Comparison between the Session Binding guarantees in this work and [31], mapping the winning conditions to each other.

compared to [31] for Direct Forgeries (1) and Indirect Forgeries (2a, d). However, this alone is not sufficient: the guarantee does not extend to the blindly signed  $rid$ , as condition (2b) is the only sub-case that strictly needs to be limited to the honest user setting.

This is where our game introduces the winning condition (2c) in a one-more unforgeability style. This condition requires that it must be infeasible to produce multiple valid tokens for different corrupt  $rid$ s. Generating multiple tokens for corrupt users and RPs would not be a direct attack on the authentication session, but it would still be undesirable behavior: it could allow corrupt RPs to present apparently IdP-attested tokens for non-existent sessions, falsely inflating their active user base. Our stronger model prevents this.

*Adding Stronger Guarantees for Honest RPs.* The second enhancement is related to (ii). The original model allows the adversary to win by producing a token for an  $rid$  that was never registered at all. Again, this was only defined for honest users, whereas our model extends this to corrupt users as well. The more significant change is that our model provides stronger guarantees when an *honest* RP is involved—which was not addressed in [31]. Our winning condition (2e) additionally requires that it must be infeasible to produce a token for some  $rid, sid$  when  $rid$  belongs to an honest RP that never authenticated for session  $sid$ .

At first glance, this might appear to be covered by the RP Accountability game in [31]. However, the RP Accountability definition only applies to fully blind authentication requests – and therefore is weaker. In our Session Binding game, we have knowledge of the finalized token including  $rid$ , and can verify whether a corresponding authentication request from this  $rid$  and session  $sid$  exists.

#### A.4 Request Authentication

This property captures the authenticity of the *request* that an IdP receives and uses to produce its token. It ensures that if an IdP receives an authenticated request  $(auth, crid, sid)$  via a user  $uid$  for which  $\text{ARes}_{\text{IdP}}(isk, auth, crid, uid, ctx, sid)$  produces an output  $\neq \perp$ , it must originate from a previously registered RP.

The game is defined by  $\text{Exp}_{\mathcal{A}, \text{OPP}\text{ID}}^{\text{REQ-AUTH}}$  (see Figure 4) and follows a classic unforgeability definition. Instead of demanding security for the final authentication token (as in Session Binding), it captures unforgeability for the RP authenticated information:  $crid, sid$ . The IdP is honest here, and the adversary can register honest RPs

through  $\mathcal{O}.\text{RegHRP}$ , obtain their authenticated requests through  $\mathcal{O}.\text{AReqRP}$ , and query tokens from the IdP through  $\mathcal{O}.\text{ARes}_{\text{IdP}}$ . The adversary wins if it outputs  $(auth^*, crid^*, uid^*, ctx^*, sid^*)$ , where the IdP accepts the authentication (i.e., does not return  $\perp$ ) and no honest RP authenticated  $crid^*$  in session  $sid^*$ .

Note that this property defines the security of the fully blind authentication request towards the IdP, which implies that all RPs must be honest — otherwise, "forging" is trivial. The original definition in [31] allowed corrupt RPs, which was possible as their work considered epoch-based renewal of membership credentials. The security model then only requires that all registered RPs of the *current* epoch must be honest, but RPs from previous epochs can be corrupt. In this sense, just as in our Session Binding definition, our Request Authentication definition can be seen as a single-epoch version of the Request Authentication definition of [31]. We do not allow any corrupt RPs in the definition as any corrupt RP in the single-epoch would allow the adversary to win trivially.

*Definition A.1 (Request Authentication).* An OPPID scheme satisfies Request Authentication if for all PPT adversaries  $\mathcal{A}$ , it holds

$$\Pr[\text{Exp}_{\mathcal{A}, \text{OPP}\text{ID}}^{\text{REQ-AUTH}}(\lambda) = 1] \leq \text{negl}(\lambda).$$

## B FULL PROOFS OF $\pi_{\text{OPP}\text{ID}}$

Here, we provide the proofs showing that our protocol described in Sec. 5 achieves all properties defined in Sec. 3. We start by analyzing the properties of the pseudonym function  $F_{\text{DL+PRF}}$ .

*Properties of  $F_{\text{DL+PRF}}$ .* We analyze the properties of our pseudonym function

$$F_{\text{DL+PRF}} = F_k(uid, rid) := H(rid)^{\text{PRF}(k, uid)}.$$

Recall that our OPPID model requires this function to provide unique, collision-free, and unlinkable pseudonyms (see Sec. 2.2).

Uniqueness requires that for every  $rid, uid$  combination, there is a unique pseudonym  $ppid = F(uid, rid)$ .  $F_{\text{DL+PRF}}$  naturally ensures uniqueness as it is a deterministic function.

Regarding collision-freeness, note that  $H(rid)^{uk}$  is a permutation for  $uk \in \mathbb{Z}_q$ . If  $uk \leftarrow \text{PRF}(k, uid)$  is an injective function, then  $F$  provides different pseudonyms for each user. That is, for all  $rid$  and  $uid \neq uid'$  it holds that  $F_k(uid, rid) \neq F_k(uid', rid)$ . This property requires that  $|\mathcal{U}| \leq \mathbb{Z}_q$ , which holds for any normal deployment, where the number of users is significantly smaller than  $\mathbb{Z}_q$ .

LEMMA B.1. *Our pseudonym function  $F_{DL+PRF}$  provides unique and collision-free pseudonyms, if  $H$  is deterministic, PRF is deterministic, injective, and  $\mathcal{U} \leq \mathbb{Z}_q$ .*

In the following, we prove that  $F_{DL+PRF}$  is a secure pseudorandom function. This immediately grants the unlinkability of the pseudonyms, which will be helpful in proving the Unlinkability property of  $\pi_{OPP-ID}$ . Note that a similar function (essentially  $F_{DL+PRF}$  with double hashing) was shown to be a secure (partially oblivious) pseudorandom function already [30], whereas we need the classic PRF property here and do not apply the outer hash.

LEMMA B.2. *Our pseudonym function  $F_{DL+PRF}$  is a secure PRF if  $H$  is a random oracle, PRF is a secure pseudorandom function, and DDH holds in  $\mathbb{G}$ .*

PROOF. Our proof has two main steps. The first step switches from PRF( $k, uid$ ) to the random values from  $uk \leftarrow_{\mathbb{R}} \mathbb{Z}_q$  while evaluating the PRF output in the oracle. This change is indistinguishable by pseudorandomness of PRF. The second step relies on the simple observation that  $H(rid)^{uk}$  is HashDH PRF [34] and shows the pseudorandomness of these values by relying on the multi-key pseudorandomness of the HashDH PRF.

*Game 0.* This game is identical to the pseudorandomness game with  $F_{DL+PRF}$ .

*Game 1.* In this game, for each PRF query ( $uid_i, rid$ ) with a new  $uid_i$ , we sample  $uk_i \leftarrow_{\mathbb{R}} \mathbb{Z}_q$  for the  $uid_i$  and answer  $F_{DL+PRF}(uid_i, rid)$  queries as  $F_{DL+PRF}(uid_i, rid) = H(rid)^{uk_i}$ . By pseudorandomness of PRF, this change is indistinguishable to the adversary.

*Game 2.* This game finalizes the proof by changing the PRF evaluations for  $F_{DL+PRF}(uid_i, rid)$  from  $H(rid)^{uk_i}$  to  $y \leftarrow_{\mathbb{R}} \mathbb{G}$ .

We show that this change is indistinguishable to the adversary by presenting sequences of indistinguishable hybrids between Games 1 and 2 where the first and last hybrids are identical to Games 1 and 2, respectively. Each hybrid changes the PRF evaluations for  $uid_i$  from  $H(rid)^{uk_i}$  to a random value. Let  $uid_1, \dots, uid_n$  represent the  $uid$  values that the adversary queries  $F_{DL+PRF}$  oracle.

*Hybrid 0.* This hybrid is identical to Game 1.

*Hybrid  $i \in \{1, \dots, n\}$ .* Hybrid $_i$  runs Hybrid $_{i-1}$  identically except for the following change. For  $F_{DL+PRF}(uid_i, rid)$  queries, Hybrid $_i$  outputs  $F_{DL+PRF}(uid_i, rid) \leftarrow_{\mathbb{R}} \mathbb{G}$  random values. Note that Hybrid  $n$  is identical to Game 2.

*Transition Hybrid  $i \rightarrow$  Hybrid  $i+1$ .* The transition between hybrids simply relies on the pseudorandomness of HashDH. In particular, Hybrid $_i$  evaluates HashDH PRF  $H(rid)^{uk_{i+1}}$  with the PRF key  $uk_{i+1}$  whereas Hybrid $_{i+1}$  changes these evaluations to the random values from  $\mathbb{G}$ . By pseudorandomness of HashDH, Hybrids  $i$  and  $i+1$  are indistinguishable and HashDH is pseudorandom if DDH assumption holds in ROM [34].  $\square$

## B.1 Proof of Theorem 5.1 (Unlinkability)

Here, we prove that  $\pi_{OPP-ID}$  satisfies Unlinkability (see Def. 3.1) if  $H$  is a random oracle, PRF is a secure pseudorandom function, and the DDH assumption holds in  $\mathbb{G}$ .

PROOF. Recall that  $\mathcal{A}$  is given the IdP's public key and oracles  $\mathcal{O} := \{\text{RegCRP}, \text{Alnit}_{\mathbb{U}}, \text{ARes}_{\text{IdP}}\}$  to register corrupt RPs, initialize honest user sessions, and obtain pseudonyms and tokens from the IdP. Eventually, it outputs  $uid_0, uid_1$ , and  $rid$ , together with  $auth, crid, ctx, sid$ , and receives a challenge token and pseudonym  $(\tau_{fin_b}, ppid_d)$  for  $uid_b$ .

We further know that if  $\mathcal{A}$  wins, it must not make any query that trivially reveals  $F_k(uid_0, rid)$  or  $F_k(uid_1, rid)$ . That is,  $\mathcal{A}$  is not allowed to make  $\mathcal{O}.\text{ARes}_{\text{IdP}}$  queries  $(auth, crid, uid_d, ctx, sid)$  for  $d \in \{0, 1\}$  where either:

- $crid$  belongs to  $rid$  from the challenge query (via previous query to  $\mathcal{O}.\text{Alnit}_{\mathbb{U}}$ ), or
- $crid$  is malicious, i.e., not an output from  $\mathcal{O}.\text{Alnit}_{\mathbb{U}}$ .

These conditions ensure that all queries  $\mathcal{A}$  makes towards  $\mathcal{O}.\text{ARes}_{\text{IdP}}$  must involve  $crid$  values that are honestly generated and for which we know the underlying  $rid$  (and blinding value  $r$ ).

This makes the proof straightforward. As we already know the blinding value of  $\bar{y}$ , we can compute the response value by relying on PRF computation instead of an oblivious PRF evaluation. We prove Unlinkability through a small sequence of games, replacing all user-dependent values for  $uid_0$  and  $uid_1$  with random values, ensuring that  $\mathcal{A}$  has no better chance than guessing to determine the bit  $b$ . Let *Game 0* be the original game.

*Game 1.* We now change the way we compute the  $\bar{y}_d$  values in  $\mathcal{O}.\text{ARes}_{\text{IdP}}$  to compute the challenge token values  $\tau_d$ . We simulate  $\mathcal{A}$ 's view by using  $F_{DL+PRF}$  as a black-box algorithm to reason about its pseudorandomness, showing the Unlinkability property of our scheme instead of its one-more pseudorandomness. Since  $\mathcal{O}.\text{ARes}_{\text{IdP}}$  queries must contain  $crid$  values generated through  $\text{Alnit}_{\mathbb{U}}$ , we can look up  $(rid, crid, orid)$  in  $Q_{rid}$  and parse  $orid = (o, r)$ . Instead of blindly computing  $\bar{y}$ , we compute  $y = F_{DL+PRF}(k, (rid, uid))$  from the known input and blind the response later:  $\bar{y}_d = y_d^r$ . This change produces outputs identical to those in the previous game, so this game hop cannot be distinguished by  $\mathcal{A}$ .

*Game 2.* In this game, we change the way we compute  $y_d = F_{DL+PRF}(k, (uid_d, rid))$  while computing  $\tau_d$  values and set  $y_d \leftarrow_{\mathbb{R}} \mathbb{G}$  as a random value. By the winning condition of the Unlinkability game, we know that there are no previous  $\mathcal{O}.\text{ARes}_{\text{IdP}}$  queries for  $(uid_d, rid)$ . Also, by the previously proven Lemma B.2, we know that  $F_{DL+PRF}$  is a secure PRF, so the change in this game is indistinguishable.

In the last game, all the bit  $d$ -related values are uniformly random, independent of  $uid_0$  and  $uid_1$ . Thus,  $\mathcal{A}$ 's chance of determining the bit  $b$  is  $1/2$ .  $\square$

*No One-More-Type Assumption.* What might be surprising at first glance is that we do not need a one-more-type assumption typically required in blind evaluation protocols, such as the (2)HashDH-OPRF, which seems equivalent to our construction. This is not surprising after closer examination. First, recall that our security property guarantees that the  $uid$  in the IdP's response remains hidden, and the  $uid$  is not blind towards the IdP but revealed in every query. Second, the guarantee can only hold for *honest* users, which is enforced throughout the game for both challenge users,

$uid_0$  and  $uid_1$ . The adversary can only obtain "blind" PRF evaluations (as part of  $\mathcal{O}.\text{ARes}_{\text{IdP}}$  queries), where the blinded input  $\bar{x}$  was honestly generated. Thus, in the security proof, the challenger is always aware of the blinded  $rid$  behind  $\bar{x}$ , i.e., knows exactly on which values the PRF  $\text{F}_{\text{DL}+\text{PRF}}$  is evaluated. This allows us to prove Unlinkability under the standard assumption that  $\text{F}_{\text{DL}+\text{PRF}}$  is a secure pseudorandom function, which holds if  $\text{H}$  is a random oracle, and the DDH assumption holds in  $\mathbb{G}$  [30].

## B.2 Proof of Theorem 5.2 (Unobservability)

We now provide a simple proof that  $\pi_{\text{OPPID}}$  satisfies Unobservability (see Def. 3.2) if COM is hiding, and the NIZK is zero-knowledge.

**PROOF.**  $\mathcal{A}$  receives  $auth = \pi$  and  $crid = (\bar{x}, com)$ , where  $com$  is a commitment to  $rid$  and  $\bar{x}$  is the blinded hash.  $\pi$  proves knowledge of an IdP-issued credential on  $rid$  and that the commitment  $com$  opens to the same  $rid$ . Unobservability follows directly from the zero-knowledge property of  $\pi$ , the perfect hiding of  $rid$  via  $\bar{x}$  with  $r$ , and the hiding property of COM with the undisclosed opening  $o$ .  $\square$

## B.3 Proof of Theorem 5.3 (Session Binding)

Here we prove that  $\pi_{\text{OPPID}}$  is Session Binding (see Def. 3.3) if the  $S_1$  and  $S_2$  schemes are both EUF-CMA secure, COM is binding, and the NIZK is zero-knowledge and simulation-extractable.

**PROOF.** We split the proof along the winning condition that the adversary must satisfy. Recall that  $\mathcal{A}$  outputs  $(rid^*, ppid^*, ctx^*, sid^*, \tau_{\text{fin}}^*)$  and wins if this is a valid yet non-trivial forgery.

The final token  $\tau_{\text{fin}}^*$  contains  $(com^*, o^*, r^*, \bar{y}^*, \sigma_\tau^*)$ . A valid token implies that  $\sigma_\tau^*$  is a valid signature, meaning

$$S_1.\text{Vf}(pk_1, (com^* || \bar{x}^* || \bar{y}^* || ctx^* || sid^*), \sigma_\tau^*) = 1$$

with  $\bar{x}^* := \text{H}(rid^*)^{r^*}$ ,  $com^* := \text{Com}(rid^*, o^*)$ , and  $ppid^* := \bar{y}^{*-r^*}$ .

We first distinguish whether we have a *direct* or *indirect forgery*, i.e., whether the information that was publicly signed by the IdP is already a forgery or not.

*Case 1: Direct Forgery with  $(\cdot, ctx^*, sid^*) \notin Q_\tau$ .* If the adversary outputs a valid forgery where the combination  $(\cdot, ctx^*, sid^*)$  was never vouched for by the honest IdP, it must have forged the IdP's signature on these values. This is infeasible if the signature scheme  $S_1$  is EUF-CMA secure.

We can build a  $S_1$  forger easily using an adversary that can perform a direct forgery as follows. We get the challenge  $S_1$  public key  $pk_1^*$  and use it as  $pk_1$  while setting the issuer public key  $ipk$ . As we do not know the corresponding secret key to  $pk_1^*$ , we need a way to simulate the  $S_1$  signatures  $\sigma_\tau$  which are part of IdP's output to  $\tau$  queries,  $\tau$ . We simulate  $\sigma_\tau$  values using the signing oracle of  $S_1$ . In particular, the oracles work as follows:

RegHRP, RegCRP, AlnitU, AReqRP, AResFin, VfRP: As they are.

AResIdP: It computes  $\bar{y}$  as before. To form the token on  $m = (com || \bar{x} || \bar{y} || ctx || sid)$ , it queries the signing oracle of  $S_1$  and gets the signature  $\sigma_\tau$  on  $m$  and outputs  $(\sigma_\tau, \bar{y})$ .

Finally, the adversary outputs  $(\sigma^*, m^*) := (\sigma_\tau^*, (com^*, \bar{x}^*, \bar{y}^*, ctx^*, sid^*))$  to the  $S_1$  unforgeability challenger as the  $S_1$  forgery for the forged token  $\tau_{\text{fin}}^* := ((com^*, o^*, r^*, \bar{y}^*, \sigma_\tau^*), y^*)$  and  $\bar{x}^* := \text{H}(rid^*)^{r^*}$ .

As the forged token is valid, we know that  $S_1.\text{Vf}(pk_1^*, \sigma^*, m^*) = 1$ . Furthermore, direct token forgeries ensure that we do not make a signing oracle query for  $(\cdot || \cdot || \cdot || ctx^* || sid^*)$ , so our forgery is on a fresh message.

*Case 2: Indirect Forgery with  $(uid, ctx^*, sid^*) \in Q_\tau$ .* If the honest IdP has previously signed the combination  $(uid, ctx^*, sid^*)$  for a session for user  $uid$ , the adversary can only win if the associated information  $(rid^*, ppid^*)$  that the IdP has blindly signed and derived is inconsistent with the expected pseudonym or behavior of honest and corrupt RPs. This inconsistency is expressed through the five sub-cases in the winning condition of the Session Binding game, and our proof branches accordingly.

What is important here is that the IdP signs additional information in  $\sigma_\tau^*$ , namely  $(com^* || \bar{x}^* || \bar{y}^*)$ . If the adversary outputs a forgery with this tuple differs from what the honest IdP has signed along with  $(ctx^* || sid^*)$ , then we can immediately turn this into a forgery of the  $S_1$  scheme. We simulate  $pk_1$  as the  $S_1$  EUF-CMA challenge public key  $pk_1^*$  and simulate the Session Binding game just as in Case 1.

Note that while the adversary does not explicitly output  $com^*$  and  $\bar{x}^*$ , these values are uniquely defined through its outputs as  $com^* = \text{Com}(rid^*, o^*)$  and  $\bar{x}^* = \text{H}(rid^*)^{r^*}$ , with  $o^*$  and  $r^*$  being part of  $\tau_{\text{fin}}^*$  and  $rid^*$ . At the end of the Session Binding game, we check the  $\text{ARes}_{\text{IdP}}$  query with  $(uid, ctx^*, sid^*) \in Q_\tau$ . If the corresponding query differs from  $(com^* || \bar{x}^* || \bar{y}^* || ctx^* || sid^*)$ , then we can output  $(\sigma_\tau^*, (com^* || \bar{x}^* || \bar{y}^* || ctx^* || sid^*))$  as a valid  $S_1$  forgery.

Thus, the rest of the proof of Case 2 is now conditioned on the fact that the full tuple  $(com^* || \bar{x}^* || \bar{y}^* || ctx^* || sid^*)$  has indeed been signed by the honest IdP in a session with user  $uid$ .

(a) *Nym Correctness:  $ppid^* \neq F(uid, rid^*)$ .* If the adversary wins by satisfying the first sub-condition, it must have produced a valid token with an incorrect pseudonym, i.e., output a  $ppid^* \neq \text{H}(rid^*)^{\text{PRF}(k, uid)}$ .

In our construction, winning under this condition is impossible (other than through manipulating the IdP's signed information, which we excluded above). Recall that the adversary's forgery must contain  $\bar{x}^* || \bar{y}^*$  along with the public session information. We have already excluded the case where  $\mathcal{A}$  manages to manipulate these values. Thus, we know that  $\bar{x}^*, \bar{y}^*$  are the values the IdP has signed in a session  $sid^*$ , where it learned the username  $uid$  and computed

$$\bar{y}^* = \bar{x}^{*\text{PRF}(k, uid)}.$$

As the forgery must pass the verification, we know that

$$ppid^* = \bar{y}^{*-r^*} \quad \text{and} \quad \bar{x}^* = \text{H}(rid^*)^{r^*}$$

Putting it all together implies that

$$ppid^* = ((\text{H}(rid^*)^{r^*})^{\text{PRF}(k, uid)})^{-r^*} = \text{H}(rid^*)^{\text{PRF}(k, uid)}.$$

Thus, for every valid token, it holds that  $ppid^* = F(uid, rid^*)$ .

(b) *RP Binding I:  $(rid, uid, ctx^*, sid^*) \in Q_{\tau_{\text{fin}}} \wedge rid \neq rid^*$ .* If an adversary wins under condition (2b), it must have "high-jacked" an honest user session. For sessions intended by honest users, we know the exact RP  $rid$  they intended to authenticate to, and the adversary wins if it can create a token for this session that

is valid for a different RP  $rid^* \neq rid$ . We can again leverage the fact that we know that  $(com^* || \bar{x}^* || \bar{y}^* || ctx^* || sid^*)$  is the original information signed by the honest IdP in the session with the honest user  $uid$ . We further know that the commitment  $com^*$  is an honestly generated commitment (through  $O.Alnit_U$ ) for  $rid$ .

As the final token contains an opening  $o^*$  and checks that  $com^* = Com(rid^*, o^*)$ , this implies that  $\mathcal{A}$  managed to open the commitment  $com^*$  to a different value  $rid^* \neq rid$ , which is infeasible under the binding property of COM. In particular, by the behavior of  $O.AResFin$ , we know that  $Q_{\tau_{fin}}$  is updated with an  $rid$  only when there is a valid opening  $o$  for the commitment  $o^*$  provided in the token  $\tau_{fin}$ . Thus, the tuple  $(com^*, rid, rid^*, o, o^*)$  breaks the binding property of the commitment.

- (c) *RP Binding II:  $(rid, ppid, ctx^*, sid^*) \in Q_{Vf} \wedge rid \neq rid^*$* . If the IdP did sign  $(com^* || \bar{x}^* || \bar{y}^* || ctx^* || sid^*)$  in a session with a *corrupt* user  $uid$  (which was ensured by the main  $S_1$  unforgeability reduction of condition (2)), we do not know the intended  $rid$  contained in  $com^*$  (or  $\bar{x}^*$ ). Thus, the adversary can "open" the token to any valid  $rid^*$  it wants. In order to create a valid forgery,  $\mathcal{A}$  must have produced (at least) two valid tokens for  $(ctx^*, sid^*)$  yet different  $(rid, ppid)$  and  $(rid^*, ppid^*)$  with  $rid \neq rid^*$ . Note that we do not make any requirements on the pseudonyms here, as the adversary already wins under condition (2a) if it can produce an invalid pseudonym.

The only remaining way for the adversary to provide a condition (2c) forgery is then providing different openings of  $com^*$  to distinct  $rid$  and  $rid^*$  values. If this occurs, the query of the adversary to  $O.Vf_{RP}$  with  $sid^*$  and  $ctx^*$  contains a valid opening of  $com^*$  to  $rid$  and  $o$  for  $rid^* \neq rid$ . Obviously, it contradicts the binding property of COM as the tuple  $(com^*, rid, rid^*, o, o^*)$  breaks the binding property.

- (d) *RP Authentication I:  $rid^* \notin HRP \cup CRP$* . If the adversary wins under this condition, it means that the adversary produced a valid token for an  $rid^*$  that was never registered with the honest IdP. As the IdP only provides an authentication token when it receives a valid registration proof  $\pi$ ,  $\mathcal{A}$  must have forged this proof in its query for  $sid^*$ .  $\mathcal{A}$  can perform such an attack either by forging a proof on an invalid statement directly, or by forging the underlying witness, which is a tuple in the form of  $(rid, o, \sigma_{rid})$ . Here, by forging a witness, we mean either forging the membership credential  $\sigma_{rid}$  on a non-registered  $rid$ , or forging the opening  $(rid^*, o^*)$  of  $com^*$  where  $(rid, o)$  is also a valid opening to  $com^*$  for a registered but corrupted  $rid$ . Neither of these cases is feasible by relying on the special soundness of NIZK, EUF-CMA of  $S_2$  under  $pk_2$ , and the binding property of the commitment scheme. To be able to reduce to a forgery under  $S_2$ , we require NIZK to be special sound and use the knowledge extractor to obtain a valid witness  $(rid, o, \sigma_{rid})$ . By relying on the binding property of the commitment scheme, we can argue that  $\sigma_{rid}$  satisfies the winning condition of the EUF-CMA game.

In more detail, our reduction in condition (2d) works as follows. We obtain a challenge public key  $pk_2^*$  from a  $S_2$  unforgeability challenger and simulate the  $pk_2$  in the identity

provider public key  $ipk$  as  $pk_2^*$ . We simulate the  $S_2$  signatures for  $O.RegCRP/O.RegHRP$  queries by relying on the signing oracle of the  $S_2$  unforgeability challenger, so the behavior of the oracles changes as follows:

Alnit<sub>U</sub>, AReq<sub>RP</sub>, ARes<sub>IdP</sub>, AResFin, Vf<sub>RP</sub>: As they are.

RegHRP: It checks that  $(rid, \cdot) \notin HRP \cup CRP$ . If it does not hold, it outputs 0. Otherwise, for the registration query for  $rid$ , it queries the  $S_2$  signing oracle, gets the signature  $\sigma_{rid}$ , updates  $HRP := HRP \cup \{(rid, \sigma_{rid})\}$  on  $rid$ , and outputs 1.

RegCRP: It checks that  $(rid, \cdot) \notin HRP \cup CRP$ . If it does not hold, it outputs 0. Otherwise, for the registration query for  $rid$ , it queries the  $S_2$  signing oracle, gets the signature  $\sigma_{rid}$ , updates  $CRP := CRP \cup \{(rid, \sigma_{rid})\}$  on  $rid$ , and outputs  $\sigma_{rid}$ .

Finally, we run the knowledge extractor for the NIZK on the proof  $\pi$  which corresponds to the  $O.ARes_{IdP}$  query for  $(uid, ctx^* sid^*) \in Q_{\tau}$  and extract a valid witness  $(rid, o, \sigma_{rid})$ . By the special soundness property of the underlying NIZK, we know that the extractor will output a valid witness with overwhelming probability, so  $S_2.Vf(pk_2^*, rid, \sigma_{rid})$  and  $Open(rid, com^* o)$ .

If  $rid \neq rid^*$ , then we break the binding property of the commitment scheme as  $(rid, o)$  and  $(rid^*, o^*)$  are distinct valid openings to the commitment  $com^*$ . Otherwise,  $\sigma_{rid}$  is a valid signature on  $rid^* = rid$ . There is no signing oracle query to the  $S_2$  challenger for  $rid^*$  by condition (2d), so  $(\sigma_{rid}, rid^*)$  is a valid and fresh  $S_2$  forgery, breaking the EUF-CMA property of  $S_2$ .

We conclude that if there is a type (2d) forger adversary, the underlying NIZK is not special sound, the underlying commitment scheme is not binding, or  $S_2$  is not EUF-CMA.

- (e) *RP Authentication II:  $rid^* \in HRP \wedge (rid^*, sid^*) \notin Q_{auth}$* . When the adversary wins by satisfying the final condition, it has created a valid token for session  $sid^*$  and honest RP  $rid^*$ , yet this RP never provided authentication for that session. As in the previous case,  $\mathcal{A}$  can do this by forging  $\pi$  directly, finding a commitment collision for  $rid^*$  and a corrupted  $rid$ , or forging the membership credential  $\sigma_{rid}$  of the honest RP. To formally prove it, we use both the zero-knowledge and simulation extractability properties here.

We aim to show that the adversary must forge a NIZK proof  $\pi$  or a witness  $(rid, o, \sigma_{rid})$  as explained in condition (2d). Thus, similar to condition (2d), we must simulate  $S_2$  signatures by relying on an EUF-CMA challenger. However, unlike condition (2d),  $rid^*$  belongs to an honest RP here, so if we make a signing query for  $rid^*$  to the EUF-CMA challenger, a signature on  $rid^*$  is not a valid forgery anymore. Thus, we do not simulate honest RP credentials with  $S_2$  signatures, but we simulate the NIZK proofs  $\pi$ 's in AReq<sub>RP</sub> oracle queries for the honest RP's without knowing/creating a valid  $\sigma_{rid}$ . As a result, we cannot simply rely on special soundness as in condition (2d), but we will need NIZK to be simulation extractable. Moreover, changing only RegHRP and RegCRP is not enough, but we also need to change AReq<sub>RP</sub> so that the honest RP authentication requests can be created using the NIZK simulator. We change the oracles' behavior as follows:

Alnit<sub>U</sub>, ARes<sub>IdP</sub>, AResFin, Vf<sub>RP</sub>: As they are.

RegHRP: It checks that  $(rid, \cdot) \notin \text{HRP} \cup \text{CRP}$ . If it holds, it updates  $\text{HRP} := \text{HRP} \cup \{(rid, \perp)\}$  and outputs 1. If not, it outputs 0. It does not make a  $S_2$  signing query in any case.

RegCRP: It checks that  $(rid, \cdot) \notin \text{HRP} \cup \text{CRP}$ . If it does not hold, it outputs 0. Otherwise, for the registration query for  $rid$ , it queries the  $S_2$  signing oracle, gets the signature  $\sigma_{rid}$ , updates  $\text{CRP} := \text{CRP} \cup \{(rid, \sigma_{rid})\}$  on  $rid$ , and outputs  $\sigma_{rid}$ .

AREqRP: Checks if  $(rid, \cdot) \in \text{HRP}$  and returns  $\perp$  if not. Runs the original  $\mathcal{O}.\text{AREqRP}$  except for computing NIZK. As we do not know a valid credential for honest  $rid$ 's, the NIZK proof is simulated using the zero-knowledge simulator.

After simulating the adversary's view as above, we run the knowledge extractor for the NIZK on the proof  $\pi$  which corresponds to the  $\mathcal{O}.\text{AResIDP}$  query for  $(uid, ctx^*, sid^*) \in Q_\tau$  and extract a valid witness  $(rid, o, \sigma_{rid})$ . By the simulation extractability property of the underlying NIZK, we know that the extractor will output a valid witness with overwhelming probability, which satisfies  $S_2.\text{Vf}(pk_2^*, rid, \sigma_{rid})$  and  $\text{Open}(rid, com^*, o)$ .

Just as in condition (2d), if  $rid \neq rid^*$ , then we break the binding property of the commitment scheme using the distinct openings  $(rid, o)$  and  $(rid^*, o^*)$  to the commitment  $com^*$ . Otherwise,  $\sigma_{rid}$  is a valid signature on  $rid^* = rid$ . There is no signing oracle query to the  $S_2$  challenger for  $rid^*$  as  $rid^* \in \text{HRP}$ , so we output a valid forgery. We conclude that if there is a type (2e) forger adversary, the underlying NIZK is not simulation extractable, the underlying commitment scheme is not binding, or  $S_2$  is not EUF-CMA.  $\square$

## B.4 Proof of Theorem 5.4 (Req. Authentication)

Here we prove that  $\pi_{\text{OPPID}}$  satisfies Request Authentication (see Def. A.1) if the  $S_2$  scheme is EUF-CMA secure, and the NIZK is zero-knowledge and simulation extractable. The proof essentially follows the RP Accountability proof from [31].

**PROOF.** The proof of Request Authentication follows from the zero-knowledge and simulation extractability of NIZK, and the unforgeability of  $S_2$ . The main proof strategy is similar to condition (2e) of the Session Binding proof. We aim to build a  $S_2$  forger using an *auth*-forger, Request Authentication adversary. For that, we simulate the  $S_2$  public key in *ipk* as an EUF-CMA challenge public key  $pk_2^*$ . Note that we will not actually need the full power of EUF-CMA, as we will not make any signing queries. Thus, the unforgeability of  $S_2$  against a key-only attack, where no signing queries are allowed, would also be sufficient.

Just as in condition (2e) of the Session Binding, we do not create credentials for honest RPs, but we simulate the corresponding NIZK proofs to generate *auth* values for these RPs. The difference from condition (2e) of the Session Binding game is that corrupted RPs are not allowed, so we do not have to simulate their credentials. Thus, we do not need to make signing queries for  $pk_2^*$  at all. In the end, we extract a valid signature on some  $rid$  value from the forged *auth*, which is a valid  $S_2$  forgery. In more detail, we simulate the adversary's view as follows:

AResIDP: As it is.

RegHRP: It checks that  $(rid, \cdot) \notin \text{HRP} \cup \text{CRP}$ . If it holds, it updates  $\text{HRP} := \text{HRP} \cup \{(rid, \perp)\}$  and outputs 1. If not, it outputs 0. It does not make a  $S_2$  signing query in any case.

AREqRP: Checks if  $(rid, \cdot) \in \text{HRP}$  and returns  $\perp$  if not. Runs the original  $\text{AREqRP}$  algorithm except for computing NIZK. As we do not know a valid credential for honest  $rid$ 's, the NIZK proof is simulated using the zero-knowledge simulator.

Finally, when the adversary outputs an authentication request forgery  $(auth^*, crid^*, uid^*, ctx^*, sid^*)$ , we run the knowledge extractor for the NIZK on the proof  $auth^* := \pi^*$ , and extract a valid witness  $(rid, o, \sigma_{rid})$ . By the winning condition, we know that there is no  $\mathcal{O}.\text{AREqRP}$  query for  $(sid^*, crid^*)$ . As honest RPs bind  $(sid^*, crid^*)$  to the created proofs, we know that  $\pi^*$  is not output by  $\mathcal{O}.\text{AREqRP}$ , so it is not a previously simulated proof by the zero-knowledge simulator. Thus, by the simulation extractability property of the underlying NIZK,  $(rid, o, \sigma_{rid})$  is a valid witness with overwhelming probability, and thus  $S_2.\text{Vf}(pk_2^*, rid, \sigma_{rid}) = 1$ . As we do not make any signing queries to the unforgeability challenger,  $(\sigma_{rid}, rid)$  is a valid forgery against the EUF-CMA property of  $S_2$ .  $\square$