

# ORIGO: Proving Provenance of Sensitive Data with Constant Communication

Jens Ernstberger\*  
 jens.ernstberger@tum.de  
 Technical University of Munich  
 Munich, Germany

Jan Lauinger\*  
 jan.lauinger@tum.de  
 Technical University of Munich  
 Munich, Germany

Yinnan Wu  
 yinnan.wu@tum.de  
 Technical University of Munich  
 Munich, Germany

Arthur Gervais  
 arthur@gervais.cc  
 University College London  
 London, United Kingdom

Sebastian Steinhorst  
 sebastian.steinhorst@tum.de  
 Technical University of Munich  
 Munich, Germany

## ABSTRACT

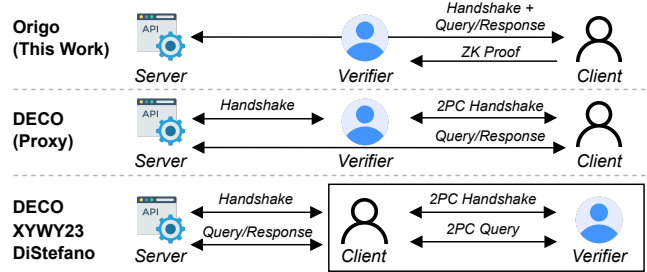
Transport Layer Security (TLS) is foundational for safeguarding client-server communication. However, it does not extend integrity guarantees to third-party verification of data authenticity. If a client wants to present data obtained from a server, it cannot convince any other party that the data has not been tampered with.

TLS oracles ensure data authenticity beyond the client-server TLS connection, such that clients can obtain data from a server and ensure provenance to any third party, without server-side modifications. Generally, a TLS oracle involves a third party, the verifier, in a TLS session to verify that the data obtained by the client is accurate. Existing protocols for TLS oracles are communication-heavy, as they rely on interactive protocols. We present ORIGO, a TLS oracle with constant communication. Similar to prior work, ORIGO introduces a third party in a TLS session, and provides a protocol to ensure the authenticity of data transmitted in a TLS session, without forfeiting its confidentiality. Compared to prior work, we rely on intricate details specific to TLS 1.3, which allow us to prove correct key derivation, authentication and encryption within a Zero Knowledge Proof (ZKP). This, combined with optimizations for TLS 1.3, leads to an efficient protocol with constant communication in the online phase. Our work reduces online communication by 375× and online runtime by up to 4.6×, compared to prior work.

## 1 INTRODUCTION

Ensuring sovereignty over digital content is imperative in contemporary computing systems and applications [33]. Recent endeavors focus on cryptographic verification of digital identities [55, 59], combating disinformation [46], and enforcing access control [47]. However, these solutions are hampered by a bootstrapping dilemma; they are rendered ineffective without a centralized root of trust that attests to the authenticity of the content.

Recent work introduces *TLS oracles*, which aim to solve this problem by cryptographically affirming the origin of digital content [74]. To do so, these solutions augment a standard TLS connection with a third party, such that the third party is oblivious to the content sent between a client and a server. The third party operates independently of the content-serving server and acts as a *verifier*, solely



**Figure 1: Systematization of system designs for TLS Oracles. Origo expands the existing set of systems by ensuring ciphertext integrity in TLS 1.3 with ZKPs.**

verifying the data exchanged. This approach empowers clients, allowing them to secure attestations for their digital content without demanding for clearance from a central application server.

Previous approaches to TLS oracles that do not require server-side changes can be divided based on whether they operate with the verifier as an external entity [23, 24, 70, 74] or with a verifier as a proxy between the client and the server [52, 74] (cf. Figure 1). All solutions aim to solve the same set of challenges – ensuring that the message presented by the client is consistent with the response sent by the server, ensuring that the verifier cannot provide an invalid response to the client, and providing confidentiality toward the verifier. To provide each of these guarantees, all protocols operate roughly in the same manner. First, a secure connection to the server is established by executing the TLS handshake and key derivation in 2PC. Second, the client and the verifier collaboratively send the request to the server, while the response is provided to the verifier *before* the client obtains the full decryption key. Third, the client generates a ZKP to prove a statement about the private payload. In proxy mode, 2PC post-handshake protocols can be discarded [74]. In particular, all previous approaches require an interactive protocol, with a linear communication in the number of gates in the 2PC circuit, for non-repetitive computations [69].

In this work, we introduce ORIGO, a TLS oracle with constant communication, independent of the size of the data requested from an application server. We leverage that in TLS 1.3, the prevalent version of TLS used to date, the IV for Authenticated Encryption

\*Both authors contributed equally to this research.

with Associated Data (AEAD) is derived from the traffic secret to facilitate a TLS oracle that does not demand for any communication intensive 2PC. By doing so, we can use a simple Succinct Non-interactive Argument of Knowledge (SNARK) to prove compliance with TLS 1.3, ensuring the same properties as achieved in previous work, under the assumption of a weaker network adversary that is unable to intercept traffic between the proxy verifier and server. With ORIGO, clients can obtain attestations on resource constrained devices, even in areas with low connectivity. Attestations can be used successively with arbitrary web applications and blockchains. Furthermore, we provide an efficient transformation of AEAD ciphertexts to SNARK friendly commitments, allowing clients to prove disjunct statements on data obtained – without demanding re-execution of the protocol.

**Technical Overview.** Thus far, the proxy setting for TLS oracles has only been sparsely described in academic works (cf. Figure 1). Therefore, we set out to clarify security assumptions, and possible attacks, to determine efficient solutions for data provenance.

**Security in the Proxy Setting.** We start by examining the security assumptions of TLS in the proxy model. Zhang *et. al* first introduce the proxy setting in the appendix of their seminal work [74]. They note that in the proxy setting, one has to assume a weaker network adversary, as the client could perform a machine-in-the-middle (MITM) attack between the client and the server. We find that this assumption is indeed a necessity, as none of the server-side messages are signed and TLS employs symmetric encryption. As a result, a MITM client can always simulate an interaction with an honest server. Further, Zhang *et. al* note that the handshake remains to be executed with *three* parties, as AES-GCM is non-committing as proven by Grubbs *et. al* [41]. Although AES-GCM is non-committing, the attack as described by Grubbs *et. al* is not trivially applicable, as a winning adversary needs to be able to choose an arbitrary ciphertext to win the security game. In Section 4, we introduce the notion of *strong receiver binding*, which applies to the proxy setting in TLS oracles. We further introduce an attack, showing that AES-GCM is *not* strong receiver binding.

**SNARKs for TLS 1.3.** Our protocol shows that a SNARK can be used to enforce strong receiver binding in the proxy setting of a TLS oracle in TLS 1.3. We derive a set of tailored optimizations that are specific to TLS 1.3 and minimize the number of constraints in the SNARK circuit. In comparison to the work of Grubbs *et. al* [40], our circuit requires  $\sim 3\times$  less constraints to derive a key in TLS 1.3.

**Client-side evidence generation.** We detail an extension to our protocol, which allows for expanded use-cases besides the issuance of trivial credential. We ensure that clients can generate evidence of ciphertext integrity by efficiently transforming AEAD plaintexts to hiding and binding commitments in TLS oracles. To do so, we let the client expand block-level nonces from the initialization vector derived from the TLS handshake in TLS 1.3. We facilitate efficiency by instantiating the HKDF.expand function with a SNARK friendly MiMC hash, and instantiate the commitment to TLS records with a salted hash. The generation of committed values augments existing circuits as an extension. We instantiate the implementation of MiMC hashes with Groth16 as the outer SNARK and recurse the inner verifier of the GKR protocol. Previous work shows that this optimization leads an order of magnitude improvement in in-circuit

hash computation [17]. We show that generating evidence for a 1kB TLS record only takes an additional 0.12 seconds when optimized with the recursive construction by Belling *et. al* [17].

In summary, our contributions are as follows:

- **A novel TLS oracle construction.** We introduce ORIGO, a protocol designed for proving data provenance with constant communication in TLS 1.3. We observe that in TLS 1.3, a simple SNARK is sufficient to ensure integrity when instantiating the third party in the TLS exchange as a proxy.
- **A novel attack to show equivocation.** We show that due to derivation of the traffic encryption key from the TLS handshake transcript, a malicious client could arbitrarily manipulate the plaintext, whilst retaining authenticity. We introduce the notion of strong receiver binding, and an attack that shows that an oracle in the proxy setting leveraging TLS 1.3 is *not* strong receiver binding.
- **Client-side evidence generation.** In TLS oracles, a client is unable to prove arbitrary statements on transmitted data once a session is completed. We describe an extension for client-side evidence generation, which allows transformation of an AEAD ciphertext to a SNARK friendly commitment to the respective plaintext, such that the client can efficiently generate proofs from attested ciphertexts, even after the session has concluded.
- **Evaluation and Application.** We implement ORIGO in a realistic setting with a real server [10]. For a server response of 1.4kB, the online communication required is 2.4kB –  $375\times$  less than the most performant related work relying on semi-honest 2PC. By globally dispersing clients, we show that ORIGO is less dependent on the client location, and up to  $4.6\times$  faster for remote clients. We leverage our real-world implementation and introduce a novel application by integrating ORIGO with Bitcoin Passport [5], extending it to be used with private and personal data (§ A).

**Limitations.** Notably, ORIGO assumes that the client cannot mount a MITM attack between the proxy verifier and the server in the TLS 1.3 session. For prover integrity, we therefore assume that the proxy can reliably connect to the server throughout the session. Although this setting assumes a weaker network adversary, we regard it as feasible due to existing detection techniques on other layers of the OSI stack [9, 64, 66, 75], and similar assumptions in other proxy applications [67, 74]. Our main focus is to enable improved performance, assuming that MITM between proxy and server is infeasible. Although we only specify the key derivation circuit for 1-RTT, one could adapt the circuit to 0-RTT by further deriving the session resumption secret in-circuit. However, previous research determined that less than 10% of connections with TLS 1.3 utilize 0-RTT [43] - further, clients can force 1-RTT to obtain a notarization.

Another limitation in the ORIGO protocol is the overhead introduced by the pre-processing required for SNARKs. In our evaluation, the pre-processing of circuits, especially to obtain the prover key, introduces a significant offline runtime overhead. Further, ORIGO inherits the randomized setup to sample a trapdoor from the respective SNARK utilized (in our case, Groth16). Plonk [35] would yield

**Table 1: Comparison of related works.  $n$  signifies the size of the request/response sent to/by the server.**

System	Setting	Legacy Compatibility	MITM Security	Communication
DECO [74]	External	✓	✓	$O(n)$
DECO/Proxy	Proxy	✓	✗	$O(n)$
XYWY23 [70]	External	✓	✓	$O(n)$
Janus [52]	Proxy	✓	✓	$O(n)$
DiStefano [23]	External	✓	✓	$O(n)$
DIDO [24]	External	✓	✓	$O(n)$
This Work	Proxy	✓	✗	$O(1)$

a universal and updatable SNARK. Fractal [25] would not require a trusted setup.

## 2 RELATED WORKS

**Proof Systems for TLS.** Many works use proof systems with TLS to filter network traffic [40], blindly obtain certificates [67], prove non-repudiation [22, 62], and provide off-chain data to blockchains [73, 74]. Many require server-side changes [22, 62], which are insufficient for adoption. We provide a comparison to the most relevant works as follows. TLS-N [62] is a server-side extension that allows for the generation of privacy-preserving proofs on the contents of a TLS session. As their work proposes a server-side extension, cooperation of servers in adoption of their protocol is crucial. Our protocol for evidence generation is similar to their proposed server-side extension, with the important difference that evidence is generated by the *client* in our setting, without the need for server-side adaptation. DECO[74] introduces the first TLS oracle for modern TLS versions without server-side adaptation (cf. Table 1). Initially derived from the PageSigner [2, 3] protocol first mentioned in 2014, their protocol asks the verifier and the client to jointly emulate a TLS client interacting with an unmodified server. In Appendix C.4 the authors of DECO describe an extension to their protocol which is similar to our protocol – the verifier acts as a proxy between the server and the client. However, their extension relies on 2PC during the TLS handshake in the extended protocol. Grubbs *et. al* [40] introduce *Zero Knowledge Middleboxes*, where clients send queries and simultaneously prove in zero knowledge that their traffic is compliant with a network policy. Here, the client only commits to its own request, rather than the response received from the server.

Further, there is a set of works conceived concurrently to our work. Xie *et. al* [70] introduce the *garble-then-prove* paradigm, which replaces authenticated garbled circuits and SNARKs with semi-honest garbled circuits and interactive ZKPs based on vector oblivious linear evaluation (VOLE). Their work introduces a  $14\times$  improvement in communication and up to  $15\times$  improvement in runtime over DECO. They also provide a protocol for converting ciphertexts to SNARK compatible commitments. Celi *et. al* [23] introduce *DiStefano* which guarantees ring privacy by using ring signatures produced over TLS certificates and relies on maliciously secure 2PC, similar to DECO. Lauinger *et. al* [52] introduce *Janus*, a TLS oracle tailored to TLS 1.3 in the proxy setting for verifying the provenance of kilobytes of data, relying on an honest verifier ZK proof system. Janus employs handshake *and* record layer 2PC on

AES encrypted counter blocks, rendering it secure against MiTM attacks, even in the proxy mode (it can also be operated in non-proxy mode, the proxy mode is outlined in the paper for operational convenience). All of the above utilize 2PC and require communication in size linear to the number of gates in the garbled circuit.

**Regular Expressions in Probabilistic Proofs.** A recent line of work focuses on optimizing regular expressions in probabilistic proofs, which can be utilized in TLS oracles to prove that the plaintext underlying a ciphertext fulfills a specified set of constraints. Although this problem is not the core focus of this work, we consider it important orthogonal work. Zombie [72] extends ZKMB [40] and allows proofs on DNS requests by transforming a regular expression to an NFA and uses the sum-check based SpartanNIZK as the probabilistic proof system. Luo *et. al* introduce zk-regex [54], which follows a similar NFA approach as Zombie but transforms the NFA to a Boolean circuit and uses MPC-in-the-head for the proof. Raymond *et. al* introduce zkreg [60] for matching regular expressions based on the Aho-Corasick automaton and use a commit-and-prove scheme to prove membership in large dictionaries of strings. Out of all, Zombie supports the most expressive expressions.

**TLS Security and Integrity.** Message Franking protocols follow a similar setting as ORIGO, as the participants engage in a protocol that similarly requires the receiver of a message to report upon the message content [29, 41, 44, 45, 65]. However, the setting of message franking differs to our setting, as the server verifying the message reported by a message receiver does not keep track of the ciphertext initially sent by the message sender. Recent works on message franking further depict attacks that allow content equivocation [41].

## 3 PRELIMINARIES

### 3.1 Cryptographic Primitives in TLS

TLS is a cryptographic protocol designed to provide secure communication over a computer network. It ensures that the data transmitted between a client and a server remains confidential and cannot be intercepted or tampered with by third parties. The most recent version is TLS 1.3 [61], which improves over TLS 1.2 [28] by reducing the number of applicable algorithms for message encryption and authentication, and further introducing a 0-RTT mode, which optionally enhances performance at the cost of security. TLS 1.3 relies on AEAD, a family of encryption algorithms which simultaneously ensure confidentiality and authenticity of transmitted data. On a high level, a TLS protocol consists of two layers – a *handshake layer*, where the server and the client negotiate the traffic encrypting keys (i.e., the Client Application Traffic Key CATK and Server Application Traffic Key SATK), and a *record layer*, which applies the traffic keys in an AEAD algorithm to securely send application data. We provide an overview of TLS 1.3 in one round-time trip (1-RTT) mode without pre-shared keys in Figure 2. The majority of the figure outlines the key derivation process, resulting in the keys that encrypt the client request and server response (client-CATK, server-SATK). Both keys are derived through a sequence of HKDF (HMAC-based Key Derivation Function [50]) invocations. We refer to [51] for rationale on the usage of HKDF in TLS 1.3.

**HKDF and HMAC.** Key Derivation from the established DHE secret HMAC-based Extract-and-Expand Key Derivation Function

(HKDF) follows the *encrypt-and-expand* paradigm, where the extraction function  $\text{HKDF.extract}(\text{SALT}, \text{IKM})$  derives a pseudorandom key from the keying material and the expansion function  $\text{HKDF.expand}(\text{PRK}, \text{INFO}, L)$  expands the pseudorandom key into another pseudorandom key of length  $L$ . Both  $\text{HKDF.extract}$  and  $\text{HKDF.expand}$  internally make use of an Hashed Message Authentication Code (HMAC) [49], which is augmented with a hash function  $H$  (usually SHA-256). An HMAC is computed as follows:

$$\text{HMAC} = H(K \oplus \text{OPAD}, H(K \oplus \text{IPAD}, M)) \quad (1)$$

The HMAC takes as input the message  $M$  and is keyed with a key  $K \in \{0, 1\}^\lambda$ . The security parameter  $\lambda$  determines padding length in the inner and outer hash [49].

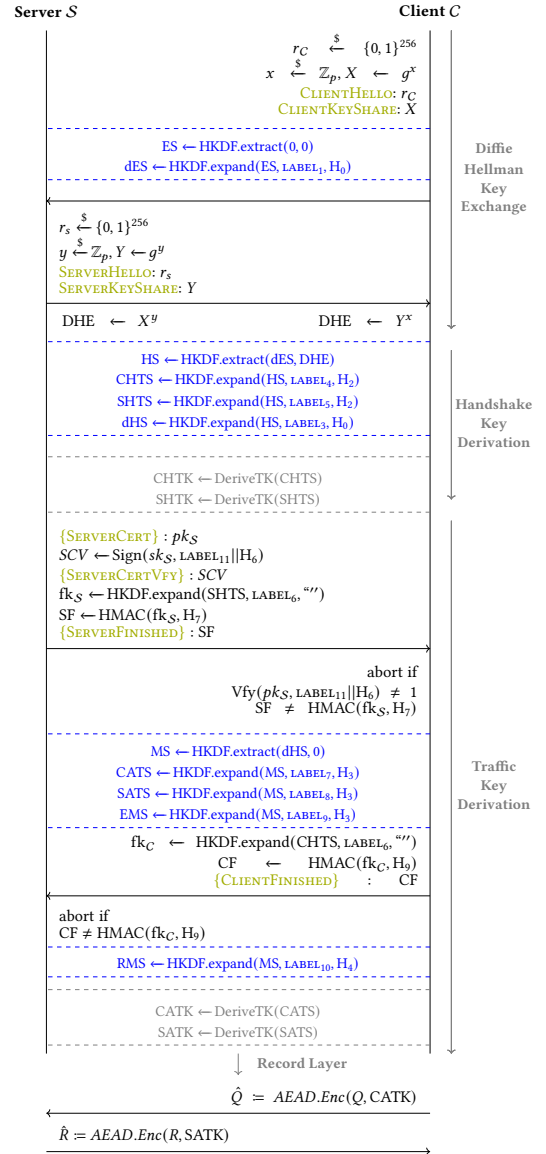
**AEAD.** In TLS 1.3, the application data is fragmented in *records*, where each record has a maximum size of 16kB, such that one record contains a maximum of 1024 chunks of size 16 byte or 128 bit. Each record is encrypted and authenticated through AEAD. TLS 1.3 supports a set of ciphersuites, namely `TLS_AES_128_GCM_SHA256`, `TLS_AES_256_GCM_SHA384` and `TLS_CHACHA_20_POLY1305_SHA256`, where support for the first is mandatory and the latter two are optional. Hence, the most common mode of operation for the underlying block cipher is AES in Galois Counter Mode (GCM) with an output length of 16 bytes or 128 bit. Importantly, the traffic sent by the client and server in the record layer are encrypted with their respective symmetric keys  $\text{CATK}$  and  $\text{SATK}$ . Both keys are known by the client and server. In the remainder of this work, we assume AES-GCM is applied for both encryption and decryption, where the ciphertext carries a 128-bit authentication tag. We detail a simplified description of GCM with AES as applied in TLS 1.3 to encrypt record layer data in Figure 3. For simplicity, we ignore associated data.  $\text{AEAD}_{\text{GCM}}(K, IV, \{P_1, P_2, \dots, P_m\})$  uses the key  $K$  with the AES blockcipher  $E$  of block size 128 bit. It encrypts the plaintext chunks  $\{P_1, P_2, \dots, P_m\}$ , where  $m$  is maximally 1024, and outputs the Tag  $T$  and the chunks of ciphertext  $\{C_1, C_2, \dots, C_m\}$ . The tag is computed with an  $\epsilon$ - $AXU$  hash function that relies on binary polynomials of degree 128 in the Galois Field  $GF(2^{128})$ . The polynomial is evaluated at the encryption of the 0-vector  $H$  and successively  $XOR$ -ed with the encryption of the initialization vector  $IV$  to obtain the tag  $T$  (for details, see the GHASH Function [1]).

### 3.2 zkSNARKs

In the realm of cryptographic protocols, zero-knowledge proofs (ZKPs) enable one party, referred to as the prover, to convince another, the verifier, of the veracity of a statement, while not disclosing any other information aside from its validity. At its core, the objective of a ZKP is to establish that a given *instance*  $x$  and its associated *witness*  $w$  conform to an NP relation  $\mathcal{R}$ . This can be formally depicted by the relationship  $(x, w) \in \mathcal{R}$ , which can be ascertained in polynomial time by a Turing machine.

A typical ZKP encompasses three algorithms:

- $\text{Setup}(\text{pp}) \rightarrow (\text{pk}, \text{vk})$ . Given public parameters  $\text{pp}$  as input, compute proving and verification keys  $\text{pk}$  and  $\text{vk}$ .
- $\text{Prove}(\text{pk}, x, w) \rightarrow \pi$ . Given the proving key  $\text{pk}$ , the instance  $x$ , and the witness  $w$ , such that  $(x, w) \in \mathcal{R}$ , as input, compute and output a proof  $\pi$ .



**Figure 2: TLS 1.3 1-RTT handshake transcript without optional messages and client certificate verification. Secrets and keys derived by both server and verifier are highlighted in blue and gray respectively. Messages sent are denoted in green. Encrypted messages are in brackets. We denote all plain messages as "handshake transcript"  $\Gamma$  (encrypted as  $\hat{\Gamma}$ ). We outline TLS 1.3 specific variable definitions in Table 4.**

- $\text{Verify}(\text{vk}, x, \pi) \rightarrow 0/1$ . Given the verification key  $\text{vk}$ , the instance  $x$ , and the proof  $\pi$  as input, output 1 if the proof is valid and 0 otherwise.

We say that  $(\text{Setup}, \text{Prove}, \text{Verify})$  is a zkSNARK if it satisfies completeness, succinctness, knowledge soundness and zero-knowledge. Completeness ensures that for any genuine statement, an honest prover should be able to convince an honest verifier. Knowledge

```

AEADGCM( $K, IV, \{P_1, P_2, \dots, P_m\}$ )
  \ In TLS 1.3 the  $IV$  is derived from the handshake as
  \  $IV \leftarrow \text{HKDF.expand}(SATS, "iv", I)$ 
   $R_i \leftarrow E_K\{IV || ctr_i\}$ 
   $C_i \leftarrow R_i \oplus P_i, \forall i \in m$ 
  -----
   $H \leftarrow E_K\{0^{128}\}$ 
   $X \leftarrow \sum_{i=1}^m C_i H^{m-i}$ 
   $T \leftarrow E_K\{IV\} \oplus X$ 
  Return  $\{C_1, C_2, \dots, C_m\}, T$ 
    
```

**Figure 3: Description of GCM, where  $E_K$  is AES. In TLS 1.3, the nonce applied in GCM is directly derived from the handshake with a hash-based key derivation function. For clarity, GCM is simplified as compared to the NIST specification [1].**

soundness, informally, describes that for every adversarial prover, there exists an *extractor*, such that whenever the prover convinces the verifier to output 1, the extractor can use the prover algorithm to output  $w$  such that  $(x, w) \in \mathcal{R}$  [16]. Further, a SNARK is succinct if the size of the proof and the verification time are  $o(|R|)$ .

### 3.3 Security Guarantees & Threat Model

Throughout this paper, we assume the same threat model as in existing literature on TLS oracles, i.e. (i) the server is assumed to behave correctly and respond honestly to any request sent by a client, (ii) the client a proxy verifier do not collude. In addition, assume that an adversarial client cannot perform a MITM attack between the server and the proxy, i.e. the proxy is able to verify the identity of the server it is connected to. Note that this assumption only applies to the connection between proxy and server, as the client has the traffic encrypting keys, and could reliably forward arbitrary messages to the proxy, without the proxy realizing. A practical implementation of this assumption may sample a notary at random from a set of potential notaries, such the client cannot predict the link between the proxy verifier and server (we discuss the practicality of this assumption in Appendix D). Beyond our additional assumption regarding MITM, we assume that security, as proven in numerous previous works [34], holds for the standard TLS protocol in version 1.3. Hence, the security properties we aim to achieve, given the above assumptions, are the following:

**Client Integrity.** A malicious client cannot convince a verifier of an untrue statement about the TLS session.

**Verifier Integrity.** A malicious verifier cannot prevent a client from receiving the correct message from the server.

**Privacy.** A malicious proxy verifier cannot obtain any knowledge beyond the size of the message that it forwards.

## 4 EQUIVOCATION IN THE STRAWMAN CASE

In a strawman case, a TLS Oracle might rely on the security of plain TLS 1.3, where the verifier simply proxies the traffic between a server and a client. In this section, we detail an attack that allows the receiver of a ciphertext  $\{C_1, C_2, \dots, C_m\}$  (i.e. encrypted TLS record received by the client) to equivocate about the underlying plaintext  $\tilde{P} = \{P_1, P_2, \dots, P_m\}$ . To do so, we reason about an adapted notion of *receiver binding*, as introduced by Grubbs *et. al* [41]. Grubbs *et. al*

```

sr - BINDAEADA
   $((K, IV, P), (K', IV', P'), C, T) \leftarrow \mathcal{A}$ 
   $b \leftarrow \text{AEAD.Ver}(K, IV, P, C, T)$ 
   $b' \leftarrow \text{AEAD.Ver}(K', IV', P', C, T)$ 
  if  $(K, IV, P) == (K', IV', P')$ :
    return false
  return  $(b = b' = 1)$ 
    
```

**Figure 4: Strong Receiver Binding Security Game. An adversary wins if it finds a differing key  $K$ , nonce  $IV$  and plaintext  $P$  for a fixed ciphertext  $C$  and authentication tag  $T$ .**

detail an attack for GCM that allows a receiver to equivocate about the underlying plaintext by arbitrarily choosing the ciphertext  $\{C_1, C_2, \dots, C_m\}$  in the context of “Message Franking” [41], where the receiver of a message intends to report the content of a received message to a third party. The third party needs to be sure that the reported message is the one initially sent by the sender. It is hence shown that GCM is not *receiver binding* due to the the applied Carter-Wegman MAC not being collision resistant [41]. In a TLS Oracle that simply proxies the communication between client and server, the proxy similarly needs to be sure that the plaintext message received by the client was initially sent by the server, whereas the proxy already holds the ciphertext.

**Strong Receiver Binding for AEAD.** To account for the proxy setting, we introduce an orthogonal notion called *strong receiver binding*, which further restricts the adversary in manipulating the ciphertext and the associated MAC. This is important in the proxy setting of a TLS oracle, as in this case the proxy obtains a ciphertext from a server, which is assumed to behave honestly. For a committing AEAD scheme to be *receiver binding*, the security game by Grubbs *et. al* requires an adversary to output a pair of triples containing associated data, a message, an opening and a so-called “franking tag” (i.e. the commitment to the ciphertext). The adversary wins if verification succeeds on both triples with the franking tag and the associated data/message pairs differ. In the following we ignore associated data for simplicity, and consider the tag of GCM as the committing portion of AES-GCM (see [41]). For GCM, verification ( $\text{AEAD.Ver}()$  in Figure 4) re-computes encryption and checks that the resulting tag *and* ciphertext match the initial tag and ciphertext. To break *strong receiver binding*, an adversary needs to find a tuple  $(K', IV', P')$  that decrypts  $(C, T)$ , which is initially derived from a differing tuple  $(K, IV, P)$ . We show in the following that AES-GCM, as applied in TLS 1.3, is not only not receiver binding, but also not strong receiver binding, and the strawman approach is insecure without further ado. To do so, we formalize the pseudocode of the security game  $\text{sr - BIND}_{\text{AEAD}}^{\mathcal{A}}$  in Figure 4. The adversary  $\mathcal{A}$  has to output a pair of triples specifying the symmetric key, nonce and plaintext. The adversary wins if verification succeeds for both tuples, given the same ciphertext and tag. The formal advantage of an adversary  $\mathcal{A}$  can be specified as

$$\text{Adv}^{\text{sr-BIND}}(\mathcal{A}) = \Pr[\text{sr-BIND}_{\text{AEAD}}^{\mathcal{A}} \rightarrow \text{true}]$$

Note that strong receiver binding implies receiver binding security.

We give an adversary  $\mathcal{A}_{GCM}$  which wins if it outputs  $((\tilde{P}, IV, K), (\tilde{P}', IV', K'), \{C_1, C_2, \dots, C_m\}, T)$ , such that

$$AEAD_{GCM}(K, IV, \tilde{P}) = AEAD_{GCM}(K', IV', \tilde{P}')$$

Choosing  $(K', IV', \tilde{P}')$  is trivial, as the adversary must solve the following equation in order to obtain a conflicting decryption to the ciphertext  $\{C_1, C_2, \dots, C_m\}$  and authentication tag  $T$ :

$$E_K\{IV\} + E_{K'}\{IV'\} + \sum_{j=1}^m (R_j \oplus P_j)H^{m-i} \\ + \sum_{j=1}^m (R'_j \oplus P'_j)H^{m-i} = 0$$

To compute a valid Tag  $T$ ,  $\mathcal{A}_{GCM}$  chooses a key  $K'$ , computes  $X$  and chooses  $IV'$  such that  $T = T'$ . To successively compute  $\tilde{C}$  and  $\tilde{C}'$ , such that  $\tilde{C} = \tilde{C}'$ ,  $\mathcal{A}_{GCM}$  computes the series  $R'_i$  from  $IV'$  and computes  $\tilde{P}$  such that the above equation is fulfilled. We elaborate on the practical feasibility of this attack in Appendix C.

**Why receiver-binding holds for AES-GCM in TLS 1.2.** In TLS 1.2, the IV for AEAD cipher suites is generated using a pseudorandom function (PRF) based on the master secret and a nonce value. The nonce value is generated by the sender and included in the record header of the TLS message. The IV is derived by applying the PRF to the master secret and the nonce value and using the first N bytes of the output as the IV, where N is the size of the IV required for the particular AEAD cipher suite being used. In TLS 1.3, the IV is generated using a key derivation function (KDF) based on the secret and the record sequence number, and it is used to initialize the AEAD cipher. The secret is a value derived from the master secret and the handshake hash, and the record sequence number is an incrementing counter that is used to differentiate each record in a TLS session. The IV is not transmitted in the header of the TLS record. Instead, it is used internally by the AEAD cipher to encrypt and decrypt the message payload. Consequently, the recipient of a message over TLS 1.3 can arbitrarily choose the IV when reporting to a third party. We ensure that this attack is not possible in our protocol tailored to TLS 1.3 by deriving the IV in a SNARK and proving correct encryption with the correct application traffic key.

## 5 THE ORIGO PROTOCOL

In this section we describe ORIGO, a novel TLS oracle in the proxy setting. The protocol operates in three phases, detailed as follows:

- **1 - Handshake & Request Execution.** First, the client engages in the handshake with the server, sends a request and obtains a response which is fragmented into records of max. 16kB. The handshake, request and response is routed through the proxy. The proxy obtains the encrypted TLS 1.3 handshake transcript, the encrypted query and the encrypted server response.
- **2 - Pre-Processing.** During pre-processing, the client generates public values from the handshake transcript to generate an optimized proof of key derivation.
- **3 - Proof Generation & Verification.** The client utilizes the pre-processed input to generate a ZKP that (i) ensures

that the key is derived correctly, and (ii) that the data transmitted provides evidence to satisfy a given predicate. The pre-processed public values and the proof are transmitted to the verifier, which verifies the server certificate and the provided ZKP to ensure non-equivocation by the client. If successful, the verifier attests that the given predicate is satisfied by returning a signature.

In the following, we first provide a formal specification of ORIGO,  $\Pi_{ORIGO}$ , where the proxy verifier provides an attestation to a pre-defined claim in the form of a signature. In Appendix E we further show that  $\Pi_{ORIGO}$  securely realizes the ideal functionality for data provenance  $\mathcal{F}_{DP}$ . In Section §5.2 we describe optimizations for proof generation by optimizing the proof of key derivation through technical optimizations to minimize the number of constraints in the final protocol. In §5.3, we describe client-side evidence generation, which allows transforming an AEAD ciphertext to a SNARK friendly commitment, such that clients can generate proofs from attested ciphertexts even after the session concluded.

### 5.1 Formal Protocol Specification

We formally specify our full protocol  $\Pi_{ORIGO}$  in Figure 5.  $\Pi_{ORIGO}$  relies on  $\mathcal{F}_{ZK}$  (cf. [68]) to abstract away the complexity of ZKPs in the protocol, and hence operates in the  $\mathcal{F}_{ZK}$ -hybrid model.  $\Pi_{ORIGO}$  is parametrized by a signature scheme  $\Sigma(KGen, Sig, Vrf)$  and an AEAD algorithm  $AEAD(Enc, Dec)$ , defined by the ciphersuite (TLS\_AES\_128\_GCM\_SHA256 & Curve P-256).

**Notation.** In our formal specification, we denote the TLS server as  $\mathcal{S}$ , the proxy verifier as  $\mathcal{V}$  and the TLS client as  $\mathcal{C}$ . In ORIGO, the verifier  $\mathcal{V}$  acts as a Proxy between  $\mathcal{C}$  and  $\mathcal{S}$ .  $\Pi_{\mathcal{S}}$  denotes the unmodified server protocol as specified in TLS 1.3. We further denote the handshake transcript between  $\mathcal{S}$  and  $\mathcal{C}$  as  $\Gamma$ . The encrypted TLS request  $\hat{Q}$  and response  $\hat{M}$  are sent by  $\mathcal{C}$  and  $\mathcal{S}$  respectively. We denote a request response pair as a tuple  $(\hat{Q}, \hat{M})$ . In general, we denote encrypted values with a caret, e.g.,  $\hat{M}$  for  $M$ .

**a) Handshake.** To send a request, the client  $\mathcal{C}$  first has to engage in a TLS 1.3 handshake with the server  $\mathcal{S}$  by initializing the protocol. As discussed in Section 3 TLS handshake we assume in our protocol is standard, besides that each message from the handshake is routed through  $\mathcal{V}$ , such that  $\mathcal{V}$  obtains the handshake transcript  $\hat{\Gamma}$ , which is encrypted after establishment of the handshake traffic keys (SHTK and CHTK in Figure 2).

**b) Request Execution.** Upon establishment of the traffic encrypting keys,  $\mathcal{C}$  is ready to send a request  $\hat{Q}$  to  $\mathcal{S}$ . The client  $\mathcal{C}$  builds the request by invoking "request". The client embeds the secret for authentication with the server  $\theta_{\mathcal{S}}$  (e.g., API key) and encrypts the request with the client application traffic key CATK. The client successively calls the "request" subroutine of the verifier  $\mathcal{V}$  with the request  $\hat{Q}$ . The verifier  $\mathcal{V}$  forwards  $\hat{Q}$  to  $\mathcal{S}$  and, under the assumption that  $\mathcal{S}$  operates honestly, receives a record  $\hat{R}$  from  $\mathcal{S}$ . The proxy stores  $\hat{R}$  and forwards the record to  $\mathcal{C}$ .

**c) Pre-Processing.** Upon invocation of "response", the client  $\mathcal{C}$  decrypts the record to obtain the plaintext record  $R$ . The client proceeds to pre-processes the record by finding the indices of AES blocks which match the provided ctx by invoking  $\Pi_{postR}$  (see Figure 6). If the record in question matches the context,  $\mathcal{C}$  outputs the

$\Pi_{ORIGO}(\Sigma, AEAD, S, \mathcal{V}, C)$

A client  $C$  interacts with a verifier  $\mathcal{V}$  to obtain data from a server  $S$ . The server  $S$  inputs a secret key  $sk_S$  and a certificate  $pk_S$  to authenticate the data sent. The context  $ctx$  describes the key/value pair of interest. Further, the protocol is parametrized with a circuit  $C$ .

01: Server  $S$ : Follow standard TLS 1.3 in 1-RTT (cf. Figure 2).

02: Client  $C$ :

03: Initialize:

04: Send **CLIENTHELLO**  $\equiv r_c \xleftarrow{\$} \{0, 1\}^{256}$  to  $S$ . Parse  $\Gamma$  from interaction with  $S$ ,

05: obtain CHTS, CATK, SHTS, SATK and verify the server certificate  $pk_S$ .

06: **On input** ( $sid$ , "request",  $\theta_S$ ) from environment  $\mathcal{Z}$ :

07: set  $Q$  with  $\theta_S$

08: set  $\hat{Q} \equiv AEAD.Enc(Q, CATK)$

09: send ( $sid$ , "request",  $\hat{Q}$ ) to  $\mathcal{V}$

10: **On receive** ( $sid$ , "response",  $\hat{R}$ ) from environment  $\mathcal{Z}$ :

11: set  $R \equiv AEAD.Dec(\hat{R}, SATK)$

12: run  $\Pi_{postR}(R, ctx)$  locally to obtain  $p_R$

13: run  $\Pi_{postH}(\Gamma)$  to obtain  $\Gamma'$

14: extract  $dHS^{in}$  and set  $w \equiv (SATK, \theta_S, dHS^{in}, R)$

15: set  $x \equiv (p_R, \Gamma', \hat{R})$

16: send ( $sid$ , "prove",  $C, x, w$ ) to  $\mathcal{F}_{ZK}$

17: send ( $sid$ , "proof",  $\Gamma', p_R$ ) to  $\mathcal{V}$

18: Proxy  $\mathcal{V}$ :

19: Initialize:

20: Initialize  $(pk_{\mathcal{V}}, sk_{\mathcal{V}}) \leftarrow \Sigma.KGen(1^\lambda)$ ,

21: transcribe encrypted handshake transcript  $\hat{\Gamma}$  (handshake proxied by  $\mathcal{V}$ ).

22: **On receive** ( $sid$ , "request",  $\hat{Q}$ ) from  $C$ :

23: store  $\hat{Q}$  and send ( $sid$ , "request",  $\hat{Q}$ ) to  $S$

24: **On receive** ( $sid$ , "response",  $\hat{R}$ ) from  $C$ :

25: store  $\hat{R}$  and send ( $sid$ , "response",  $\hat{R}$ ) to  $C$

26: **On receive** ( $sid$ , "proof",  $\Gamma', p_R$ ) from  $C$ :

27: Parse  $\Gamma'$ , derive SHTK from  $SHTS^{in}$

28: set  $\Gamma \equiv AEAD.Dec(\hat{\Gamma}, (SHTK))$

29: invoke  $\Pi_{Auth}$  and verify the server certificate  $pk_S$

30: set  $x \equiv (p_R, \Gamma', \hat{R})$

31: send ( $sid$ , "verify",  $C, x$ ) to  $\mathcal{F}_{ZK}$

32: wait for  $\mathcal{F}_{ZK}$  to respond, if *false* return  $\perp$

33: send  $\sigma \equiv \Sigma.sign(ctx, \hat{R})$  to  $C$

**Figure 5: Formal Specification of ORIGO.** The protocol is parametrized by a signature scheme  $\Sigma(KGen, Sig, Vf)$  and an AEAD scheme  $AEAD$ , specified by the ciphersuite in TLS 1.3. The circuit utilized proves (i) derivation of the server key (SATK) of the TLS session, (ii) encryption of the plaintext to the ciphertext with the server key and (iii) matching a specific part of the plaintext against a publicly available policy.  $\theta_S$  is the value that allows the client to obtain access to the server (e.g. an API key).

pre-processed transcript  $\Gamma'$ . The pre-processed values in  $\Gamma'$  are sent to  $\mathcal{V}$  in plain and decrease the number of HMAC evaluations in the Key Derivation Circuit (cf. Figure 10). The secret input to the ZKP is  $dHS^{in}$ , as specified in line 11 of Figure 10, which depicts all computations for key derivation, including circuit optimizations introduced in the following Section 5.2 (note that "in" denotes a value that results from the inner compression function in HMAC).

**d) Proof Generation & Verification.** To prove data provenance,  $C$  needs to prove the integrity of the plaintext underlying the obtained TLS record by (i) proving the key derivation (cf. Figure 10), (ii) proving that the authentication tag is correct, and (iii) proving that the private input  $R$  encrypts to the encrypted response  $\hat{R}$  stored by the verifier  $\mathcal{V}$ . Additionally, the client proves that the TLS record in question satisfies the statement as defined by the context

$\Pi_{postR}(R, ctx)$

firstBlock  $\leftarrow 0$ , lastBlock  $\leftarrow 0$ , offset  $\leftarrow 0$ , found  $\leftarrow 0$

for  $i \leftarrow 1$  to  $n$ :

if  $ctx[0] \in B_i$ :

firstBlock  $\leftarrow B_i$ ; offset  $\leftarrow B_i[ctx[0]] - 16 \cdot i$ ; found  $\leftarrow 1$

if  $ctx["."] \in B_i$ :

pos  $\leftarrow B_i[ctx["."]]$

if  $ctx[|ctx|] \in B_i$ :

lastBlock  $\leftarrow B_i$

set  $p_R \leftarrow (firstBlock, lastBlock, offset, pos)$

return  $p_R$

**Figure 6: Protocol  $\Pi_{postR}$  for post-processing the record to identify the correct indices of AES Blocks containing the context.** Gray values identify floating point values.

$\Pi_{Auth}(SHTS, LABEL_6, LABEL_{11}, H_6, H_7, SF')$

1:  $fk_S \leftarrow HKDF.expand(SHTS, LABEL_6, "")$

2:  $SF \leftarrow HMAC(fk_S, H_7)$

3:  $SF' \stackrel{?}{=} SF$

4:  $Vfy(pk_S, LABEL_{11} || H_6) \neq 1$

**Figure 7: Out of circuit verification of the server-side certificate to ensure authenticity of the obtained ciphertext and handshake transcript.**  $SF'$  is derived from the disclosed SHTS.

$ctx$ . Formally, we adopt the ZKP functionality  $\mathcal{F}_{ZK}$  as described in [68]. In  $\mathcal{F}_{ZK}$ ,  $C$  and  $\mathcal{V}$  provide their private and public inputs, and  $\mathcal{F}_{ZK}$  sends *true* to  $\mathcal{V}$  if the private witness provided by  $C$  satisfies the circuit  $C$ . The client utilizes  $\mathcal{F}_{ZK}$  by assigning the private input  $w$ , which includes the server application traffic key SATK, the authentication secret  $\theta_S$ ,  $dHS^{in}$  and the plaintext record  $R$ . The public input  $x$  is assigned accordingly, and the client  $C$  invokes "prove" in  $\mathcal{F}_{ZK}$ . The client sends the post-processed transcript and the post-processed record by calling "proof" at the verifier  $\mathcal{V}$ .

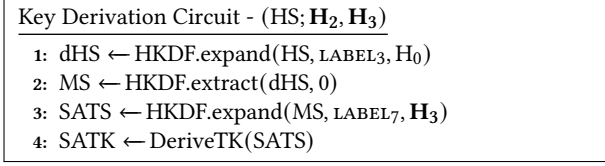
The verifier parses the pre-processed transcript  $\Gamma'$  and obtains  $SHTS^{in}$ , through which it can derive the handshake encrypting key SHTK. Successively, the verifier  $\mathcal{V}$  decrypts the encrypted handshake. The client utilizes  $\Pi_{Auth}$  as a subroutine to verify the server-side certificate and check the validity of  $SF$ . Upon successful verification,  $\mathcal{V}$  sets as public input the post-processed values  $p_R$  obtained through the context, the post-processed transcript  $\Gamma'$  and the stored encrypted record  $\hat{R}$ . The verifier  $\mathcal{V}$  calls "verify" in  $\mathcal{F}_{ZK}$ , and attests with a signature once  $\mathcal{F}_{ZK}$  returns.

We characterize the security of  $\Pi_{ORIGO}$  through the following theorem (formal proof in Appendix E):

**THEOREM 1.** *Assuming ECDSA is EUF-CMA secure and the compression function  $f$  is a random oracle,  $\Pi_{ORIGO}$  UC-realizes the ideal functionality for data provenance ( $\mathcal{F}_{DP}$ ) in the ( $\mathcal{F}_{ZK}$ )-hybrid world.*

## 5.2 Optimizations for Proof Generation

Proof generation for ORIGO involves proving (i) correct key derivation, (ii) correct encryption in-circuit and (iii) that a specific part of the plaintext matches against a publicly available policy. Whereas



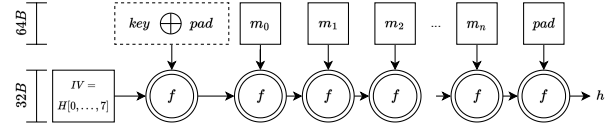
**Figure 8: ZKP sub-circuit to prove derivation of the server application traffic key SATK in TLS 1.3. HS is the private input, public inputs are bold. Other variables are constants.**

proving correct encryption is also a concern of orthogonal solutions to the TLS oracle problem, proving key derivation is unique to the proxy setting ORIGO operates in as it presents an alternative way to prevent equivocation by the client (cf. Section 4). We first describe the naive approach to proving key derivation (cf. Figure 3, which depicts unoptimized key derivation), and successively show how to “strip down” the number of in-circuit computations in several steps (cf. Figure 10 which depicts fully optimized key derivation).

**Naive Approach.** In the TLS 1.3 handshake, the client and server first execute a Diffie-Hellman Key Exchange to derive the shared DHE. To do so, e.g. the client samples a 256-bit nonce  $r_C \leftarrow \{0, 1\}^{256}$  which is sent as **CLIENTHELLO** including algorithm selections, and the **CLIENTKEYSHARE** extension which includes the public key shares  $x \leftarrow \mathbb{Z}_p$  and  $X \leftarrow g^x$  for the key exchange. The shared secret on the server side is derived as  $\text{DHE} \leftarrow X^y$ . Now, given DHE, the *first* secret derived from a pseudorandom function is HS, which allows server and client to further derive the handshake encrypting keys CHTK and SHTK, which encrypt messages in the authentication phase. In the authentication phase, the server authenticates to the client with its long-term PKI certificate, and computes the authentication value  $\text{SCV} \leftarrow \text{Sign}(sk_S, \text{LABEL}_{11} || H_6)$ , which is a signature over the continuously updated session hash. Given authentication, both parties proceed to derive application traffic keys.

A naive approach to proving key derivation would be to reimplement the functions invoked in the TLS handshake in an arithmetic circuit to successively apply a ZKP for circuit satisfiability. This would involve in-circuit certificate verification to ensure authentication of  $Y$ , and thus, HS, encryption of all handshake messages under CHTK and SHTK as well as derivation of the application traffic keys CATK and SATK. Given the number of hash invocations, as well as the complexity of in-circuit encryption which for AES demands for emulation of  $GF(2^8)$ , this approach would be prohibitively expensive. In the following, we show how to minimize the number of in-circuit operations without forfeiting security. As we want to prove properties on the server response, we only optimize the sub-circuit that computes the key that encrypts the server response with the server application traffic key SATK. The client has knowledge of all derivations and exchanged messages.

**1 - Certificate Verification & Handshake Encryption.** Due to the key independence property of TLS 1.3 [30], the client can disclose the server handshake traffic secret SHTS to the proxy without compromising HS and record layer application traffic secrets. This is possible, as HKDF.expand utilizes a one-way collision resistant hash function over HS, which is derived from sufficient entropy. The proxy utilizes SHTS to derive the handshake encrypting key



**Figure 9: Circuit optimization based on Merkle-Damgård structure. Dotted boxes are private inputs, solid boxes are public inputs. For a private key, the prover computes the first  $f$  in-circuit. All remaining values can be computed out of circuit and compared with the public hash  $h$ .**

SHTK to decrypt handshake traffic. As such, the proxy also decrypts the initially sent  $\text{SF}'$ , which is used as an input in the out-of-circuit verification algorithm to ensure valid certificates (cf. SHTS and signature verification in Figure 7).

**2 - Minimal Key Derivation.** Given the decrypted handshake transcript, we depict the minimal circuit for derivation of the server application traffic key in Figure 8, following the notation as introduced in Section 3 and Dowling *et. al* [31]. Notice that transcripts hashes  $H_2, H_3, H_7$  can be computed locally by the proxy verifier, given that SHTS is disclosed by the client. The circuit depicted still involves 12 invocations of SHA-256, given that each HKDF.expand invokes one HMAC for the given input sizes and SHA-256 with output size 32 Byte, HKDF.extract constantly invokes one HMAC and DeriveTK involves two invocations of HKDF.expand.

**3 - HMAC Optimization.** We further optimize the key derivation circuit by leveraging the structure of SHA-256, as applied in HMAC and TLS 1.3. We consider the optimization of general interest, as HKDF.expand invokes HMAC iteratively, increasing in number of rounds given a larger key size. With `TLS_AES_128_GCM_SHA256` as a ciphersuite, 128 bit encryption keys are used such that both functions HKDF.extract and HKDF.expand call HMAC once internally. With `TLS_AES_128_GCM_SHA256`, the concatenation of the inner hash  $H((K' \oplus ipad) || m)$  (32 bytes) and  $K' \oplus opad$  (64 bytes) yields a 96 byte output, which in turn, is the input to the outer hash function. The input to the inner hash function is of size  $64 + \text{len}(m)$  bytes. Thus, both hash input sizes in HMAC are above 64 bytes. If the hash input of SHA256 is above 64 bytes, SHA256 applies the Merkle-Damgård structure which repetitively invokes a compression function  $f$  to reduce the input to a fixed sized output (cf. Figure 9). In SHA-256,  $f$  is instantiated by the compression blockcipher SHACAL-2, which uses 64 computation rounds to hide its input [53]. Depending on whether the inner or outer hash is computed, the first call of the one-way compression blockcipher inside SHA256 already hides inputs  $(K \oplus ipad)$  or  $(K \oplus opad)$  of 64 bytes and thereby also hides the secret input  $K$  [74]. As a result, only the first input of the iteratively applied compression function in SHA-256 can be used as a private input, reducing the number of constraints of each SHA-256 invocations by 50%. Note, that this optimization assumes that  $f$  is PRF-secure and collision resistant, which is an assumption that is sufficient in our case as it is also assumed to show that HMAC is a dual-PRF [14].

We summarize all required computations to prove the TLS 1.3 key derivation in Figure 10. Computations executed at the proxy (Verifier) and client (Prover) are denoted as  $\mathcal{V}$  and  $\mathcal{P}$  respectively.



Constants	
1:	$m^{H_2} = 32 \parallel \text{len}(\text{LABEL}_2) \parallel \text{LABEL}_2 \parallel \text{len}(\mathbf{H}_2) \parallel \mathbf{H}_2 \parallel 1;$
2:	$m^{H_3} = 32 \parallel \text{len}(\text{LABEL}_3) \parallel \text{LABEL}_3 \parallel \text{len}(\mathbf{H}_3) \parallel \mathbf{H}_3 \parallel 1;$
3:	$m^{H_0} = 32 \parallel \text{len}(\text{LABEL}_0) \parallel \text{LABEL}_0 \parallel \text{len}(\mathbf{H}_0) \parallel \mathbf{H}_0 \parallel 1;$
4:	$m_f = 32 \parallel \text{len}(\text{LABEL}_f) \parallel \text{LABEL}_f \parallel 0 \parallel 1;$
5:	$m_k = 16 \parallel \text{len}(\text{LABEL}_k) \parallel \text{LABEL}_k \parallel 0 \parallel 1;$
6:	$m_{iv} = 12 \parallel \text{len}(\text{LABEL}_{iv}) \parallel \text{LABEL}_{iv} \parallel 0 \parallel 1;$
Protocol	
1:	$\mathcal{P} - h^{\text{HS}, \text{ipad}}, l_x = f(\text{IV}, 0, \text{HS} \oplus \text{ipad})$
2:	$\mathcal{P} - \text{SHTS}^{\text{in}}, _ = f(h^{\text{HS}, \text{ipad}}, l_x, m^{H_2})$ // disclosed to $\mathcal{V}$
3:	$\mathcal{P} - h^{\text{HS}, \text{opad}}, l = f(\text{IV}, 0, \text{HS} \oplus \text{opad})$ // disclosed to $\mathcal{V}$
4:	$\mathcal{V} - \text{SHTS}, _ = f(h^{\text{HS}, \text{opad}}, l, \text{SHTS}^{\text{in}})$
5:	$\mathcal{V} - \text{fk}_{\mathcal{S}}^{\text{in}}, _ = f(f(\text{IV}, 0, \text{SHTS} \oplus \text{ipad}), m_f)$
6:	$\mathcal{V} - \text{fk}_{\mathcal{S}}, _ = f(f(\text{IV}, 0, \text{SHTS} \oplus \text{opad}), \text{fk}_{\mathcal{S}}^{\text{in}})$
7:	$\mathcal{V} - \text{SF}^{\text{in}}, _ = f(f(\text{IV}, 0, \text{fk}_{\mathcal{S}} \oplus \text{ipad}), H_7 \parallel 1)$
8:	$\mathcal{V} - \text{SF}', _ = f(f(\text{IV}, 0, \text{fk}_{\mathcal{S}} \oplus \text{opad}), \text{SF}^{\text{in}})$
9:	$\mathcal{V} - \text{SF}' \stackrel{?}{=} \text{SF}$
<hr/>	
10:	$\mathcal{P} - \text{dHS}^{\text{in}}, l = f(h^{\text{HS}, \text{ipad}}, l_x, m^{H_0})$ // private input
11:	ZK - $\text{dHS}, _ = f(h^{\text{HS}, \text{opad}}, l, \text{dHS}^{\text{in}})$
12:	$\mathcal{P} - h^{\text{dHS}, \text{ipad}}, l = f(\text{IV}, 0, \text{dHS} \oplus \text{ipad})$ // disclosed to $\mathcal{V}$
13:	$\mathcal{V} - \text{MS}^{\text{in}}, _ = f(h^{\text{dHS}, \text{ipad}}, l, 0 \text{ bytes})$
14:	ZK - $\text{MS}, _ = f(f(\text{IV}, 0, \text{dHS} \oplus \text{opad}), \text{MS}^{\text{in}})$
15:	$\mathcal{P} - h^{\text{MS}, \text{ipad}}, l = f(\text{IV}, 0, \text{MS} \oplus \text{ipad})$ // disclosed to $\mathcal{V}$
16:	$\mathcal{V} - \text{SATS}^{\text{in}}, _ = f(h^{\text{MS}, \text{ipad}}, l, m^{H_3})$
17:	ZK - $\text{SATS}, _ = f(f(\text{IV}, 0, \text{MS} \oplus \text{opad}), \text{SATS}^{\text{in}})$
<hr/>	
18:	$\mathcal{P} - h^{\text{SATS}, \text{ipad}}, l = f(\text{IV}, 0, \text{SATS} \oplus \text{ipad})$ // disclosed to $\mathcal{V}$
19:	$\mathcal{V} - \text{SATK}^{\text{in}}, _ = f(h^{\text{SATS}, \text{ipad}}, l, m_k)$
20:	ZK - $h^{\text{SATS}, \text{opad}}, l = f(\text{IV}, 0, \text{SATS} \oplus \text{opad})$
21:	ZK - $\text{SATK}', _ = f(h^{\text{SATS}, \text{opad}}, l, \text{SATK}^{\text{in}})[:16]$
22:	ZK - $\text{Assert}(\text{SATK}' == \text{SATK})$

**Figure 10: Computations required to execute the optimized circuit for key derivation. The sub-circuit for proving key derivation is highlighted in gray. Computations executed by the prover and verifier are denoted by  $\mathcal{P}$  and  $\mathcal{V}$ .  $f$  is the one-way compression function SHACAL-2, which takes a 32 byte input as the first argument, a padding length helper as the second argument, and a 64 byte input as the third argument. The function returns the length  $l$  (required for padding) of the hashed input, and the output of SHACAL-2. The initialization vector of  $f$  computes as  $\text{IV} = H[0, \dots, 7]$ .**

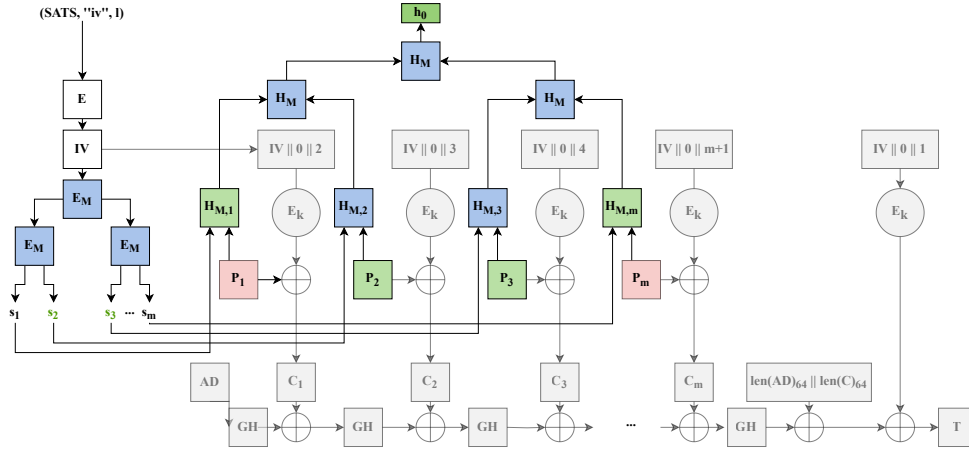
“ZK” refers to computations within the arithmetic circuit being proven by the general-purpose ZKP, highlighted in gray. We describe all computations as invocations of compression functions  $f$ . We denote values that result from the inner compression function in HMAC with the superscript  $x^{\text{in}}$ . Figure 10 shows that if the client discloses  $h^{\text{HS}, \text{opad}}$  and  $\text{SHTS}^{\text{in}}$  (line 2 and 3), the verifier can derive SHTS (line 4) out of circuit to verify SF. Note

that  $\text{dHS}^{\text{in}}$  is the only private input (line 10), which allows the in-circuit derivation of dHS (line 11). Deriving either one of the application traffic keys CATK or SATK requires 7 in-circuit invocations of the SHACAL-2 compression function  $f$ . Formally, the ZKP circuit formally proves that the witness tuple  $(x, w)$  is a satisfying assignment to the circuit describing the relation  $\mathcal{R}^{\text{KDC}}$ , where  $\mathcal{R}^{\text{KDC}}(x, w) = 1$  only if  $\text{SATK}' == \text{SATK}$ , for a given  $x = \{h^{\text{HS}, \text{opad}}, \text{dHS}^{\text{in}}, \text{MS}^{\text{in}}, \text{SATS}^{\text{in}}, \text{SATK}^{\text{in}}\}$  and  $w = \{\text{dHS}^{\text{in}}\}$ . Note that for completeness, the SATK is disclosed in line 22, which compromises privacy. Retaining privacy requires a circuit that proves correct key derivation, and successively proves correct encryption, including an assertion that a plaintext satisfies a predicate. The circuit used in ORIGO leverages the KDC circuit as a sub-circuit, without line 22 and without disclosing the SATK (cf. Section 5.1).

### 5.3 Client Side Evidence Generation

In this section, we outline how ORIGO can be extended beyond simple attestations to generate a commitment to the plaintext records transmitted in a TLS session. The rationale of evidence generation follows the paradigm of evidence generation as introduced by Ritzdorf *et. al* [62], with the important difference that evidence is generated by the *client* rather than the *server*. This preserves that no server-side changes are demanded to generate evidence of data provenance. In Figure 11 we depict the process of client side evidence generation for 4 plaintext chunks of size 16B. We realize the conversion of plaintext to hiding and binding commitments in two steps. First, the client expands the per-record nonce into a salt tree. Second, the client generates a salted Merkle Tree, where the tree leafs are instantiated as hash functions over the plaintext AEAD blocks salted with the per-block salts generated by the salt tree. We assume that a hash function salted with a sufficiently large salt therein provides a hiding and binding commitment. Importantly, the initial nonce is the initialization vector  $\text{IV}$ , which is derived from SATS in TLS 1.3 and directly binds the nonce utilized to generate salt values to the TLS handshake.

**Generating a Proof of Evidence.** To generate evidence, the client needs to convince the proxy verifier that the root hash of the salted merkle tree  $h_0$  was correctly derived with respect to the salt tree expansion and Merkle tree computation. Therefore, the client computes a ZKP, attesting the correct derivation of the root hash  $h_0$ , in addition to the proof of key derivation, authentication tag and encryption of records as described in Section 5.1. As in TLS 1.3, we instantiate the expansion function that is used in the salt tree with  $\text{HKDF.expand}$ , which is a common pseudorandom function for deriving keying material that is cryptographically secure, contextually bound, and properly domain-separated. However, due to the requirement of computing  $\text{HKDF.expand}$  inside a SNARK, and the omitted demand to base HMAC on SHA-256, the MiMC hash [12] can be used in  $\text{HKDF.expand}$  for HMAC. The output size of a MiMC hash is determined by the field that it operates over, as it is designed to work in a finite field setting. In our implementation, we instantiate a SNARK over BN254, which operates over a prime field of 254 bit. This is differing to SHA-256, which outputs a 256 Bit hash value. Therefore, each salt in the leafs of the salt tree is of size 127 bit. As the maximum number of blocks in TLS record is 1024, the maximum depth of the salt expansion tree is 10, which results



**Figure 11: Client side evidence generation for a single TLS record.** The example shows 4 blocks encrypted with AES-GCM (components of plain AES-GCM [1] are highlighted in light gray). The maximum number of blocks  $m$  is 1024, where each block has a size of 16B.  $E$  is the expansion function  $\text{HKDF.expand}$  instantiated with SHA-256, whereas  $E_M$  is  $\text{HKDF.expand}$  instantiated with MiMC. To disclose  $P_2, P_3$  to a third party, the client reveals all values highlighted in green. Additionally, the third party verifies the signature of the proxy verifier, which verifies that  $h_0$  is authentic.

in a maximum of 1023 invocations of  $\text{HKDF.expand}$ . When instantiated with MiMC, this is equivalent to 1022 HMAC evaluations with a pseudorandom key of 127 bit and a single HMAC evaluation with a pseudorandom key of 96 bit (the initial IV). As each HMAC internally composes two hash evaluations (cf. Section 3), this demands 2046 MiMC invocations over BN254.

The same optimizations for proof generation as in the key derivation circuit can be applied for the generation of the salt and merkle tree, such that the number of MiMC evaluations can be divided by two, as described in Section 5.2. Verification of hashes optimized for SNARKs remains a significant bottleneck in client-side evidence generation. To address this, alternative SNARK constructions have been proposed. Recent work by Belling *et al.*[17] demonstrates how hashing can be significantly accelerated by embedding the GKR[38] protocol within a Groth16 prover. Specifically, their approach uses a public-coin interactive protocol (GKR) to prove the base circuit and then recursively proves the GKR verification within Groth16. This method significantly reduces the number of constraints compared to proving a native hash function like MiMC directly in Groth16 - generating an in-circuit Merkle tree with 1024 leaves (depth 10) requires 52.2k constraints using MiMC, but only 6.2k constraints when using the recursive GKR construction (cf. Section 7.2).

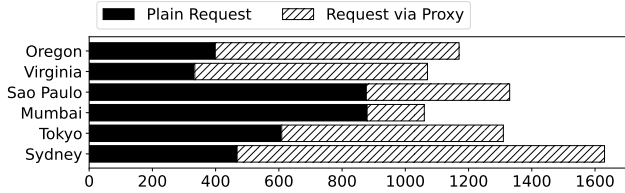
**Presenting Evidence to a Third Party.** The client submits  $h_0$  as an additional public input to the verifier, such that  $x := (p_R, \Gamma', \hat{R}, h_0)$  (cf. Figure 5). The private input remains  $w := (\text{SATK}, \theta_S, \text{dHS}^{\text{in}}, R)$ . The verifier derives the SATK, and additionally derives the salt and computes the salted merkle tree to check the validity of the provided Merkle Root. Once verified, the prover obtains a signature from the verifier which attests to the validity of the provided  $h_0$ . If the client intends to disclose a full block in plain, it can open the commitment to the plaintext by disclosing a number of hashes and salts that is logarithmic in the number of elements in the salt and merkle tree. An example of disclosing  $P_3, P_4$  to a third party in plain can be observed in Figure 11, where all disclosed values are highlighted in green.

Similarly, the client can prove an arbitrary statement on any non-disclosed block by proving  $\text{ZK-PoK}\{s_i, P_i : h_0 = \text{MT}(H_{M,i}(s_i, P_i))\}$  and disclosing all other outputs of  $H_{M,j}$  where  $j \neq i$ . Here, the statement or context to be attested by the proxy verifier does not need to be fixed. Instead, the verifier simply authenticates the evidence by the client, who can reveal a custom amount of information whenever demanded, even after the TLS session concluded.

**Further Optimization.** To optimize the performance of client side evidence generation, one may resort to an alternative commitment based on the counter block  $\text{CB}_i = E_k(\text{IV} || 0 || i + 1)$ . This is possible as the AES encryption of the counter values in GCM mode act as a pseudorandom function. As long as the counter values are unique for each block of plaintext under a given key, the resulting keystream will be pseudorandom, making the ciphertext appear random and ensuring that repeated plaintexts do not result in repeated ciphertexts. Hence, a client can utilize  $\text{CB}_i$  as a sufficiently random value to build the Merkle Tree obtain the hash  $h_0$ . This demands for additional storage of counter blocks by the client, in addition to the plaintext and ciphertext. After proving correct derivation of the counter block Merkle Tree, the client can prove an arbitrary statement on any non-disclosed block by computing  $\text{ZK-PoK}\{k : h_0 = \text{MT}(H_{M,i}(\text{CB}_i))\}$ , where  $k$  is SATK for the server response, or CATK for the request. This introduces a significant performance improvement, as it does not demand the generation of any salt values, rendering the salt tree expansion superfluous.

## 6 IMPLEMENTATION

We implement ORIGO in  $\sim 5k$  lines of golang code, out of which 2.7k lines of code are the implementation of SNARK circuits. This number excludes the adapted Golang TLS standard library. To facilitate parsing of HTTP requests and responses over TLS, as well as to transcribe the TLS handshake and extract relevant intermediate values, we customize the standard TLS package as provided in `go/crypto` [6]. We rely on *gnark* [18] for our implementation of



**Figure 12: Overhead of tunneling a TLS session through a proxy. The proxy is fixed in Munich. The execution time is in milliseconds. The TLS server is hosted by PayPal [7].**

general-purpose SNARKs. Gnark is a library that provides a programming interface to facilitate the development of ZKP applications. It relies on gnark-crypto [19], a library for arithmetic in finite fields and operations over common (pairing-friendly) elliptic curves. Gnark supports implementations of three SNARKs - Groth16, Plonk with the KZG polynomial commitment and Plonk with FRI as a commitment scheme and interactive proof. Our implementation includes circuits for proving SHA-256, HMAC, HMAC with MiMC, AES-128, AES in Galois Counter Mode, and the full TLS key derivation.

We open-source our implementation of ORIGO in a public repository [10]. ORIGO can either be run in a local environment with a mock server, or in a realistic setting with an externally hosted server. For a realistic setting, we integrated ORIGO with PayPal [7].

## 7 EVALUATION

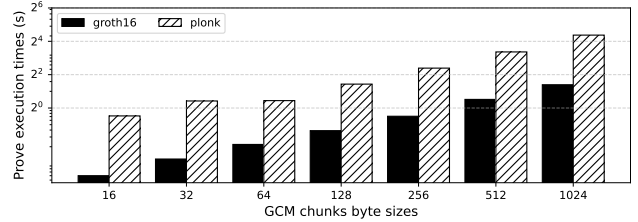
We evaluate ORIGO by evaluating (i) the overhead of ORIGO over the TLS baseline handshake and request execution, (ii) the individual performance of techniques introduced in § 5.2-§ 5.3, (iii) a global deployment and (iv) a holistic evaluation with related works.

**Experimental Setup.** Experiments to measure the circuit runtime are conducted on a MacBook Air configured with 24GB of RAM and the Apple M2 chip. In the global deployment, each client is executed on an AWS EC2 c5.4xlarge instance with 32GB of RAM and Intel Xeon Platinum 8124M CPU with 16 vCPUs. We leverage the insight of Ernstberger *et. al* [32], which states that circuits in gnark are 40% – 50% faster on compute optimized hardware. The proxy runs on an Intel server with 128GB of RAM and an Intel(R) Xeon(R) Gold 6130 CPU @ 2.10GHz. All results are averaged over 20 executions. In our experiments with adjusted bandwidth and throughput, we throttle bandwidth and throughput with tc.

**Preprocessing.** We do not evaluate preprocessing on the client and proxy side, which are out of circuit operations as described in Section 5.2. This is acceptable, as preprocessing is performed out of circuit and only requires a minimal compute and bandwidth.

### 7.1 Overhead over TLS in request execution

Figure 12 shows that the overhead introduced by the proxy setting is minimal with regard to the overall protocol runtime. We measure the latency averaged over 20 invocations with the real-world API provided by PayPal. We deploy the client in 6 distinct locations around the globe, and execute (i) a plain request to the API server as well as (ii) a handshake and request through the proxy, where the proxy save the relevant handshake transcript as well as the query and response transmitted by the client and server respectively. We



**Figure 13: Prover Time for AES-GCM. The x-axis depicts the plaintext byte size (in further experiments in this section, such as global experiments, we use 32 byte).**

find that the overhead introduced by the proxy ranges between 1.2× in Mumbai and 3.5× in Sydney. The overall end-to-end runtime for an HTTP request with ORIGO is 1s to 1.6s, which is well below application-specific TLS timeouts ( 10 – 15 seconds).

### 7.2 Performance of Circuit Implementations

Table 3 in the appendix provides microbenchmarks for all ORIGO circuits. We measure the size of the circuit in number of constraints, offline and online runtime, and offline and online communication. Note, that offline communication differs depending on the protocol setting and proof system. For Groth16, there is no separation of a randomized and deterministic setup procedure. Hence, the trusted setup would be executed by an external party, and the prover key and verifier are dispersed to the prover and verifier respectively. In Plonk, the randomized setup procedure can be executed separately (see e.g. [57] or [8] for a practical example). The SRS can be dispersed, whereas the prover and verifier derive the prover and verifier key locally in a deterministic setup. The complete end-to-end circuit for ORIGO counts 400.2k constraints. The bottleneck in generating a proof is AES-GCM, which requires 77.7k constraints and 0.23 seconds to prove the encryption of a 32 Byte plaintext (cf. Figure 13, which shows that computing the proof for 512 bytes takes > 1 second). In our evaluation, the AES-GCM circuit is optimized utilizing the logarithmic derivative lookups as introduced by Haboeck *et. al* [42]. Utilizing lookup arguments to prove AES as described by Orru *et. al* [58], or proving AES in the native field (e.g. [27]) can further lower the cost of AES. Currently, our circuit applies logarithmic derivative lookups as introduced by Haboeck [42]

**Client Side Evidence Generation.** To measure the performance of client side evidence generation, we evaluate the time it takes to (i) generate the salt tree and (ii) generate the salted merkle tree. In an unoptimized setting the HMAC circuit with the MiMC hash function takes up 52.2k constraints and 0.41s, whereas a comparable SHA-256 based HMAC evaluation takes up 184.9k constraints and 0.75s. Each expansion function requires a single HMAC evaluation for MiMC over BN254, which totals the prover time to ~ 419s - a significant overhead requiring further optimization.

With optimized evidence generation, only the merkle tree over the encrypted counter blocks needs to be computed in-circuit, in addition to the circuits described in Table 3. The generation of an in-circuit Merkle Tree with 1024 leaves of depth 10 takes 2.04s with

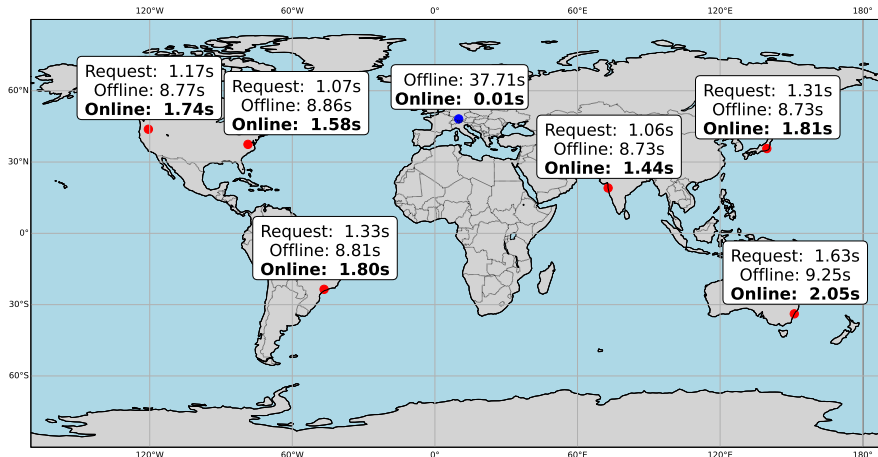


Figure 14: Evaluating globally dispersed clients requesting data from PayPal through a proxy verifier in Munich. Proxy “Offline” time includes the setup of  $pk$  &  $vk$ , and compiling the constraint system. Client “Offline” time is the compilation of the constraint system. Client “Online” time is the time to create a proof (cf. § 5), verifying it by the proxy, and receiving a response.

Table 2: DECO [74], XYWY23 [70], Janus [52] and this work under WAN conditions. For communication cost, values in brackets signify the communication cost for prover/verifier keys respectively. For runtime, we report online operations. Values in brackets are offline operations for the SNARK setup and circuit compilation. The proof system used for ORIGO is Groth16.

	Communication Cost			WAN (100Mbps, RTT=50ms)		
	256 B	1 KB	2 KB	256 B	1 KB	2 KB
DECO [74]	206 MB	345 MB	476 MB	24 s	36.3 s	51.6 s
XYWY23 [70]	15.2 MB	22.9 MB	33.3 MB	3.19 s	3.96 s	4.9 s
Janus [52]	249 MB	321 MB	409 MB	3.27 s	3.9 s	4.76 s
This Work (c5.4xlarge)	<b>2.4 kB</b> (111.9 MB/10.9 kB)	<b>2.4 kB</b> (111.9 MB/10.9 kB)	<b>2.4 kB</b> (111.9 MB/10.9 kB)	<b>1.50 s</b> (28.6 s/9.2 s)	<b>1.53 s</b> (28.6 s/9.2 s)	<b>1.51 s</b> (28.6 s/9.2 s)
This Work (m5.large)	”	”	”	<b>9.82 s</b> (216.37 s / 11.1 s)	<b>9.79 s</b> (214.23 s / 10.9 s)	<b>9.83 s</b> (217.52 s / 11.3 s)

a MiMC hash, and 0.12s when instantiated recursively with the GKR construction by Belling *et. al* [17].

### 7.3 Global Deployment

In Figure 14, we evaluate the execution time for globally distributed clients accessing a PayPal server [7] through a proxy verifier in Munich. All clients are running on AWS c5.4xlarge. The request sent aims to create an order for a product valued at a specific price. An order in PayPal is an agreement between a buyer and a seller, the integration with ORIGO allows the client to prove that a product has been purchased via PayPal payment for a certain price. This API has a request size of 873 byte and a response size of 1429 byte. We deploy the client in 6 different cities and observe that the online execution time of our protocol ranges from 1.44s to 2.05s. Xie *et. al* [70] conduct a similar experiment by deploying the client in 18 different cities, reporting an online runtime between 0.3s and 9.6s for a similarly sized request / response pair, which

indicates that ORIGO enjoys a more stable online execution time, independent of the geographical location of the TLS session verifier. This is primarily caused by the number of online communication rounds, which is 1 for ORIGO, and 31 in the case of Xie *et. al*.

### 7.4 Comparison with Related Works

We compare the performance of ORIGO to DECO [74], the concurrent works of Xie *et. al* [70], and Lauinger *et. al* [52] All works rely on a three-party handshake, which replaces the traditional TLS handshake, and allow for a stronger network adversary as opposed ORIGO. For DECO, we rely on the performance evaluation conducted by Xie *et. al* [70], which instantiate DECO with the state-of-the-art protocol for malicious 2PC, Ferret [71]. The evaluation of DECO only includes the time needed in malicious 2PC, which includes computing the TLS session keys and 4 pairs of AES-GCM ciphertext/tags [70]. Measurements are generated using

AWS c5.4xlarge and m5.large (Note, that m5.large is less powerful than e.g. a mobile phone). We summarize results in Table 2.

As expected, we observe that our protocol reports constant communication and execution time, independent of the size of the request and response. For communication cost, we can observe that the online communication in ORIGO is smaller than in both DECO and the work of Xie *et. al* [70]. For a query size of 587 bytes and response size of 894 bytes, Xie *et. al* report a communication cost of 17.6 MB offline and 0.9 MB online, whereas ORIGO requires a constant 2.4 kB, which is 375× less. For a request and response of 256 Byte, our solution completes in 1.50 seconds, which is significantly faster than DECO’s 24 seconds. This presents a  $\sim 16\times$  speedup. Note, that the online execution time highly depends on the size of the proven statement. In our example, the statement fits into two chunks of 128 bits. If in-circuit encryption of multiple chunks is required, the online time of ORIGO scales as described in Figure 13.

## 8 DISCUSSION AND FUTURE WORK

This paper introduces ORIGO, a TLS oracle with constant communication from general purpose ZKPs. The main challenge was to minimize the in-circuit operations to achieve practical prover runtimes (§ 5.2), and ensure mitigation of an attack on integrity (§ 4). Another challenge was achieving flexibility for clients, such that they can create arbitrary proofs post-session (§ 5.3). We show that the total online communication of ORIGO is 2.4kB. We evaluate ORIGO with globally dispersed clients, showing that the online runtime for proof generation and presentation ranges between 1.44 seconds and 2.05 seconds. The performance benefit is primarily obtained by the proxy setting, and the assumption that MITM attacks are infeasible, which allows us to use a SNARK in a public coin protocol. This is interesting, as it allows outsourcing the prover to an untrusted server without forfeiting confidentiality [37].

There are three limitations that demand further exploration. First, assuming that MITM is impossible demands for additional mitigation on other OSI layers. To perform the attack, a client would need to manipulate the DNS resolution at the proxy, or intercept the traffic between the proxy and the server. This could be attempted through DNS hijacking or spoofing, which requires compromising the DNS resolver by the proxy. With DNSSEC, the proxy can ensure that the domain name resolution is secure. To ensure that the client doesn’t supply a malicious domain name in its request, the proxy can pre-specify a set of whitelisted domain names. Certain applications, may deem ORIGO sufficient for their application, as practical attacks are non-trivial. We further discuss the feasibility of an MITM attack for traffic injection in Appendix D. Second, the size of the SRS is large ( $\sim 100 - 300\text{MB}$  for Groth16/Plonk) and requires offline communication. If each proxy verifier creates its own SRS, clients need to download large files before a session. A public SRS may be a suitable alternative to lessen this issue [8].

## REFERENCES

- [1] 2007. Recommendation for Block Cipher Modes of Operation: Galois/Counter Mode (GCM) and GMAC. <https://csrc.nist.gov/publications/detail/sp/800-38d/final>.
- [2] 2014. TLS Notary - Initial Version. [https://old.tlsnotary.org/how\\_it\\_works](https://old.tlsnotary.org/how_it_works).
- [3] 2014. tlsnotary - cryptographic proof of fiat transfer for p2p exchanges. <https://bitcointalk.org/index.php?topic=173220.0>.
- [4] 2023. Ceramic Network. <https://ceramic.network/>.
- [5] 2023. Bitcoin Passport. <https://passport.gitcoin.co/>.
- [6] 2023. Golang Crypto Package. <https://pkg.go.dev/crypto>.
- [7] 2023. PayPal API specification. <https://developer.paypal.com/api/rest/>.
- [8] 2023. PROTO-DANKSHARDING AND THE CEREMONY. <https://ceremony.ethereum.org/>.
- [9] 2023. ThousandEyes - Best BGP Route Network Monitoring Solution. <https://www.thousandeyes.com/solutions/bgp-and-route-monitoring>.
- [10] 2023. TLS Oracle Demo. <https://github.com/opex-research/tls-oracle-demo>.
- [11] Behzad Abdolmaleki, Noemi Glaeser, Sebastian Ramacher, and Daniel Slamanig. 2023. Universally Composable NIZKs: Circuit-Succinct, Non-Malleable and CRS-Updateable. *Cryptology ePrint Archive* (2023).
- [12] Martin Albrecht, Lorenzo Grassi, Christian Rechberger, Arnab Roy, and Tyge Tiessen. 2016. MiMC: Efficient encryption and cryptographic hashing with minimal multiplicative complexity. In *International Conference on the Theory and Application of Cryptology and Information Security*. Springer, 191–219.
- [13] Oussama Amine, Karim Baghery, Zaira Pindado, and Carla Ràfols. 2023. Simulation extractable versions of Groth’s zk-SNARK revisited. *International Journal of Information Security* (2023), 1–15.
- [14] Matilda Backendal, Mihir Bellare, Felix Günther, and Matteo Scarlata. 2023. When messages are keys: Is HMAC a dual-PRF?. In *Annual International Cryptology Conference*. Springer, 661–693.
- [15] Karim Baghery, Markulf Kohlweiss, Janno Siim, and Mikhail Volkov. 2021. Another look at extraction and randomization of groth’s zk-SNARK. In *Financial Cryptography and Data Security: 25th International Conference, FC 2021, Virtual Event, March 1–5, 2021, Revised Selected Papers, Part I 25*. Springer, 457–475.
- [16] Mihir Bellare and Oded Goldreich. 1992. On defining proofs of knowledge. In *Annual International Cryptology Conference*. Springer, 390–420.
- [17] Alexandre Belling, Azam Soleimanian, and Olivier Bégassat. 2022. Recursion over Public-Coin Interactive Proof Systems; Faster Hash Verification. *Cryptology ePrint Archive* (2022).
- [18] Gautam Botrel, Thomas Piellard, Youssef El Housni, Ivo Kubjas, and Arya Tabaie. 2023. *ConsensSys/gnark: v0.8.0*. <https://doi.org/10.5281/zenodo.5819104>
- [19] Gautam Botrel, Thomas Piellard, Youssef El Housni, Arya Tabaie, Gus Gutoski, and Ivo Kubjas. 2023. *ConsensSys/gnark-crypto: v0.11.2*. <https://doi.org/10.5281/zenodo.5815453>
- [20] Ran Canetti. 2001. Universally composable security: A new paradigm for cryptographic protocols. In *Proceedings 42nd IEEE Symposium on Foundations of Computer Science*. IEEE, 136–145.
- [21] Ran Canetti and Marc Fischlin. 2001. Universally composable commitments. In *Advances in Cryptology—CRYPTO 2001: 21st Annual International Cryptology Conference, Santa Barbara, California, USA, August 19–23, 2001 Proceedings 21*. Springer, 19–40.
- [22] Srđjan Capkun, Ercan Ozturk, Gene Tsudik, and Karl Wüst. 2021. ROSEN: RObst and SElective Non-repudiation (for TLS). In *Proceedings of the 2021 on Cloud Computing Security Workshop*. 97–109.
- [23] Sofia Celi, Alex Davidson, Hamed Haddadi, Gonçalo Pestana, and Joe Rowell. 2023. Distefano: Decentralized infrastructure for sharing trusted encrypted facts and nothing more. *Cryptology ePrint Archive* (2023).
- [24] Kwan Yin Chan, Handong Cui, and Tsz Hon Yuen. 2023. DIDO: Data Provenance from Restricted TLS 1.3 Websites. *Cryptology ePrint Archive* (2023).
- [25] Alessandro Chiesa, Dev Ojha, and Nicholas Spooner. 2020. Fractal: Post-quantum and transparent recursive proofs from holography. In *Advances in Cryptology—EUROCRYPT 2020: 39th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Zagreb, Croatia, May 10–14, 2020, Proceedings, Part I 39*. Springer, 769–793.
- [26] World Wide Web Consortium et al. 2019. Verifiable credentials data model 1.0: Expressing verifiable information on the web. [https://www.w3.org/TR/vc-data-model/#\\_core-data-model](https://www.w3.org/TR/vc-data-model/#_core-data-model) (2019).
- [27] Benjamin E Diamond and Jim Posen. 2024. Polylogarithmic Proofs for Multitowers over Binary Towers. *Cryptology ePrint Archive* (2024).
- [28] Tim Dierks and Eric Rescorla. 2008. The transport layer security (TLS) protocol version 1.2. (2008).
- [29] Yevgeniy Dodis, Paul Grubbs, Thomas Ristenpart, and Joanne Woodage. 2018. Fast message franking: From invisible salamanders to encryption. In *Annual International Cryptology Conference*. Springer, 155–186.
- [30] Benjamin Dowling, Marc Fischlin, Felix Günther, and Douglas Stebila. 2015. A cryptographic analysis of the TLS 1.3 handshake protocol candidates. In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*. 1197–1210.
- [31] Benjamin Dowling, Marc Fischlin, Felix Günther, and Douglas Stebila. 2021. A cryptographic analysis of the TLS 1.3 handshake protocol. *Journal of Cryptology* 34, 4 (2021), 1–69.
- [32] Jens Ernstberger, Stefanos Chaliasos, George Kadianakis, Sebastian Steinhorst, Philipp Jovanovic, Arthur Gervais, Benjamin Livshits, and Michele Orrù. 2023. zk-Bench: A Toolset for Comparative Evaluation and Performance Benchmarking of SNARKs. *Cryptology ePrint Archive* (2023).
- [33] Jens Ernstberger, Jan Lauinger, Fatima Elsheimi, Liyi Zhou, Sebastian Steinhorst, Ran Canetti, Andrew Miller, Arthur Gervais, and Dawn Song. 2023. SoK: Data

- Sovereignty. In *2023 IEEE 8th European Symposium on Security and Privacy (EuroS&P)*, 122–143. <https://doi.org/10.1109/EuroSP57164.2023.00017>
- [34] Marc Fischlin, Felix Günther, Benedikt Schmidt, and Bogdan Warinschi. 2016. Key confirmation in key exchange: A formal treatment and implications for TLS 1.3. In *2016 IEEE Symposium on Security and Privacy (SP)*. IEEE, 452–469.
- [35] Ariel Gabizon, Zachary J Williamson, and Oana Ciobotaru. 2019. Plonk: Permutations over lagrange-bases for ocumenical noninteractive arguments of knowledge. *Cryptology ePrint Archive* (2019).
- [36] Chaya Ganesh, Yashvanth Kondi, Claudio Orlandi, Mahak Panchoi, Akira Takahashi, and Daniel Tschudi. 2023. Witness-succinct universally-composable snarks. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 315–346.
- [37] Sanjam Garg, Aarushi Goel, Abhishek Jain, Guru-Vamsi Policharla, and Sruthi Sekar. 2023. zkSaaS: Zero-Knowledge SNARKs as a Service. *Cryptology ePrint Archive* (2023).
- [38] Shafi Goldwasser, Yael Tauman Kalai, and Guy N Rothblum. 2015. Delegating computation: interactive proofs for muggles. *Journal of the ACM (JACM)* 62, 4 (2015), 1–64.
- [39] Jens Groth. 2006. Simulation-sound NIZK proofs for a practical language and constant size group signatures. In *Advances in Cryptology—ASIACRYPT 2006: 12th International Conference on the Theory and Application of Cryptology and Information Security, Shanghai, China, December 3–7, 2006. Proceedings 12*. Springer, 444–459.
- [40] Paul Grubbs, Arasu Arun, Ye Zhang, Joseph Bonneau, and Michael Walfish. 2022. {Zero-Knowledge} Middleboxes. In *31st USENIX Security Symposium (USENIX Security 22)*, 4255–4272.
- [41] Paul Grubbs, Jiahui Lu, and Thomas Ristenpart. 2017. Message franking via committing authenticated encryption. In *Annual International Cryptology Conference*. Springer, 66–97.
- [42] Ulrich Haböck. 2022. Multivariate lookups based on logarithmic derivatives. *Cryptology ePrint Archive* (2022).
- [43] Ralph Holz, Johanna Amann, Abbas Razaghpanah, and Narseo Vallina-Rodriguez. 2019. The era of TLS 1.3: Measuring deployment and use with active and passive methods. *arXiv preprint arXiv:1907.12762* (2019).
- [44] Rawane Issa, Nicolas Alhaddad, and Mayank Varia. 2022. Hecate: Abuse reporting in secure messengers with sealed sender. In *31st USENIX Security Symposium (USENIX Security 22)*, 2335–2352.
- [45] Peng Jiang, Baoqi Qiu, and Liehuang Zhu. 2022. Report When Malicious: Deniable and Accountable Searchable Message-Moderation System. *IEEE Transactions on Information Forensics and Security* 17 (2022), 1597–1609.
- [46] Daniel Kang, Tatsunori Hashimoto, Ion Stoica, and Yi Sun. 2022. ZK-IMG: Attested Images via Zero-Knowledge Proofs to Fight Disinformation. *arXiv preprint arXiv:2211.04775* (2022).
- [47] Eleftherios Kokoris Kogias, Enis Ceyhan Alp, Linus Gasser, Philipp Svetolik Jovanovic, Ewa Syta, and Bryan Alexander Ford. 2021. Calypso: Private data management for decentralized ledgers. *Proceedings of the VLDB Endowment* 14, CONF (2021), 586–599.
- [48] Ahmed Kosba, Zhichao Zhao, Andrew Miller, Yi Qian, Hubert Chan, Charalampos Papamanthou, Rafael Pass, Abhi Shelat, and Elaine Shi. 2015. Coco: a framework for building composable zero-knowledge proofs. *Cryptology ePrint Archive, Report 2015/1093* (2015).
- [49] Hugo Krawczyk, Mihir Bellare, and Ran Canetti. 1997. HMAC: Keyed-hashing for message authentication.
- [50] Hugo Krawczyk and Pasi Eronen. 2010. *HMAC-based extract-and-expand key derivation function (HKDF)*. Technical Report. RFC 5869, May.
- [51] Hugo Krawczyk and Hoeteck Wee. 2016. The OPTLS protocol and TLS 1.3. In *2016 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 81–96.
- [52] Jan Lauinger, Jens Ernstberger, Andreas Finkenzeller, and Sebastian Steinhorst. 2025. Janus: Fast Privacy-Preserving Data Provenance For TLS. *Proceedings on Privacy Enhancing Technologies* (2025).
- [53] Jiqiang Lu and Jongsung Kim. 2008. Attacking 44 rounds of the SHACAL-2 block cipher using related-key rectangle cryptanalysis. *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences* 91, 9 (2008), 2588–2596.
- [54] Ning Luo, Chenkai Weng, Jaspal Singh, Gefei Tan, Ruzica Piskac, and Mariana Raykova. 2023. Privacy-Preserving Regular Expression Matching using Nondeterministic Finite Automata. *Cryptology ePrint Archive* (2023).
- [55] Deepak Maram, Harjasleen Malvai, Fan Zhang, Nerla Jean-Louis, Alexander Frolow, Tyler Kell, Tyrone Lobban, Christine Moy, Ari Juels, and Andrew Miller. 2021. Candid: Can-do decentralized identity with legacy compatibility, sybil-resistance, and accountability. In *2021 IEEE Symposium on Security and Privacy (SP)*. IEEE, 1348–1366.
- [56] Sanketh Menda, Julia Len, Paul Grubbs, and Thomas Ristenpart. 2023. Context Discovery and Commitment Attacks: How to Break CCM, EAX, SIV, and More. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 379–407.
- [57] Valeria Nikolaenko, Sam Ragsdale, Joseph Bonneau, and Dan Boneh. 2022. Powers-of-tau to the people: Decentralizing setup ceremonies. *Cryptology ePrint Archive* (2022).
- [58] Michele Orrù, George Kadianakis, Mary Maller, and Greg Zaverucha. 2024. Beyond the circuit: How to Minimize Foreign Arithmetic in ZKP Circuits. *Cryptology ePrint Archive* (2024).
- [59] Deevashwer Rathee, Guru Vamsi Policharla, Tiancheng Xie, Ryan Cottone, and Dawn Song. 2022. ZEBRA: Anonymous Credentials with Practical On-chain Verification and Applications to KYC in DeFi. *Cryptology ePrint Archive, Paper 2022/1286*. <https://eprint.iacr.org/2022/1286> <https://eprint.iacr.org/2022/1286>.
- [60] Michael Raymond, Gillian Evers, Jan Ponti, Diya Krishnan, and Xiang Fu. 2023. Efficient Zero Knowledge for Regular Language. *Cryptology ePrint Archive* (2023).
- [61] Eric Rescorla. 2018. *The transport layer security (TLS) protocol version 1.3*. Technical Report.
- [62] Hubert Ritzdorf, Karl Wüst, Arthur Gervais, Guillaume Felley, and Srdjan Capkun. 2017. TLS-N: Non-repudiation over TLS enabling-ubiquitous content signing for disintermediation. *Cryptology ePrint Archive* (2017).
- [63] Martin Schläffer. 2011. *Cryptanalysis of AES-based hash functions*. na.
- [64] Chris McMahon Stone, Tom Chothia, and Flavio D Garcia. 2017. Spinner: Semi-automatic detection of pinning without hostname verification. In *Proceedings of the 33rd annual computer security applications conference*, 176–188.
- [65] Nirvan Tyagi, Paul Grubbs, Julia Len, Ian Miers, and Thomas Ristenpart. 2019. Asymmetric message franking: Content moderation for metadata-private end-to-end encryption. In *Annual International Cryptology Conference*. Springer, 222–250.
- [66] Yinxi Wan, Kuai Xu, Guoliang Xue, and Feng Wang. 2020. Iotargos: A multi-layer security monitoring system for internet-of-things in smart homes. In *IEEE INFOCOM 2020-IEEE Conference on Computer Communications*. IEEE, 874–883.
- [67] Liang Wang, Gilad Asharov, Rafael Pass, Thomas Ristenpart, and Abhi Shelat. 2019. Blind certificate authorities. In *2019 IEEE Symposium on Security and Privacy (SP)*. IEEE, 1015–1032.
- [68] Chenkai Weng, Kang Yang, Jonathan Katz, and Xiao Wang. 2021. Wolverine: fast, scalable, and communication-efficient zero-knowledge proofs for boolean and arithmetic circuits. In *2021 IEEE Symposium on Security and Privacy (SP)*. IEEE, 1074–1091.
- [69] Chenkai Weng, Kang Yang, Zhaomin Yang, Xiang Xie, and Xiao Wang. 2022. AntMan: interactive zero-knowledge proofs with sublinear communication. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, 2901–2914.
- [70] Xiang Xie, Kang Yang, Xiao Wang, and Yu Yu. 2023. Lightweight Authentication of Web Data via Garble-Then-Prove. *Cryptology ePrint Archive* (2023).
- [71] Kang Yang, Chenkai Weng, Xiao Lan, Jiang Zhang, and Xiao Wang. 2020. Ferret: Fast extension for correlated OT with small communication. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, 1607–1626.
- [72] Collin Zhang, Zachary DeStefano, Arasu Arun, Joseph Bonneau, Paul Grubbs, and Michael Walfish. 2023. Zombie: Middleboxes that Don’t Snoop. *Cryptology ePrint Archive* (2023).
- [73] Fan Zhang, Ethan Cecchetti, Kyle Croman, Ari Juels, and Elaine Shi. 2016. Town crier: An authenticated data feed for smart contracts. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, 270–282.
- [74] Fan Zhang, Deepak Maram, Harjasleen Malvai, Steven Goldfeder, and Ari Juels. 2020. DECO: Liberating web data using decentralized oracles for TLS. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, 1919–1938.
- [75] Jian Zhang, Jennifer Rexford, and Joan Feigenbaum. 2005. Learning-based anomaly detection in BGP updates. In *Proceedings of the 2005 ACM SIGCOMM workshop on Mining network data*, 219–220.

## A APPLICATION TO GITCOIN PASSPORT

To demonstrate the practical use of our construction, we integrated ORIGO with Gitcoin Passport [5]. With Gitcoin Passport, users collect *stamps* from authenticators such as Facebook, Twitter or Github. By collecting stamps, a user can prove *sybil-resistance*. The more stamps a user collects, the higher the trust score, which indicates how difficult it is to forge an identity. To issue a stamp, a user has to authorize the Gitcoin Passport application to write to its personal Ceramic datastream<sup>1</sup> (“Passport”). The user connects to the Gitcoin Passport app, which communicates with the Gitcoin Identity and Access Management (IAM) server, to verify that an account with one of the previously mentioned platforms exists. The

<sup>1</sup>A data storage network for Web 3.0 applications which identifies data streams by a *Decentralized Identifier (DID)*.

**Table 3: Benchmarks for SNARK circuits in ORIGO. We assume the context spans 32 Byte, i.e. two TLS record blocks need to be encrypted in-circuit. We separate offline and online runtime and communication. CCS is the compiled constraint system, SRS is the reference string established in the randomized setup of the SNARK, PK and VK are the prover key and verifier key which are derived in a deterministic setup, PW is the public witness, and proof  $\pi$ .**

ZKP Circuit	# Constraints ( $\times 10^3$ )	Execution Offline		Execution Online		Communication Offline			Communication Online	
		Compile (s)	Setup (s)	Prove (s)	Verify (ms)	CCS, SRS (MB)	PK (MB)	VK (KB)	$\pi$ (KB)	PW (B)
<b>Groth16 backend</b>										
KDC	322.4	2.0	12.43	0.58	1.0	9.03	47.8	4.93	0.128	128
Tag	37.6	0.175	2.13	0.116	2.0	1.31	7.05	2.54	0.128	44
Record	40.1	0.187	2.3	0.117	2.0	1.77	7.35	2.19	0.128	72
<b>Plonk backend</b>										
KDC	564.56	0.454	1.98	9.52	1.8	8.58, 33.5	67.1	34.3	0.552	128
Tag	67.29	0.055	0.301	1.36	1.6	1.03, 4.2	8.42	34.3	0.552	44
Record	72.2	0.059	0.307	1.33	1.6	1.1, 4.2	8.4	34.3	0.552	72

Gitcoin passport application, in collaboration with the Gitcoin IAM server, acts as centralized issuer of a verifiable credential (“stamps”) as specified by the standardization of the W3C [26]. Once a stamp is issued, anyone can verify its validity. Note that presently, all stamps issued through Gitcoin Passport rely on publicly available data.

With ORIGO, users can create stamps from their private data, without relying on a centralized entity which is trusted for confidentiality. We integrated ORIGO with Gitcoin Passport to create a stamp from a successful PayPal payment through the PayPal API. To initiate the session, the client signs up with the PayPal API to generate an access token. The client executes  $\Pi_{ORIGO}$  with the proxy as a verifier, the PayPal API acts as the server and Gitcoin/Ceramic acts as an external third party. After obtaining a response to the HTTPS GET request, the client proves that the response submitted via TLS contains information about a payment, that is greater than a certain threshold, to the proxy verifier. Upon verification of message integrity and the correctness of values, the proxy verifier creates a signed verifiable credential<sup>2</sup> and sends it to the client. We leverage an adapted version of the writer SDK of Gitcoin Passport<sup>3</sup>, such that the client can create and write to its personal Ceramic Tile Stream [4]. The verifiable credential is then included as a stamp under the previously created DID. As such, our application allows for stamps from private data, where the user does not give up control over its “passport” to another entity. By issuing verifiable credentials from private data, the proxy verifier can act as a complementary issuer to the Gitcoin Passport application, without being trusted for confidentiality.

## B PRIVACY RIGHTS & REGULATIONS

In our practical implementation, each session in ORIGO is governed by a *policy*. The purpose of a policy for data feeds for confidential data is twofold — to provide a set of constraints that enables the verifier to validate the correctness of statements of the client, and

to provide an interface that ensures the validity of data sources providing qualitatively equivalent data (e.g., age verification through either the DMV or governmental API endpoints). In our current implementation, the policy specifies applicable data sources, response formats and value constraints to be proven by the client.

Further, a policy can also implicitly or explicitly guarantee data accuracy. To address liability concerns, the policy could include terms that specify the extent of the data provider’s liability for data accuracy. Moreover, including a hash of a Terms of Service (ToS) or Disclaimer within the policy can serve as an immutable reference to the terms under which the data is provided and attested.

Thus, ORIGO can enforce terms under which data is attested, and deny attestations that violate these terms. We expect that a set of policies will provide clients with a choice of data providers and endpoints, from which attested data can be obtained — similar to how passport stamps can be obtained in Gitcoin Passport [5].

## C FEASIBILITY OF THE ATTACK ON TLS 1.3

The attack presented in Section 4 challenges the integrity and uniqueness aspects of AES-GCM, aiming to compromise its security by finding collisions in encryption and authentication. The attack applies in a setting, where the adversary aims to convince an entity about an underlying plaintext, whereas the entity verifying the statement already holds the ciphertext encrypted with AES-GCM. The attack leverages that in TLS 1.3, the *IV* is *not* explicitly included in the transmitted messages, but rather derived from the handshake transcript. We begin to outline the trivial attack on an arbitrary plaintext, and successively discuss the feasibility of the attack given an adversarial chosen, structured plaintext.

First, the adversary  $\mathcal{A}_{GCM}$  a key  $K'$ . Successively,  $\mathcal{A}_{GCM}$  computes the evaluation of the binary polynomial at  $H$ , by computing

$$X \leftarrow \sum_{i=1}^m C_i H^{m-i}$$

where  $H \leftarrow E_K\{0^{128}\}$  is the encryption of the zero vector. Here,  $\mathcal{A}_{GCM}$  is unable to choose  $X$ , as each  $C_i$  is fixed (i.e. known by the verifying entity). Now,  $\mathcal{A}_{GCM}$  chooses  $IV'$  such that  $T = T'$ , i.e. the

<sup>2</sup>We integrate a JavaScript SDK for W3C Verifiable Credentials <https://github.com/digitalbazaar/vc-js>

<sup>3</sup><https://github.com/gitcoinco/passport-sdk/tree/main/packages/writer>

**Table 4: Acronyms for messages, secrets, transcript hashes and labels used in the TLS 1.3 handshake. A detailed description can be found in RFC 8446 [61? ].**

Acronym	Definition
HS	Handshake Secret
SHTS	Server Handshake Traffic Secret
CHTS	Client Handshake Traffic Secret
dHS	derived Handshake Secret
MS	Master Secret
SATS	Server Application Traffic Secret
CATS	Client Application Traffic Secret
H <sub>0</sub>	$H(“”)$
H <sub>1</sub>	$H(\text{CLIENTHELLO})$
H <sub>2</sub>	$H(\text{CLIENTHELLO}  \text{SERVERHELLO})$
H <sub>3</sub>	$H(\text{CLIENTHELLO}  \dots  \{\text{SERVERFINISHED}\})$
H <sub>4</sub>	$H(\text{CLIENTHELLO}  \dots  \{\text{CLIENTFINISHED}\})$
H <sub>6</sub>	$H(\text{CLIENTHELLO}  \dots  \{\text{SERVERCERT}\})$
H <sub>7</sub>	$H(\text{CLIENTHELLO}  \dots  \{\text{SERVERCERTVIFY}\})$
H <sub>9</sub>	$H(\text{CLIENTHELLO}  \dots  \{\text{CLIENTCERTVIFY}\})$
LABEL <sub>0</sub>	“ext binder” / “res binder”
LABEL <sub>1</sub>	“c e traffic”
LABEL <sub>2</sub>	“e exp master”
LABEL <sub>3</sub>	“derived”
LABEL <sub>4</sub>	“c hs traffic”
LABEL <sub>5</sub>	“s hs traffic”
LABEL <sub>6</sub>	“finished”
LABEL <sub>7</sub>	“c ap traffic”
LABEL <sub>8</sub>	“s ap traffic”
LABEL <sub>9</sub>	“exp master”
LABEL <sub>10</sub>	“res master”
LABEL <sub>11</sub>	“TLS 1.3, server CertificateVerify”
LABEL <sub>f</sub>	“tls13 finished”
LABEL <sub>k</sub>	“tls13 key”
LABEL <sub>iv</sub>	“tls13 iv”

newly computed tag is equivalent to the one known by the verifying entity for a differing, adversarial chosen initialization vector.  $T'$  is computed as  $T' \leftarrow E_{K'}\{IV'\} \oplus X$ . Now that an adversarial tag for a differing  $IV$  is computed, it remains to show that the adversary can derive a colluding ciphertext with a differing plaintext, given  $IV'$  and  $K'$ . To do so, the adversary computes  $R' \leftarrow E_{K'}\{IV' || ctr\}$  and chooses an arbitrary plaintext  $P$  such that  $C \leftarrow R' \oplus P'$  is satisfied. Given an arbitrary plaintext, this attack is hence trivial.

For a chosen plaintext,  $\mathcal{A}_{GCM}$  needs to find a second preimage attack where they generate a different set of inputs  $(K', IV', P')$  that produce the same encrypted output and authentication tag as the original  $(K, IV, P)$ . While an adversary can freely choose a new key and IV, ensuring these selections result in the same tag and ciphertext for a different plaintext requires reversing the GHASH and AES. For  $\mathcal{A}_{GCM}$  to revert AES-GCM, the expected number of attempts is in the order of  $2^{128}$ , which is the size of the key space. Previous work on cryptanalytic techniques demonstrate

a rebound attack that can, under certain conditions, significantly lower the complexity of finding a set of inputs that satisfy a given cryptographic condition than what a brute-force approach would suggest [63]. For AES, this lowers the complexity to about  $2^{48}$  for constructing valid pairs that follow a specific differential trail.

Concurrently, Menda *et. al* [56] introduce the generalized notion of *context discoverability* (CDY), where the adversary aims to discover the decryption context  $(K, N, A)$  for a fixed ciphertext  $C$ . They introduce the formal security notion  $CDY_n^*$ , which lets the adversary choose the nonce and fixes the key and associated data. They show that GCM does *not* fulfill  $CDY_n^*$  in an attack similar to ours in Section 4, and lower bound the adversarial advantage to  $2^{-32}$ . The property that we define as *strong receiver binding* would be equivalent to  $CDY_{(n,k)}^*$ , such that the associated data is fixed and nonce and key can be chosen arbitrarily. It holds in TLS 1.2, as an adversarial client is not able to arbitrarily choose the nonce  $n$ .

## D ON THE POSSIBILITY OF MITM ATTACKS

Assuming that MITM by a malicious client is not possible demands for additional mitigation on another OSI layer (without any server-side changes). For example, a proxy can employ IP Address Monitoring [66] and Secure BGP Practices on the network layer [9, 75] or Certificate Pinning on the transport layer to ensure correct certificates [64]. Further, external verification of the server’s public key, separate to the proxy session, can ascertain additional security.

Similar assumptions, which assume additional care to account for traffic injection through e.g. BGP attacks, have been acknowledged and embraced in similar settings that require a proxy [67, 74]. The reason for this is that BGP attacks are (i) are difficult to mount in practice, and (ii) various detection techniques have been proposed to detect BGP attacks [9, 75]. To additionally enhance the security of the approach introduced in this work, one may rely on an N-out-of-N multisignature with globally dispersed proxies, or a similar approach that relies on N-version programming, following the initial specification outlined in this work. However, we acknowledge that additional MiTM mitigation mechanisms introduce an additional overhead to our protocol when compared with approaches that do not face this issue [23, 52, 70, 74], which our evaluation does not account for. We leave further exploration of extending the security guarantees in the proxy mode as future work.

## E FORMAL SECURITY ANALYSIS

In this section, we analyze and prove the security of  $\Pi_{ORIGO}$  in the UC paradigm by Cannetti *et. al* [20]. In short, the UC model employs a real-ideal world paradigm, where security follows from indistinguishability between an execution of the real protocol in the real world and execution of an ideal version in the ideal world. The reasoning is that the ideal functionality specifies the ideal security properties and fulfills them by definition. If the real protocol execution is indistinguishable from the execution with the ideal world, it follows that it yields exactly the same security properties.

Formally, a protocol  $\Pi$  UC-realizes and ideal functionality  $\mathcal{F}$  if  $\Pi$  emulates  $I_{\mathcal{F}}$ , the ideal protocol of  $\mathcal{F}$ . A protocol  $\Pi$  emulates protocol  $I_{\mathcal{F}}$  if for any polynomial time (PT) adversary  $\mathcal{A}$ , there exists a PT adversary  $\text{Sim}$  such that for all PT environments  $\mathcal{Z}$  that



```


$\mathcal{F}_{DP}(\lambda, \mathcal{R}, \mathcal{S}, \mathcal{V}, C)$


01: \ \ Initialize
02: Init: Set STORAGE  $\equiv \emptyset$  and initialize the corruption state as uncorrupted.
03: \ \ Client sends request
04: On receive ( $\text{sid}$ , "request",  $\theta_S$ ) from  $C$ :
05:   set  $Q = \text{Request}(\theta_S)$ 
06:   notify  $\mathcal{A}$  of ( $\text{sid}$ , "request",  $|Q|$ )
07:   set STORAGE[ $\text{sid}$ ]  $\equiv (Q, \_)$ 
08: \ \ Server send response
09: On receive ( $\text{sid}$ , "response",  $Q, R$ ) from  $S$ :
10:   set  $Q'$  from STORAGE[ $\text{sid}$ ]
11:   check that  $Q' = Q$ , else return  $\perp$ 
12:   notify  $\mathcal{A}$  of ( $\text{sid}$ , "response",  $|R|$ )
13:   update STORAGE[ $\text{sid}$ ]  $\equiv (Q, R)$ 
14: \ \ Forward any data
15: On receive ( $\text{sid}$ , "forward",  $\mathcal{P}$ ) from  $\mathcal{V}$ :
16:   if  $\mathcal{P} = S$  set out  $= Q$  from STORAGE[ $\text{sid}$ ]
17:   if  $\mathcal{P} = C$  set out  $= R$  from STORAGE[ $\text{sid}$ ], else return  $\perp$ 
18:   notify  $\mathcal{A}$  of ( $\text{sid}$ , "forward",  $|\text{out}|, \mathcal{P}$ ), block until  $\mathcal{A}$  replies
19:   send ( $\text{sid}$ , out) to  $\mathcal{P}$ 
20: \ \ Client proves provenance (and statement, optional)
21: On receive ( $\text{sid}$ , "prove",  $R$ ) from  $C$ :
22:   set  $R'$  from STORAGE[ $\text{sid}$ ]
23:   check that  $R' = R$ , else return  $\perp$ 
24:   send ( $\text{sid}$ ,  $\mathcal{R}(R)$ ) to  $\mathcal{V}$ 
25: On receive ("corrupted") from  $\mathcal{V}$ 
26:   set corruption state to corrupted and return an element
27:   randomly sampled from  $\{0, 1\}^\lambda$  when invoked with "forward"
    
```

**Figure 15: The ideal functionality for data provenance via TLS with a proxy verifier  $\mathcal{V}$ .**

only output one bit the following executions are indistinguishable:

$$\forall \mathcal{Z}; \text{EXEC}_{\text{PROT}, \mathcal{A}, \mathcal{Z}} \approx \text{EXEC}_{\mathcal{I}_f, \text{Sim}, \mathcal{Z}} \quad (2)$$

## E.1 Ideal Functionality $\mathcal{F}_{DP}$

We model the security of proving data provenance as an ideal functionality  $\mathcal{F}_{DP}$  in Figure 15. Each participant in the ideal functionality is identified by a unique identifier, denoted as  $\mathcal{S}$ ,  $\mathcal{V}$ ,  $C$  and  $\mathcal{P}$  respectively. Further,  $\mathcal{F}_{DP}$  is parametrized with a relation  $\mathcal{R}$  to be satisfied by the data provided by the server  $\mathcal{S}$ . The unique identification of participants models the requirement of uniquely identifiable communication partners, which is a requirement to account for the possibility of adversarial MITM clients. Note that in the UC model, identifiers assumed to be provided via an underlying lower-level protocol that also guarantees authenticated communication [21]. The network adversary is formalized through the delayed output terminology [20]. At a high level, the client  $C$  first calls the "request" subroutine of  $\mathcal{F}_{DP}$  by providing an authentication secret  $\theta_S$ . The functionality internally stores the query  $Q$  for the session identifier  $\text{sid}$ . Note that only once the proxy verifier  $\mathcal{V}$  calls "forward", the request is forwarded to the server  $\mathcal{S}$ . However, the verifier  $\mathcal{V}$  never learns the content request. The adversary solely learns the size of the request  $Q$ . Once the server  $\mathcal{S}$  receives the response, it checks whether  $\mathcal{V}$  actually forwarded the request sent by  $C$ . Similarly, once the client receives the response  $R$ , it checks whether  $\mathcal{V}$  actually forwarded the response sent by  $\mathcal{S}$ . Intuitively,  $\mathcal{F}_{DP}$  describes implies the following security properties, which are similarly described in related works [70, 74]: **Client Integrity**. A malicious client cannot convince a verifier of an untrue statement

about the TLS session. **Verifier Integrity**. A malicious verifier cannot prevent a client from receiving the correct message from the server. **Privacy**. A malicious proxy verifier cannot obtain any knowledge beyond the size of the message that it forwards.

## E.2 Security Proof

In the following, we show the validity of Theorem 1. For clarity, we highlight that Theorem 1 presumes EUF-CMA security of ECDSA for verification of certificates, and security of the SHACAL-2 compression function  $f$  underlying the SHA-256 hash function, as only the input to the first compression function of the Merkle-Damgard structure is proven in our circuit. To prove that  $\Pi_{\text{ORIGO}}$  UC-realizes  $\mathcal{F}_{DP}$ , we specify a simulator  $\text{Sim}$ , such that no environment  $\mathcal{Z}$  can distinguish an interaction between  $\Pi_{\text{ORIGO}}$  and  $\mathcal{A}$  from an interaction with  $\mathcal{F}_{DP}$  and  $\text{Sim}$ .

*Construction of Sim.* We construct the simulator  $\text{Sim}$  by distinguishing the two main cases for adversarial entities in  $\Pi_{\text{ORIGO}}$  — an adversarial  $C$  and an adversarial  $\mathcal{V}$ . We assume the ideal functionality as presented in Figure 15. As the server  $\mathcal{S}$  is always presumed to be honest, we distinguish between the simulator for (i) an adversarial  $C$  and (ii) an adversarial  $\mathcal{V}$  below.  $\text{Sim}$  runs  $\mathcal{A}_C$  (and for 2)  $\mathcal{A}_V$ ), as well as  $\mathcal{F}_{ZK}$ , internally.  $\text{Sim}$  forwards any input from  $\mathcal{Z}$  and tracks messages sent from and to  $\mathcal{A}_C / \mathcal{A}_V$ . The simulator  $\text{Sim}$  simulates the compression function  $f$  as a random oracle by sending uniformly random strings from  $\{0, 1\}^{256}$  when queried by the adversary.

**1) Sim for malicious  $C$ .** For a malicious client  $C$  we intend to prove *client integrity*.

(1) When  $\mathcal{A}$  sends ( $\text{sid}$ , "request",  $Q$ ) to  $\mathcal{S}$ ,  $\text{Sim}$  forwards ( $\text{sid}$ , "forward",  $\mathcal{S}$ ) to  $\mathcal{F}_{DP}$ . When server sends ( $\text{sid}$ , "response",  $R$ ) to  $C$ ,  $\text{Sim}$  forwards ( $\text{sid}$ , "forward",  $C$ ) to  $\mathcal{F}_{DP}$  and ( $\text{sid}$ , out) to  $\mathcal{A}$ .  $\text{Sim}$  obtains  $(\hat{Q}, \hat{R})$ , where  $\hat{Q}$  is the encrypted request and  $\hat{R}$  is the encrypted response.

(2) When receiving ( $\text{sid}$ , "prove",  $C, x, w$ ) from  $\mathcal{A}$ ,  $\text{Sim}$  checks that  $k := \text{KDC}(\Gamma)$  with the TLS transcript and  $\hat{R} := \text{AEAD.Enc}(R, \text{SATK})$  with AES-GCM.

(3) If the previous checks pass,  $\text{Sim}$  forwards ( $\text{sid}$ , "prove",  $R$ ) to  $\mathcal{F}_{DP}$ , which successively outputs to  $\mathcal{V}$ .  $\text{Sim}$  outputs whatever  $\mathcal{A}$  outputs.

**Indistinguishability.** We now argue that the execution with  $\text{Sim}$  and  $\mathcal{F}_{DP}$  is indistinguishable from the real-world execution with  $\mathcal{A}$  and  $\Pi_{\text{ORIGO}}$  through a set of hybrids. We start off with hybrid  $H1$ , which represents the real-world execution of the protocol.  $H2$  lets  $\text{Sim}$  emulate  $\mathcal{F}_{ZK}$ .  $H3$  lets  $\text{Sim}$  simulate  $\mathcal{V}$ .  $H4$  adds the integrity checks from step 3).

**Hybrid  $H1$ .** Is the real-world execution of  $\Pi_{\text{ORIGO}}$ .

**Hybrid  $H2$ .** proceeds as in  $H1$ , but in addition the simulator  $\text{Sim}$  emulates  $\mathcal{F}_{ZK}$ . All communication between  $\mathcal{A}_C$  and  $\mathcal{S}$  or  $\mathcal{V}$  is sent through  $\text{Sim}$ , which forwards it as in the real execution. It is trivial to see that  $H1$  and  $H2$  are indistinguishable as the simulation of  $\mathcal{F}_{ZK}$  is perfect.

**Hybrid  $H3$ .**  $H3$  is the same as  $H2$ , but  $\text{Sim}$  internally simulates  $\mathcal{V}$ .  $\text{Sim}$  only forwards the client request and server response if it was recorded.  $H3$  is indistinguishable from  $H2$ , as the distribution in the ideal and real world are the same, i.e. both  $\mathcal{V}$  and  $\text{Sim}$  would

abort under the same circumstances when  $C$  presents an invalid proof.

**Hybrid  $H4$ .**  $H4$  is the same as  $H3$ , with the change that the additional check for record existence is performed by  $\text{Sim}$ . To show that  $H3$  is indistinguishable from  $H4$ , we show that both  $\mathcal{V}$  and  $\text{Sim}$  would abort in the same cases. In the real world, the misbehavior of the client would lead to an abort in two cases – the client deviates during the handshake with the server or in the proof presentation. We argue that the same two would lead to an abort in  $H4$  as follows. First, the client may deviate during handshake. This will be detected by  $\text{Sim}$  and will result in exactly the same behavior as in the real world. Second, the client may provide a proof on data that is not associated with the record transmitted. This would lead similarly to an abort in the real and ideal world as  $\text{Sim}$  internally checks whether the sent queries and records match the session id.

It remains to argue that  $H4$  is indistinguishable from the ideal execution. As  $H4$  only accepts a valid statement if the client is in possession of  $\theta_S$  and obtained the records from  $\mathcal{S}$ ,  $H4$  and the ideal world are indistinguishable.

**2)  $\text{Sim}$  for malicious  $\mathcal{V}$ .** For a malicious  $\mathcal{V}$ , we intend to show verifier integrity, such that a malicious  $\mathcal{V}$  cannot prevent a client from receiving correct messages from the server. Recall that  $\text{Sim}$  runs  $\mathcal{A}_{\mathcal{V}}$  internally and forwards all traffic between  $\mathcal{Z}$  and  $\mathcal{A}_{\mathcal{V}}$ .

(1) Upon reception of any message sent by  $\mathcal{A}_{\mathcal{V}}$ ,  $\text{Sim}$  forwards it to  $\mathcal{F}_{DP}$ . During initialization,  $\text{Sim}$  obtains the encrypted handshake transcript and forwards it to  $\mathcal{A}_{\mathcal{V}}$ . The client  $C$  decrypts the handshake and verifies  $pk_S$ .

(2) After  $C$  calls  $\mathcal{F}_{DP}$  with "request" or  $\mathcal{S}$  calls  $\mathcal{F}_{DP}$  with "response",  $\text{Sim}$  obtains the leaked output sizes  $|Q|$  and  $|R|$ .

(3) Upon request of  $\mathcal{A}_{\mathcal{V}}$ ,  $\text{Sim}$  runs the handshake with  $\mathcal{S}$  by simulating  $C$ , obtaining the handshake secrets and traffic encryption secrets negotiated with  $\mathcal{S}$ .

(4)  $\text{Sim}$  sends a random request  $Q$  to  $\mathcal{A}_{\mathcal{V}}$  and forwards to  $\mathcal{F}_{DP}$  whatever  $\mathcal{A}_{\mathcal{V}}$  responds.

(5)  $\text{Sim}$  asserts that  $pk_S$  verifies with  $pk_S$ .

(6)  $\text{Sim}$  prepares the alternative handshake transcript by pre-computing internal hash functions and forwarding them to  $\mathcal{A}_{\mathcal{V}}$ .  $\text{Sim}$  discloses the handshake keys to  $\mathcal{A}_{\mathcal{V}}$ .

(7) In the proof presentation phase,  $\text{Sim}$  sends  $(\text{sid}, \mathcal{R}(R))$  to  $\mathcal{A}_{\mathcal{V}}$  and responds whatever  $\mathcal{A}_{\mathcal{V}}$  responds.

**Indistinguishability.** We argue that the ideal execution with  $\text{Sim}$  and  $\mathcal{F}_{DP}$  is indistinguishable from the real-world execution with  $\mathcal{A}$  and  $\Pi_{ORIGO}$  through a set of hybrids.  $H1$  represents the execution of the real-world protocol.  $H2$  lets  $\text{Sim}$  additionally emulate  $\mathcal{F}_{ZK}$ .  $H3$  additionally lets  $\text{Sim}$  simulate  $C$ .

**Hybrid  $H1$ .** is the real-world execution of  $\Pi_{ORIGO}$ .

**Hybrid  $H2$ .**  $H2$  is the same as  $H1$ , but  $\text{Sim}$  internally simulates  $\mathcal{F}_{ZK}$ . Further,  $\text{Sim}$  calls  $\mathcal{F}_{DP}$  with "request" to receive  $|Q|$  and  $|R|$ . Since the simulation of the ideal functionality is perfect,  $H2$  and  $H1$  are indistinguishable.

**Hybrid  $H3$ .**  $H3$  is the same as  $H2$ , but  $H3$  additionally internally simulates  $C$ .  $\text{Sim}$  engages in the handshake with  $\mathcal{S}$ .  $\text{Sim}$  emulates a handshake transcript encrypted with random keys and chooses  $Q$ ,  $R$  at random and provides them to  $\mathcal{A}_{\mathcal{V}}$ . Finally,  $\text{Sim}$  directly sends  $(\text{sid}, \text{"prove"}, \text{data}, \mathcal{S})$  to  $\mathcal{A}_{\mathcal{V}}$  alongside the pre-computed adapted handshake transcript and the handshake encrypting keys.

It remains to argue that  $H3$  is the same as the ideal world. We argue that this is the case because: When providing the alternative transcript,  $\mathcal{A}_{\mathcal{V}}$  only learns a hash function in the real world, as  $f$  is a random oracle. Further, as EUF-CMA security holds, signed certificates cannot be forged. Due to key independence,  $\mathcal{A}_{\mathcal{V}}$  does not learn anything about the record layer from the handshake layer. Records and Queries appear equally random. In the record layer,  $Q$  and  $R$  are encrypted and of the same size as in the real world, such that they're indistinguishable. Finally,  $\mathcal{A}_{\mathcal{V}}$  receives the same output as in the real-world execution.

*Instantiating  $\mathcal{F}_{ZK}$ .* Universally Composable security of SNARKs has been deeply investigated in previous works [11, 13, 15, 36, 48]. The core problem with ensuring universal composability for SNARKs is that they commonly rely on the Fiat-Shamir transform for non-interactivity and therefore demand a non-black-box extractor. For Groth16, it has been shown that it only satisfies weak simulation extractability [15], a weaker notion than universal composability. Achieving universal composability for e.g. Groth16 demands for custom designs that achieve a stronger notion of simulation extractability with black-box simulation[39]. Recent work shows how to obtain strong simulation extractability for Groth16 with minimal overhead [13]. Further, Abdolmaleki *et. al* [11] provide a generic UC compiler for NIZKs with updatable CRS, such as Plonk, that introduces a  $\approx 1.2\times$  overhead in concrete instantiation. Note, that the implementation of Groth16 and Plonk in gnark is not adapted, and hence they do not satisfy  $\mathcal{F}_{ZK}$ . We leave an instantiation with universally composable SNARKs to future work.