

# SSI, from Specifications to Protocol? Formally Verify Security!

Christoph H.-J. Braun

braun@kit.edu

Karlsruhe Institute of Technology

Karlsruhe, Germany

Tobias Käfer

tobias.kaefer@kit.edu

Karlsruhe Institute of Technology

Karlsruhe, Germany

Ross Horne

ross.horne@strath.ac.uk

University of Strathclyde

Glasgow, United Kingdom

Sjouke Mauw

sjouke.mauw@uni.lu

University of Luxembourg

Esch-sur-Alzette, Luxembourg

## ABSTRACT

We evaluate a bundle of specifications from the Self-Sovereign Identity (SSI) paradigm to construct an authentication protocol for the Web. We demonstrate how relevant standards such as W3C Verifiable Credentials (VC), W3C Decentralised Identifiers (DIDs), and components of the Hyperledger Aries Framework are to be assembled methodologically into a protocol. We make those assumptions from standard trust models explicit that underlie the derived protocol, and verify security and privacy properties, notably secrecy, authentication, and unlinkability. This enables us to formally justify the additional precision that we urge these specifications to consider, to ensure that implementors of SSI-based systems do not neglect security-critical controls.

## CCS CONCEPTS

- **Security and privacy** → **Web protocol security**; *Authentication*;
- **Information systems** → **World Wide Web**.

## KEYWORDS

Web Standards; Self-Sovereign Identity; Formal Verification

### ACM Reference Format:

Christoph H.-J. Braun, Ross Horne, Tobias Käfer, and Sjouke Mauw. 2024. SSI, from Specifications to Protocol? Formally Verify Security! . In *Proceedings of the ACM Web Conference 2024 (WWW '24)*, May 13–17, 2024, Singapore, Singapore. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3589334.3645426>

## 1 INTRODUCTION

The *Self-Sovereign Identity* (SSI) paradigm, popularised in Christopher Allen's seminal blog post [2], refers to the idea of placing users, or more generally *agents*, in control of their digital identity. That is, agents should be able to create digital identities and use them on the Web, without involving a third party when identities are requested, presented or verified. In SSI, security and privacy are declared paramount to protect the user [2, 16, 41]: users must be in control, data minimisation should be observed, and privacy

**Table 1: Layers of common SSI specifications.**

Layer	Function	Most common specification
4	Cryptographic envelopes	DIF DIDComm Messaging [16]
3	Message exchange protocols	Hyperledger Aries [22, 27, 45]
2	Attribute assertion	W3C Verifiable Credentials [41]
1	Agent identification	W3C Decentralised Identifiers [39]

preservation measures are desired. SSI systems are being considered internationally [8, 15, 17, 28, 34] and in various domains, e. g., finance [14], healthcare [26], and the public sector [21, 34].

Standards and specifications underlying these SSI systems are still incomplete, but SSI systems based on these are already being built today: In 2021, hackers from the Chaos Computer Club (CCC), Europe's largest hacker association, demonstrated a flawed application of a common SSI standard in the German driver's license app *ID-Wallet* (cf. <https://github.com/Fluepke/ssi-poc>) showcasing that authentication of the entity that verifies credentials was not guaranteed by the protocol derived from SSI specifications [5]. The app was therefore cancelled just before roll-out on a national scale. This prominent failure urges us to scrutinise security considerations in the specifications regarding their corresponding level of completeness.

To realise an identity and access management system, the following building blocks are required: agent identification, attributes for authorisation, and protocols for authentication. In SSI, the specifications for those building blocks are subject to ongoing and mostly separate standardisation efforts by the World Wide Web Consortium (W3C), the Decentralised Identity Foundation (DIF), and the Hyperledger Community. In Tab. 1, we summarise function and layering of the specifications for W3C Decentralised Identifiers (DIDs) [39], the W3C Verifiable Credentials (VC) Data Model [41], Hyperledger Aries protocols [22, 27, 45], and DIF DIDComm [16].

When we examine the individual specifications, we observe that information on how to implement the standards is scattered across supplementary material, and without elaborating on the security implications. What is more, information on how to properly combine the different standards is also limited: each specification focuses on their domain of interest with little considerations of the other layers of Tab. 1, resulting in the fragmentation of specifications. While this layered thinking is commonplace in software engineering, problems may arise from a security perspective.

That being said, the specifications including DIDComm [16] and the W3C VC data model [41] provide tools that indeed can be used to create secure and privacy-preserving applications, if it's done



This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike International 4.0 License.

WWW '24, May 13–17, 2024, Singapore, Singapore

© 2024 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-0171-9/24/05.

<https://doi.org/10.1145/3589334.3645426>

right. However, when combining these standards, there should be a methodology for ensuring that mistakes compromising security are avoided. We thus formulate the following research questions:

- (1) Can a formal model of an authentication protocol be derived from the SSI specifications?
- (2) What essential security requirements can be distilled from the SSI specifications and related documents?
- (3) Does the formal model satisfy the desired security properties?
- (4) What essential design decisions **MUST** be made in order to guarantee security which are not evident from the standards?

To answer these research questions, we apply the following methodology: First, we review the relevant specifications and standards (Sec. 3.1, 3.2, 4.1) and provide an illustrating example (Sec. 3.3). Next, we present an authentication protocol constructed from the specifications (Sec. 4.2) and provide a formal mapping between the abstract protocol model and the specifications (Sec. 4.3). Based on the combined knowledge on specifications and security best practices, we define necessary trust relations between agents (Sec. 5.1) and map formal security properties back to the informal desire of those from the SSI specifications and related documents (Sec. 4.3). We verify secrecy and authentication properties using the verification tool *Proverif* [6, 7] and privacy properties (unlinkability) using the tool *DeepSec* [12] (Sec. 5.2). Finally, we summarize the essential design decisions required to guarantee the security of SSI protocols as feedback to the specifications (Sec. 6).

## 2 SECURITY METHODOLOGIES AND SSI

We review Allen’s SSI principles in the light of established security and privacy methodologies. To this end, we first present such methodologies. We next use them to identify the following properties as most relevant for SSI: secrecy, authentication, and unlinkability, which we subsequently introduce. We last validate the relevance of those properties by linking them to some of Allen’s SSI principles.

Security properties to be considered can be derived from security methodologies such as STRIDE. Of the threats in STRIDE, Spoofing of user identity, Repudiability, and Information disclosure can be partially addressed by establishing the security properties *secrecy* and *authentication*, as initially articulated by Lowe [33]. Elevation of privilege is also impacted by authentication, since authentication is often established in order for access control to function. Tampering and Denial of service are more perpendicular threats.

Next to those security properties, we consider the privacy property of *unlinkability*, which means that two uses of a system cannot be linked. A system satisfying unlinkability, c.f. ISO/IEC 15408 security standard for information systems (aka. Common Criteria [35]), offers stronger privacy assurance than one satisfying anonymity.

*Secrecy and forward secrecy.* Most threats are impacted by the *secrecy* of long-term keys during regular execution of the protocol. Information disclosure within a session is also impacted by the *secrecy* of material specific to a session of the protocol, such as session keys and nonces. Secrecy holds, if whenever a session involving honest agents completes, it is impossible that an attacker can obtain the secrets in that session. Secrecy can also be evaluated in the face of long-term keys being compromised, e.g. via a data breach, by checking *forward secrecy*. Formally, forward secrecy is modelled

as two phases, the first where the protocol runs normally, and the second where the long-term keys are revealed. We check whether then the secrets in the first phase remain secret. This ensures that information shared before a data breach remains secret.

*Authentication.* Threats such as spoofing of identity and repudiation can be addressed by formal authentication properties, where *agreement* is among the strongest. Agreement ensures that, when one party completes the protocol, we can assume that the other parties performed all previous actions in the protocol, and, for corresponding pairs of send and receive actions, the data was the same. Agreement is, in addition, said to be *injective*, if for every session completed, there is a unique session for each of the other parties involved. Injectivity is required to prevent *replay attacks*.

*Unlinkability.* is sometimes referred to as non-correlability in anonymous credential systems. We will explain that, even for credentials that are not anonymous, unlinkability can be achieved *from the perspective of the issuer*. Specifically, after issuing a credential, the issuer cannot track how it is used with honest verifiers [20]. In this way we formalises data sovereignty, in the sense that the issuer cannot monitor the usage of a credential after issuance.

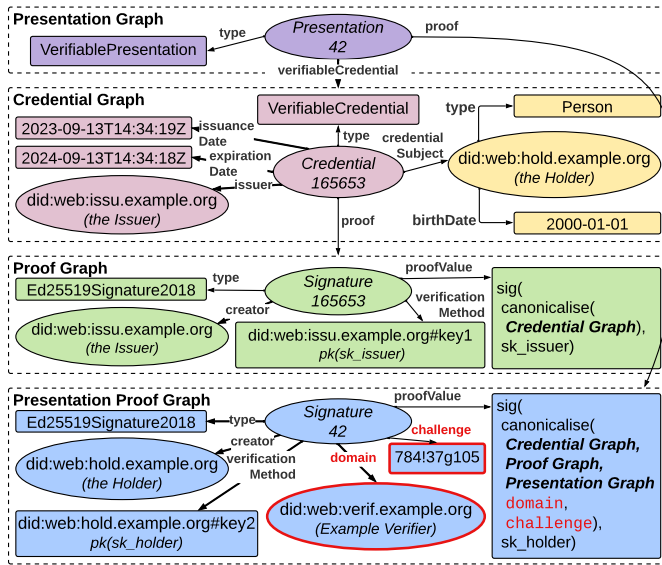
*Security & privacy for SSI.* We reinforce the need for the presented properties by connecting them to Allen’s SSI principles [2]. Allen’s 2nd SSI principle, **CONTROL**, refers to users being able to control what information about them is revealed and what is kept secret. Thus control desires *secrecy* and *authentication*, where authentication lends assurance regarding the context in which data is revealed. Allen’s 3rd SSI principle, **ACCESS**, refers to users being able to access information about themselves, but at the same time keeping it *secret* from others. Allen’s 5th SSI principle, **PERSISTENCE**, stipulates that identities must be long-lived in the sense that identities may be retained when keys they map to are rotated. Thus persistence is supported by *forward secrecy*, which ensures secrecy in scenarios where key rotation is necessary. Allen’s 9th SSI principle, **MINIMISATION**, desires non-correlability, aka. *unlinkability*, while acknowledging that it is difficult to fully achieve. This is consistent with our observation that some but not all unlinkability properties hold. Allen’s 10th SSI principle, **PROTECTION**, demands identity authentication to occur independently from potential interference by a third party to ensure the rights of individual users. While this principle focuses on protecting the rights of users, it hints at *unlinkability from the perspective of issuers*, to which we want to draw attention in this work. It also hints at *forward secrecy* since that mitigates against obtaining keys via coercion. *Authentication* and hence *agreement* is explicitly mentioned in this principle.

## 3 WEB STANDARDS AS BASIC BUILDING BLOCKS OF AN SSI PROTOCOL

We introduce DIDs and VCs, recommended by the W3C as Layer 1 and 2 of an SSI protocol, and sketch their application using a simplified example protocol, omitting technologies for Layers 3 and 4.

### 3.1 W3C Decentralised Identifiers (DID)

Contrary to centrally managed identifiers, a Decentralised Identifier (a DID) is under control of the agent that it refers to. The Decentralised Identifiers W3C recommendation [39] specifies how



**Figure 1: The graph-based VC data model according to its JSON-LD context [41]. Notice the signature-relevant attributes of challenge (a nonce) and domain from the presentation proof graph marked with red.**

such an agent proves control over a DID. A DID maps to a *DID document* that incorporates information about cryptographic public keys, which may be used by the agent controlling the DID to prove their control over the DID. The mapping between the DID and its DID document is defined by a *DID method*. The DID methods, which can be defined and registered by anyone [43], define how to retrieve a DID document from a DID. DID methods typically involve a form of secure lookup to obtain the DID document, such as `did:web` [36], which defines provisioning of DID documents via TLS, or `did:ethr` [44] and `did:sov` [29], which define blockchain-based DIDs on Ethereum or Hyperledger Indy, respectively. Locally resolvable DID Methods include `did:key` [42], where a DID document or associated public key is encoded in the DID itself.

For example, in Fig. 1, the DID `did:web:issu.example.org` resolves to a DID document containing the public keys of the agent issuer. Furthermore, the URI `did:web:issu.example.org#key1` identifies which of the agent’s keys are employed. We clarify that, while not explicitly specified, an explicit check is required to ensure that the given key URI is among those listed in the DID document.

### 3.2 W3C Verifiable Credential (VC) data model

A Verifiable Credential (VC) according to the W3C recommendation [41] is a Resource Description Framework (RDF) dataset comprised of two RDF graphs: the *credential graph* containing claims and attributes, which links to a *proof graph* containing the credential’s digital signature and metadata concerning its interpretation. The *claims* of the credential graph are the statements about the `credentialSubject` that are asserted by the issuer (e.g. the birthdate in Fig. 1). Claims include metadata such as a DID identifying the issuer who signs the credential, its `expirationDate` and `issuanceDate`. The proof graph indicates the type of proof

explained generically in a separate Verifiable Credential Data Integrity working draft [40]. The Data Integrity specification permits multiple types, e.g. `Ed25519Signature2018`, which define how to establish a signature’s validity using a specific signature scheme and algorithm for obtaining a canonical representation of the RDF graph underlying the credential [3, 23, 30, 31]. Fig. 1 suggests, using typically notation employed in security (Dolev-Yao style [19]), how signatures are generated using secret keys (*sk*) of agents and the graphs. A term of form  $\text{sig}(M, sk)$  denotes a *signature* on bitstring *M* using a private key *sk*, and  $\text{pk}(sk)$  denotes the corresponding public key. This symbolic approach to security abstracts away the implementation of signature schemes and canonicalise functions.

The VC specification [41] also defines a data model for Verifiable Presentations (VPs), which are necessary for the credential subject, holding the VC, to prevent trivial replay attacks when a credential is presented to another agent. A VP ties a VC cryptographically to a particular session, by signing the VC along with session information to certify that the holder has approved that the VC may be used in the specified session only. The VP is described in a *presentation graph* which links to the *credential graph* of the relevant VC and the *presentation proof graph*, which describes a digital signature on the presentation graph. Fig. 1 shows a VP comprised of the four graphs: the *presentation graph* and *presentation proof graph* and the two graphs of the VC, i.e., the *credential graph* and *credential proof graph*. As suggested symbolically in Fig. 1, the session information signed by the signature in the presentation proof graph is a serialisation of the presentation graph **and** optional attributes in the proof graph, notably the domain and challenge. This arguably confusing decision, i.e., to place some information signed (the domain and challenge) outside the presentation graph, is mandated by the Data Integrity draft [40], and implicitly in the VC specification’s examples [41].

### 3.3 An example of authentication using VCs

Vcs enable an agent (the *holder*) to prove to a second agent (the *verifier*) that a third agent (the *issuer*) has asserted and signed some claims about the holder. In other words, with a presentation of a VC, an agent proves to a verifying agent that:

- they are in possession of the VC;
- the VC was issued by a particular issuer;
- the VC contains some claims, e.g., attributes of the holder;
- the VC was presented by the holder itself, for the purpose the verifier intended, e.g., authenticated resource access on the Web.

As example, consider the use case of a student accessing online teaching material of a guest professor. The student’s university (issuer) provides the student (holder) with a digital student credential. It is signed by the university and asserts that the holder is a student.

A guest professor at the university provides online teaching material, served from their personal Web server, to the university’s students. To access the online material, students have to prove that they are really enrolled at the university by creating a VP asserting a signature on the VC along with an identifier (DID) for the teacher and other session-specific information. The professor verifies the signature on the VC and VP using the public keys of the university and student respectively, and checks that the claims and session parameters are as expected.

Students verify that they are really talking to the professor by looking up a trusted mapping between identifier (e. g. DID) and public keys. In the Web-based case at hand, the professor's homepage may advertise the professor's DID. Thereby, the here-described system relies on the Domain Name System (DNS) and the transport protocol HTTPS to ensure the integrity of the identifier-agent mapping. Other approaches for maintaining such identification mapping may include government registries, governed blockchains or smart contracts. This mapping is commonly (and commonly implicitly) deemed out-of-scope by SSI authentication protocols as a system-level governance challenge. We make this assumption explicit in Sec. 5.1 as the *DID document and proof method assumptions*.

We note that the W3C Working Group Note on VC Use Cases [37] proposes 30 use cases from 7 domains. Among them, 25 concern authentication of the holder, similarly to the one described above, while 5 concern transferability and revocation. Trust assumptions or protocols are not made explicit, as our paper addresses in Sec. 5.1.

## 4 CONSTRUCTING AUTHENTICATION PROTOCOLS FOR SSI

To construct authentication protocols using the presented Web standards (for Layers 1 and 2), we need to add protocol components for Layers 3 and 4. We therefore first present potential protocol components for Layers 3 and 4, which we find in the Hyperledger Aries protocols and DIDComm. We then present the thus derived SSI authentication protocol for Web resource access, and map the protocol in detail to the SSI specifications.

### 4.1 Potential protocol components

The Hyperledger Aries community is in the process of defining a framework of protocols for creating, transmitting and storing verifiable digital credentials. We cover the three protocols intended for building VC-based authentication protocols. These protocols are agnostic to the specific data models and formats of the transmitted payload, e. g., DIDs, VCs and VPs.

**The Aries DID Exchange protocol** [45] is a protocol for establishing a session between agents using DIDs. It defines three message types: a *request* communicating the DID of the requesting agent; a *response* completing the exchange from the responding agent; and a *complete* message confirming the exchange from the requesting agent to the responding agent.

**The Aries Issue Credential protocol** [22] defines two message types: a *request-credential* message for a holder to request issuance of a VC, and an *issue-credential* message containing the VC.

**The Aries Present Proof protocol** [27] defines two message types: a *request-presentation* message from the verifier requesting a verifiable presentation, and a *presentation* message containing a VP from the holder. These protocols list "attachment registries" linking to possible data models for messages, including VCs and VPs.

**The DIF DIDComm Messaging.** The above protocols define only a message flow and do not consider how messages are encrypted on the wire – the intention being that the protocols are meant to be used on top of DIDComm Messaging [16]. Specified by the Decentralised Identity Foundation (DIF), DIDComm Messaging [16] is a methodology for encrypting and signing messages, using keys in the DID documents of communicating agents.

### 4.2 A thus constructed authentication protocol

From the SSI specifications, we aim to derive an authentication protocol for Web resource access, and specify it at a level of precision amenable to symbolic verification. We notice that application-specific extensions to the message exchange protocols (Layer 3) are required to construct a functional protocol. The result is a three-party protocol comprised of two two-party protocols: First, the **issuance of a VC** is conducted, such that this VC can be used in multiple sessions of the **provenance of a VC**. A recently proposed architecture [9] implements this protocol, without formal verification of its security properties<sup>1</sup>.

The protocol is informally illustrated in Fig. 2 as a message sequence chart. For the protocol's complete applied  $\pi$ -calculus specification [1] amenable to formal verification, see Appendix A Tab. 3 - 4. In keeping with Dolev-Yao style symbolic notation,  $\{M\}_K$  denotes the *encryption* of bitstring  $M$  with a public key  $K$ . We notate by  $\text{proj}_i M$  the  $i^{\text{th}}$  projection of a tuple  $M$  of the form  $\langle M_1, \dots, M_i, \dots \rangle$ ,  $\text{dec}(M, K)$  is the *decryption* of ciphertext  $M$  using private key  $K$ , and predicate  $\text{check}(M, K)$  checks a signature  $M$  against a public key  $K$ . We use the following notation:  $I, P, V$  are the DIDs of Issuer, holder (aka. **Prover**), and **Verifier**, respectively.  $(ssk_X, spk_X)$  is a session key pair of agent  $X$  with  $spk_X = pk(ssk_X)$ ; and the long-lived key pair  $(sk_X, pk_X)$  with  $pk_X = pk(sk_X)$  correspondingly.  $n$  are nonces and  $s$  are signature values calculated via  $s_k = \text{sig}(m'_k, sk_X)$ . Constant *attr* serves as an attribute to assert, *URI* is the URI of a desired Web resource under access control, and *RULE* is an access control rule expressing expected claims in a VC.

Both two-party (sub-)protocols consist of two sub-sub-protocols: A "handshake protocol" and an application-specific "follow-up protocol". Both two-party protocols start with the Aries DID Exchange Protocol as the handshake protocol. Seamlessly, the issuance of a VC continues with the Aries Issue Credential Protocol as follow-up. Similarly, the provenance of a VC follows then an extended version of the Aries Present Proof Protocol.

In particular, the Aries Present Proof protocol is extended with an additional *access-request* message that includes the URI of a resource the holder wishes to access, and also includes an *access-response* message communicating an access token if the present proof protocol succeeds. The application-specific messages wrap the Aries Present Proof protocol. Similarly, the last message of the DID Exchange Protocol between two agents – a *complete* message – is simultaneously the first message of the follow-up protocol in Fig. 2. For the issuance, the *complete* message is simultaneously a *request-credential* message; while for provenance, the *complete* message is simultaneously an *access-request* message.

### 4.3 Connecting protocol and specifications

We map the protocol (cf. Fig. 2) to the specification layers (cf. Tab. 1).

#### Layer 1: Agent identification

- We interpret  $(sk_X, pk(sk_X))$  to be the long-lived key pair of an agent  $X$ .
- We interpret the long-lived identifiers of an agent,  $X$  with  $X \in (I, P, V)$ , to be DIDs, e. g., using `did:web` [36]. By Sec. 5.1, we

<sup>1</sup>[9] is all about the interoperability of SSI standards, especially different flavours of VC and DID, and proposes an architecture that addresses those interoperability challenges.

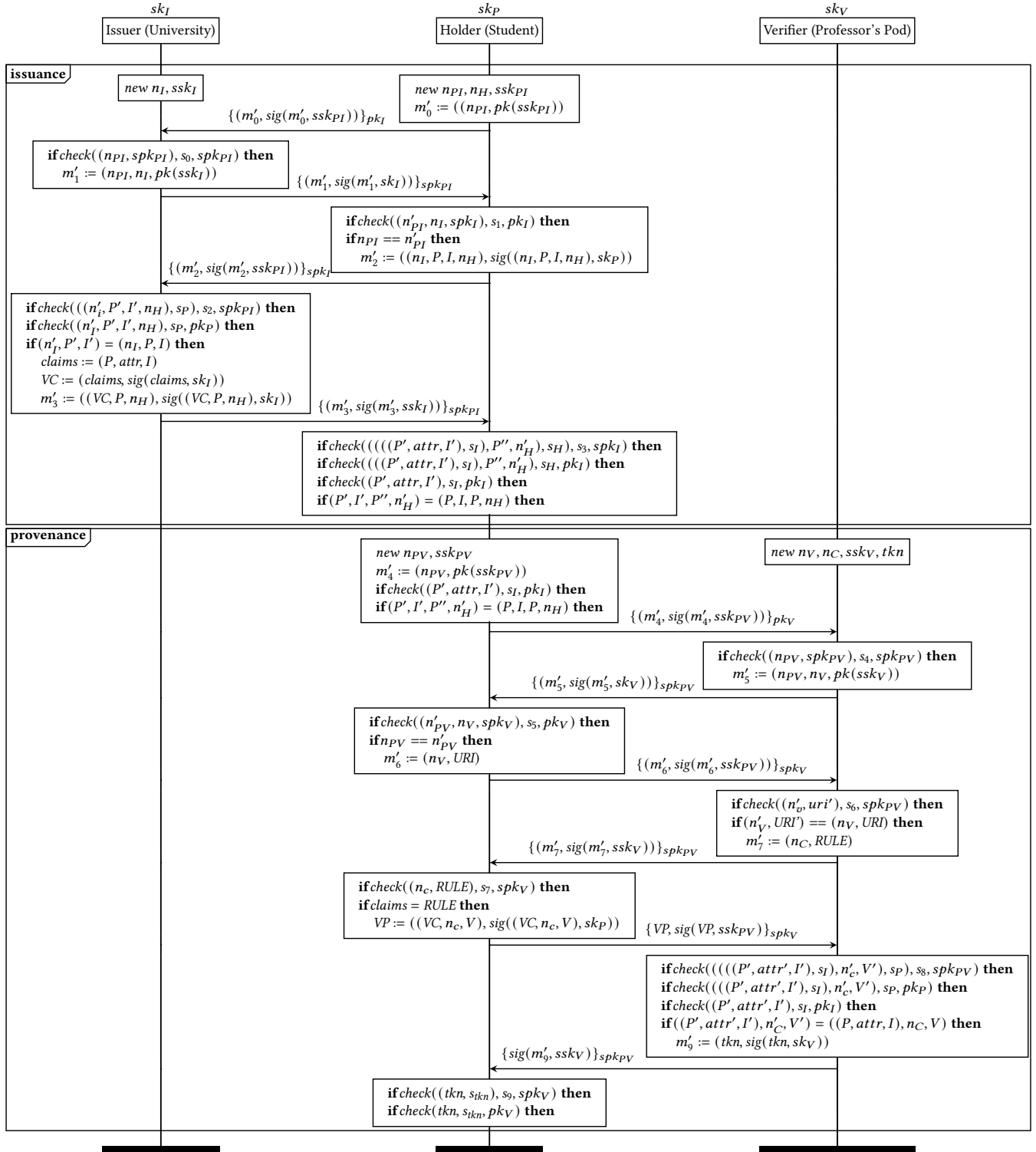


Figure 2: An SSI authentication protocol consisting of issuance and provenance sub-protocols.



assume the public key to be obtainable from an agent's DID;  $pk(sk_X) = getPubKey(X)$ .

- We interpret  $(ssk_X, pk(ssk_X))$  to be a short-lived session key pair of an agent  $X$ . We interpret the public key of a session key pair  $spk_X = pk(ssk_X)$  to be encoded in a DID, e. g., using `did:key` [42], when transmitted in messages.

### Layer 2: Attribute assertion

- By Sec. 5.1, we assume the issuer to have verified that the holder is actually exhibiting the attribute *attr* to be asserted.
- We interpret  $claims = (P, attr, I)$  according to the W3C VC data model [41], with  $P$  being the `credentialSubject`, *attr* being some claim, e. g., `birthDate` from Fig. 1, and  $I$  being the issuer.
- We interpret  $VC = (claims, sig(claims, sk_I))$  according to the W3C VC data model [41], where the *claims* are RDF triples forming the *credential graph* and the signature value  $sig(claims, sk_I)$  is the *proofValue* of the *proof graph*. The *verificationMethod* is  $pk(sk_I)$ . For the signature, we assume *claims* in canonicalised form (i. e. after *canonicalise(Credential Graph)* from Fig. 1).
- We interpret  $VP = ((VC, n_C, V), sig((VC, n_C, V), sk_P))$  according to the W3C VC data model [41], where the  $VC$  is linked via *verifiableCredential* from the presentation. In the *presentation proof graph*,  $sig((VC, n_C, V), sk_P)$  is the *proofValue* with  $pk(sk_P)$  the *verificationMethod*,  $n_C$  the challenge, and  $V$  the domain. For the signature, we assume  $(VC, n_C, V)$  in canonicalised form (i. e. after *canonicalise(Credential Graph, Proof Graph, Presentation Graph, domain, challenge)* from Fig. 1).

### Layer 3: Message exchange protocols

- We interpret the first three messages exchanged between any two agents,  $m'_0 - m'_2$  and  $m'_4 - m'_6$ , according to the Aries DID Exchange Protocol [45]. We notice that nonce is **not required** by the specification. Let  $X \in \{PI, PV\}$  and  $Y \in \{I, V\}$ :
  - $m'_0$  and  $m'_4$  are *request* messages with  $pk(ssk_X)$  as *did* and  $n_X$  as nonce.
  - $m'_1$  and  $m'_5$  are *response* messages with  $pk(ssk_Y)$  as *did* and  $n_Y$  and  $n_X$  as nonce.
  - $m'_2$  and  $m'_6$  are *complete* messages with  $n_Y$  as a nonce.
- In the interaction between Issuer and Holder, we interpret the two last messages,  $m'_2$  and  $m'_3$ , according to the Aries Issue Credential Protocol [22]. We notice that domain and challenge are **not required** by the specifications.
  - $m'_2$  is a *request-credential* message that contains an attachment of:  $n_I$  as nonce,  $n_H$  as challenge,  $P$  as holder,  $I$  as issuer, and a corresponding signature value as *proofValue*.
  - $m'_3$  is a *issue-credential* message that contains an issued credential as an attachment: We interpret the attachment to be a VP according to the W3C VC data model [41]. It includes the freshly issued VC as *verifiableCredential*,  $n_H$  as challenge and  $P$  as domain.
- In the interaction between Holder and Verifier, we interpret the two messages  $m'_7$  and  $VP$  (corresponding to  $m'_8$ ) according to the Aries Present Proof Protocol [27]. We notice that domain and challenge are **not required** by the specifications.
  - $m'_7$  is a *request-presentation* message that contains a Verifiable Presentation Request as an attachment. We interpret the attachment according to some attachment data model definition, e. g., provided by [9]: It includes  $n_C$  as challenge,  $V$  as

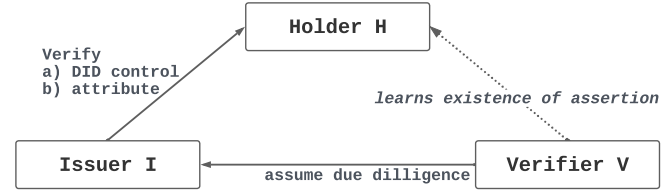


Figure 3: The trust assumption triangle.

domain and *RULE*, the definition of the VC to present, e. g., as *requiredCredential*.

- $VP$  (i. e.  $m'_8$ ) is interpreted as a *presentation* message containing the actual VP as an attachment: We interpret the attachment as a VP according to the W3C VC data model [41] with the VC as *verifiableCredential*,  $n_C$  as challenge and  $V$  as domain.
- The two messages,  $m'_6$  and  $m'_9$ , are interpreted according to the mentioned extension of Aries Present Proof.
  - $m'_6$  is a *access-request* message that includes *URI* as target.
  - $m'_9$  is a *access-response* message; includes *tkn* as *accessToken*.

### Layer 4: Cryptographic envelopes

- We interpret  $m = \{(m', sig(m', sk_S))\}_{pk(sk_R)}$  to model a message  $m$  with payload  $m'$  subject to signature using the sender's private key  $sk_S$  and encryption using the receivers' public key  $pk(sk_R)$ , i. e., *authcrypt* as defined by DIDComm [16].
- We assume automatic decryption of a message  $m$  to its payload  $m'$  if possible for an agent. We then interpret *check*( $m', s, pk(sk_S)$ ) as explicitly checking the signature value  $s$  of payload  $m'$  using the sender's public key  $pk(sk_S)$ , as required by *authdecrypt* defined by DIDComm [16].

## 5 TRUST, SECURITY AND PRIVACY, FORMALLY VERIFIED

We formally justify the correctness of the constructed protocol. Firstly, we make explicit the trust assumptions that must be respected for secure functioning of the protocol. We then use standard tools to verify a comprehensive range of security and privacy properties. We highlight throughout how trust assumptions and properties verified are connected to SSI principles and systems.

### 5.1 Trust assumptions necessary for SSI

In order to reason about security and privacy it is essential to make explicit the underlying trust assumption against which we verify the protocol. We declare the following assumptions about the agents and the underlying infrastructure of identifiers and cryptographic keys. These assumptions are often not stated by SSI protocols and specifications, but relied on implicitly.

- (1) The **Self-Sovereign Identity (SSI) assumption**: All agents can mint and manage key pairs in a self-sovereign manner and honest agents never intentionally publish their private keys.
- (2) The **DID document assumption**: All agents assume the integrity of the link between the DID of an honest agent and a DID document containing public keys of the honest agent. Thus the infrastructure employed by the honest agents for their DIDs must be trusted.

- (3) The **proof method assumption**: If the proof graph contains a URI indicated by `proofMethod` then the *key* extracted from the URI must also appear in the DID document obtained via the DID of the relevant agent (issuer in a VC or holder in a VP).<sup>2</sup>
- (4) The **Verifier-Issuer assumption**: An honest verifier assumes that an honest issuer has conducted due diligence when validating the assertions signed by the issuer, e.g. that the holder is actually a student (which may be out-of-band).
- (5) The **well-specified assumption**: Honest parties assume that parties they trust follow the protocol, even during a data breach. Importantly, honest agents may also engage in sessions with agents that do not follow the protocol [19, 32], which is reasonable since attackers assuming roles within the system can co-exist with honest participants and those attackers may exploit their position to interfere with sessions between honest agents.

The goal with any SSI authentication protocol is to establish a trust relationship from the verifier to the holder, transitively via the issuer by means of asserting and signing claims, as suggested in Fig. 3. The verifier trusts the issuer to have asserted correct information about the holder. The fact that the issuer has some relationship with the holder, is only revealed by the holder upon presentation of the credential. Then, transitively, the verifier may trust the holder to exhibit a certain attribute which has been attested by the issuer. Inversely, the holder must be willing to present the credential to the verifier. This case is similar to the holder trusting the issuer to be the (honest) issuer when revealing private information for validation of attributes to include in the credential.

## 5.2 Results of security and privacy analysis

Based on a model reflecting these trust assumption we return to the security and privacy properties laid down in Sec. 2 to verify that they hold. Proofs of all properties are summarised in Tab. 2.

*Forward secrecy.* Relevant secrecy and forward secrecy properties are formally verified in rows 1-2, respectively, of Tab. 2. Our formal analysis shows that session secrets in the past are preserved even if the long-term private keys of all agents are revealed. Some information may be leaked without compromising other properties, notably the VC itself and information about the access control policy is leaked to an attacker posing as a verifier or holder, respectively. The VC is leaked, as the holder may present the credential to a compromised verifier, and the VP contains the VC (this is not an attack, since the attacker cannot use the VC). The access control rule is leaked, as the protocol explains to anyone who asks what credential is required to access a resource via a URI it controls.

*Authentication.* There are multiple agreement properties [33] to check for the protocol. Between two parties we have: if the issuer completes the protocol, then it injectively agrees with the holder regarding the first three messages of the protocol. This ensures that the issuer really issued the credential to the holder it believes it did. If the holder reaches the fourth message in the protocol, then it injectively agrees with the first four messages of the issuer. This ensures that the holder really received a credential from the intended

issuer. If the holder completes a session with a verifier, then it injectively agrees with all five message exchanges in a session with a verifier. This ensures that the holder really presented a credential to the intended verifier. If the verifier completes the protocol, then it injectively agrees with its first four messages with the holder. This ensures that the presentation was really received from the agent concerned. The final property fails if `domain` is omitted in the VP. See row 6 of Tab. 2 and further explanation in Appendix A.1.

Since three parties are involved in SSI, additional assurance regarding authentication can be achieved if we check a multi-party agreement property [13], where one agent establishes a belief about two or more other agents. If the verifier completes the protocol, then it agrees with all messages of both the holder and issuer (excluding the final message sent). This ensures that if a credential was presented by a holder, then that credential originates in a legitimate session with an issuer. This property is non-injective, since a credential may be issued once and used many times, meaning there is not a one-to-one correspondence between verifier sessions and issuer sessions. Perhaps surprisingly, the above multi-party agreement property does not follow from the two-party agreement properties. Indeed we were able to uncover the presence of an attack on multi-party agreement, which cannot be detected using two-party agreement if the holder were to neglect to check the signature on a VC they are issued. Specifically, an attacker may pose as an issuer and re-issue VCs of an honest issuer. See row 8, Tab. 2 and also Appendix A.2 for details on multi-party agreement.

All authentication properties above hold even if VPs in previously completed sessions are leaked (e.g. due to a data breach or a requirement to reveal logs). This compromise situation is important to note, since if we mistakenly did not include the challenge in the VP then all authentication properties from the perspective of the verifier fail once VPs are revealed. This compromise situation with the challenge present and missing is presented in respective rows 3 and 6 of Tab. 2. For a complete picture regarding agreement, we also check that, even if the holder is compromised, there is a non-injective agreement between the verifier and the issuer, regarding their common data, namely the VC (see row 5 Tab. 2). This means that credentials from honest issuers cannot be forged.

*Unlinkability.* We formulate unlinkability of honest holders in the presence of malicious issuers as an equivalence problem between (A) a process where honest holders and verifiers exchange the same VC twice and (B) an idealised process where each session between honest holders and verifiers involves a fresh VC. In order to model the issuer as an attacker, the long-term keys of the issuer are revealed to the attacker. A proof established using the DeepSec tool supports row 9 Tab. 2. See also Appendix A.3.

A stronger property is unlinkability of the holder against colluding issuers and verifiers, which further ensures that verifiers cannot link two uses of the same credential. This property cannot be achieved for regular VCs, since the identity of the holder appears in each verifiable presentation (row 10 Tab. 2). However, it is achieved for anonymous credentials that hide the identity of the holder using zero-knowledge proofs (row 13 Tab. 2).<sup>3</sup>

<sup>2</sup>For a clarification of “checked against” we refer to: “Dereferencing a public key URL reveals information about the controller of the key, which can be checked against the issuer of the credential.” [41]

<sup>3</sup>Related work on cryptographic schemes for anonymous credentials analyses unlinkability in the face of colluding issuers and verifiers using computational cryptographic methods such as universal composability [10]. Computational methods further check

**Table 2: Structure of our code repository for formal verification of security properties for our instance of an SSI authentication protocol. Paths relative to <https://github.com/uvdsl/ssi-protocol-verify/tree/main/>. See <https://doi.org/10.5281/zenodo.10654423>.**

Protocol	Property	No.	Relative File Path in Repository	OK	Attack
Plain VCs (PlainVCs/DIDComm/)	Secrecy	1	ssipv.pv#L287	✓	
		2	archive/ssipv_forward_secrecy.pv	✓	
	Agreement	3	ssipv.pv#309	✓	
		4	ssipv_ok_VP_leaked.pv	✓	
		5	ssipv_unforgeable_VC.pv	✓	
		6	ssipv_attack_domain_missing_replay.pv	×	<i>masquerade as prover</i>
	Unlinkability	7	ssipv_attack_no_nonce_VP_leaked.pv	×	<i>replay credential</i>
		8	ssipv_attack_VC_reissued.pv	×	<i>reissue old credential</i>
		9	ssipv_unlinkable.dps	✓	
		10	ssipv_attack_verifier_unlinkability.dps	×	<i>verifier tracks prover</i>
Anon VCs (AnonVCs/DIDComm/)	Secrecy	11	ssipv.pv#L297	✓	
	Agreement	12	ssipv.pv#L319	✓	
	Unlinkability	13	ssipv_unlinkability_ok_wrt_verifier.dps	✓	
Plain VCs + Diffie Hellmann (PlainVCs/DIDComm+DH/)	Secrecy	14	ssipv.pv#L302	✓	
	Agreement	15	ssipv.pv#L324	✓	
Anon VCs + Diffie Hellmann (AnonVCs/DIDComm+DH/)	Secrecy	16	ssipv.pv#L312	✓	
	Agreement	17	ssipv.pv#L334	✓	

*Further protocols and anonymous credentials.* The examined protocol is, of course, not a unique solution for SSI on the Web. For example, we have verified variants that open with a Diffie-Hellman handshake in place of the DIDComm Exchange DID handshake (see lines 14-17 of Tab. 2). A natural question is why, given privacy is reflected in SSI principles, have we mainly discussed non-anonymous verifiable credentials that reveal the DID of a prover to verifiers, rather than anonymous credentials. The reason is that the security properties of anonymous credentials only hold if trust assumptions are strengthened. In particular, the security of an individual agent depends on the honesty of the *entire group of agents holding the same credential*, (this is reflected in the model employed to verify rows 11-12 of Tab. 2). Those proofs involve a richer message theory modelling BBS+ zero-knowledge proofs and a modified protocol (not shown). Due to this weakened trust assumption, for anonymous credentials, Allen’s SSI principle of **CONTROL** is weakened, that is, control becomes a collective responsibility not entirely ones own. This explains our focus on cryptographically simpler VCs. To strengthen trust in order for the security of anonymous credential systems to function, adequate wallet management measures of a group of agents must be made explicit, e.g. the issuer should authenticate an attested wallet rather than the holder directly, or all employees holding an attribute need adequate security training. On the other hand, anonymous credentials do strengthen unlinkability, as explained in Sec. 5.2, and in turn the SSI principle of minimisation. Thus there is trade-off between trust and privacy when choosing between anonymous and regular credentials.

## 6 CONCLUSION

We presented in Sec. 4.3 a mapping between a symbolic model of an SSI protocol (Fig. 2) and specifications for SSI in Sec. 3 and 4.1. This has enabled us to use symbolic security tools to verify a range of

that properties of schemes instantiating signature abstractions,  $\text{sig}(\dots)$ , are robust against brute-force attacks. Computational proofs therefore complement our symbolic proofs that concern logical flaws in the usage of schemes in the flow of a protocol [4, 24]. Symbolic methods can also be compositional [25].

security and privacy properties summarised in Table 2. These properties formally support the argument that the constructed protocol is indeed in alignment with the principles of SSI. The most important insight that we reinforce throughout the paper is that certain parameters marked as optional in specifications are not optional. Notably, omitting the domain and challenge in the VP leads to critical attacks allowing attackers to authenticate themselves using the credentials of honest agents (rows 6-7 Tab. 2). Trust clarifications that do not appear explicitly in specifications, notably the *proof method assumption* in Sec. 5.1, are critical for all properties. The role of trust assumptions in ensuring properties verified underscores the importance of spelling out such trust assumptions and protocol design decisions to mitigate vulnerabilities in SSI-based systems.

We believe that our methodology, which is to connect elements of SSI to standard security models, is general enough to be applied to evaluate protocols tailored to other SSI use cases. In particular, we have explained that DIDs map to identities, as they typically appear in security protocols, and their resolution to a public key, is the typical trust assumption that the honest agents know the mapping between honest identities and public keys. We have also explained how elements of VC standards and signatures in proofs can be represented symbolically in a protocol specification, and how layers provided by Hyperledger Aries and DIDComm may be assembled. We acknowledge that different use cases may require a slightly different assembly of the standards, some of which we touch on in Sec. 4.3. We also explained how such mappings can be used to provide genuine insight in the compliance with SSI principles by connecting those principles to standard security properties.

## ACKNOWLEDGMENTS

This research was partially funded by the COST (European Cooperation in Science and Technology) Action on Distributed Knowledge Graphs (CA19134) and partially supported by the German federal ministry of education and research (BMBF) in MANDAT (FKZ 16DTM107B).



## REFERENCES

- [1] Abadi, M., Blanchet, B., Fournet, C.: The applied pi calculus: Mobile values, new names, and secure communication. *J. ACM* **65**(1), 1:1–1:41 (2018). <https://doi.org/10.1145/3127586>
- [2] Allen, C.: The path to self-sovereign identity (2016), <http://www.lifewithalacrity.com/2016/04/the-path-to-self-sovereign-identity.html>
- [3] American National Standards Institute: Public key cryptography for the financial services industry: the Elliptic Curve Digital Signature Algorithm (ECDSA). X9.62, ANSI (2005)
- [4] Arapinis, M., Chothia, T., Ritter, E., Ryan, M.: Analysing unlinkability and anonymity using the applied pi calculus. In: Proceedings of the 23rd IEEE Computer Security Foundations Symposium (CSF). pp. 107–121 (2010). <https://doi.org/10.1109/CSF.2010.15>
- [5] Biselli, A.: Konzeptionell kaputt und ein riesiger Rückschritt (2021), <https://netzpolitik.org/?p=338612>
- [6] Blanchet, B.: An Efficient Cryptographic Protocol Verifier Based on Prolog Rules. In: Proceedings of the 14th IEEE Computer Security Foundations Workshop (CSFW). pp. 82–96 (2001)
- [7] Blanchet, B., Cheval, V., Cortier, V.: ProVerif with lemmas, induction, fast subsumption, and much more. In: Proceedings of the 43rd IEEE Symposium on Security and Privacy (S&P). pp. 205–222 (2022). <https://doi.org/10.1109/SP46214.2022.9833653>
- [8] Boysen, A.: Decentralized, self-sovereign, consortium: The future of digital identity in Canada. *Frontiers Blockchain* **4**, 624258 (2021)
- [9] Braun, C.H.J., Papanchev, V., Käfer, T.: SSSI: an architecture for semantic interoperable self-sovereign identity-based access control on the Web. In: Proceedings of the 32nd Web Conference (WWW). p. 3011–3021. ACM (2023). <https://doi.org/10.1145/3543507.3583409>
- [10] Camenisch, J., Dubovitskaya, M., Haralambiev, K., Kohlweiss, M.: Composable and modular anonymous credentials: Definitions and practical constructions. In: Proceedings of the 21st International Conference on the Theory and Application of Cryptology and Information Security (ASIACRYPT). pp. 262–288. Springer (2015). [https://doi.org/10.1007/978-3-662-48800-3\\_11](https://doi.org/10.1007/978-3-662-48800-3_11)
- [11] Chen, J., Paxson, V., Jiang, J.: Composition kills: A case study of email sender authentication. In: Proceedings of the 29th USENIX Security Symposium (USENIX Security 20). pp. 2183–2199 (2020)
- [12] Cheval, V., Kremer, S., Rakotonirina, I.: DEEPSEC: deciding equivalence properties in security protocols theory and practice. In: Proceedings of the 39th IEEE Symposium on Security and Privacy (S&P). pp. 529–546 (2018). <https://doi.org/10.1109/SP.2018.00033>
- [13] Cremers, C., Mauw, S.: Operational Semantics and Verification of Security Protocols. *Information Security and Cryptography*, Springer (2012). <https://doi.org/10.1007/978-3-540-78636-8>
- [14] de Cristo, F.S., Shbair, W.M., Trestioreanu, L., State, R., Malhotra, A.: Self-Sovereign Identity for the financial sector: A case study of PayString service. In: Proceedings of the 3rd International Conference on Blockchain. pp. 213–220. IEEE (2021). <https://doi.org/10.1109/Blockchain53845.2021.00036>
- [15] Cucko, S., Turkanovic, M.: Decentralized and Self-Sovereign Identity: Systematic mapping study. *IEEE Access* **9**, 139009–139027 (2021). <https://doi.org/10.1109/ACCESS.2021.3117588>
- [16] Curren, S., Looker, T., Terbu, O.: DIDComm messaging. Editor's draft, DIF: Decentralized Identity Foundation (2021), <https://identity.foundation/didcomm-messaging/spec/>
- [17] Darnell, S.S., Sevilla, J.: 3 stages of a pan-African identity framework for establishing Self-Sovereign Identity with blockchain. *Frontiers Blockchain* **4**, 631640 (2021)
- [18] Dingle, P., Hammann, S., Hardman, D., Winczewski, C., Smith, S.: Alice attempts to abuse a verifiable credential. In: White Papers from the 9th Workshop on Rebooting the Web of Trust (RWOT) (2019), <https://github.com/WebOfTrustInfo/rwot9-prague/blob/master/final-documents/alice-attempts-abuse-verifiable-credential.pdf>
- [19] Dolev, D., Yao, A.: On the security of public key protocols. *IEEE Transactions on Information Theory* **29**(2), 198–208 (1983). <https://doi.org/10.1109/TIT.1983.1056650>
- [20] Esposito, C., Horne, R., Robaldo, L., Buelens, B., Goesart, E.: Assessing the solid protocol in relation to security and privacy obligations. *Inf.* **14**(7), 411 (2023). <https://doi.org/10.3390/INFO14070411>
- [21] Freytsis, M., Barclay, I., Radha, S.K., Czajka, A., Siwo, G.H., Taylor, I.J., Bucher, S.L.: Development of a mobile, Self-Sovereign Identity approach for facility birth registration in Kenya. *Frontiers Blockchain* **4**, 631341 (2021). <https://doi.org/10.3389/fbloc.2021.631341>
- [22] Glastra, T., Aristy, G.: Aries RFC 0453: Issue credential protocol 2.0. RFC, Hyperledger Aries Community (2021), <https://github.com/hyperledger/aries-rfcs/tree/main/features/0453-issue-credential-v2>
- [23] Hogan, A.: Canonical forms for isomorphic and equivalent RDF graphs: Algorithms for leaning and labelling blank nodes. *ACM Trans. Web* **11**(4), 22:1–22:62 (2017). <https://doi.org/10.1145/3068333>
- [24] Horne, R., Mauw, S.: Discovering ePassport vulnerabilities using bisimilarity. *Logical Methods in Computer Science* **17** (2021). [https://doi.org/10.23638/LMCS-17\(2:24\)2021](https://doi.org/10.23638/LMCS-17(2:24)2021)
- [25] Horne, R., Mauw, S., Yurkov, S.: Unlinkability of an improved key agreement protocol for EMV 2nd gen payments. In: Proceedings of the 35th IEEE Computer Security Foundations Symposium (CSF). pp. 364–379 (2022). <https://doi.org/10.1109/CSF54842.2022.9919666>
- [26] Houtan, B., Hafid, A.S., Makrakis, D.: A survey on blockchain-based Self-Sovereign patient identity in healthcare. *IEEE Access* **8**, 90478–90494 (2020). <https://doi.org/10.1109/ACCESS.2020.2994090>
- [27] Khateev, N., Curran, S.: Aries RFC 0454: Present proof protocol 2.0. RFC, Hyperledger Aries Community (2021), <https://github.com/hyperledger/aries-rfcs/blob/main/features/0454-present-proof-v2/README.md>
- [28] Kudra, A.: Self-sovereign identity (SSI) in Deutschland. *Datenschutz und Datensicherheit* **46**(1), 22–26 (2022)
- [29] Lodder, M., Hardman, D.: Sovrin DID method specification. Editor's draft (2023), <https://sovrin-foundation.github.io/sovrin/spec/did-method-spec-template.html>
- [30] Longley, D., Kellogg, G., Yamamoto, D.: RDF dataset canonicalization a standard RDF dataset canonicalization algorithm. Candidate recommendation draft, W3C (2023), <https://www.w3.org/TR/rdf-canon/>
- [31] Longley, D., Sporny, M.: RDF dataset canonicalization. Final community group report, W3C (2022), <https://www.w3.org/community/reports/credentials/CG-FINAL-rdf-dataset-canonicalization-20221009/>
- [32] Lowe, G.: Breaking and fixing the Needham-Schroeder public-key protocol using FDR. *Softw. Concepts Tools* **17**(3), 93–102 (1996). [https://doi.org/10.1007/3-540-61042-1\\_43](https://doi.org/10.1007/3-540-61042-1_43)
- [33] Lowe, G.: A hierarchy of authentication specifications. In: Proceedings of the 10th IEEE Computer Security Foundations Workshop (CSFW). pp. 31–44 (1997). <https://doi.org/10.1109/CSFW.1997.596782>
- [34] Mahula, S., Tan, E., Crompvoets, J.: With blockchain or not? opportunities and challenges of Self-Sovereign Identity implementation in public administration: lessons from the Belgian case. In: Proceedings of the 22nd Annual International Conference on Digital Government Research (DG.O). pp. 495–504. ACM (2021). <https://doi.org/10.1145/3463677.3463705>
- [35] National Security Agency: Common Criteria for information technology security evaluation (CCMB-2017-04-002) (2017), <https://www.commoncriteriaportal.org/files/ccfiles/CCPART2V3.1R5.pdf>
- [36] Prorock, M., Steele, O., Terbu, O.: did:web method specification. Editor's draft (2023), <https://w3c-ccg.github.io/did-method-web/>
- [37] Sambra, A.: Verifiable credentials use cases. Working group note, W3C (2019), <https://www.w3.org/TR/vc-use-cases/>
- [38] Sambra, A.: Verifiable credentials implementation guidelines 1.0. Editor's draft, W3C (2023), <https://w3c.github.io/vc-imp-guide/>
- [39] Sporny, M., Guy, A., Sabadello, M., Reed, D.: Decentralized Identifiers (DIDs). W3C recommendation, W3C (2022), <https://www.w3.org/TR/did-core/>
- [40] Sporny, M., Longley, D., Prorock, M.: Verifiable credential data integrity 1.0: Securing the integrity of verifiable credential data. Candidate recommendation snapshot, W3C (2023), <https://www.w3.org/TR/2023/CR-vc-data-integrity-20231121/>
- [41] Sporny, M., Noble, G., Longley, D., Burnett, D.C., Zundel, B., Hartog, K.D.: Verifiable credentials data model v1.1. W3C recommendation, W3C (2022), <https://www.w3.org/TR/vc-data-model/>
- [42] Sporny, M., Zagidulin, D., Longley, D., Steele, O.: The did:key method v0.7. Unofficial draft (2022), <https://w3c-ccg.github.io/did-method-key/>
- [43] Steele, O., Sporny, M.: DID specification registries. Note, W3C DID Working Group (2023), <https://www.w3.org/TR/did-spec-registries/#did-methods>
- [44] Veramo core team: ETHR DID method specification. Editor's draft (2022), <https://github.com/decentralized-identity/ethr-did-resolver/blob/master/doc/did-method-spec.md>
- [45] West, R., Bluhm, D., Hailstone, M., Curren, S., Curran, S., Aristy, G.: Aries RFC 0023: DID exchange protocol 1.0. RFC, Hyperledger Aries Community (2021), <https://github.com/hyperledger/aries-rfcs/tree/main/features/0023-did-exchange/README.md>

## A FORMAL DEFINITION PROTOCOL ROLES, AND ELABORATION

The applied  $\pi$ -calculus processes in Table 3 and 4 are formal representations of the roles in Figure 2. Notice that once the holder has been issued a VC, unboundedly many prover sessions using the VC may be called, as represented by replication operator '!'. These process definitions are used in the Proverif and DeepSec specifications of secrecy, agreement and unlinkability properties

verified. They are assembled in various network configurations, to investigate the impact of threats described throughout Sec. 5.2.

### A.1 Security-critical ambiguities in specs

Information implementing SSI standards is scattered across supplementary material, and without elaborating on the security implications. For example, the VC Data Model specification [41] does contain examples that include security-relevant data fields (namely: domain, challenge), which are required to, e. g., prevent replay attacks. Those fields are in fact defined in a working draft on VC Data Integrity [40], where they are marked as optional. Neither the VC Data Integrity draft [40] nor the VC Implementation Guidelines [38] are sufficient to understand that not using these optional fields can result in critical authentication vulnerabilities. Elsewhere, both Aries protocols for credential issuance [22] and presentation [27] define protocol messages, where the payload of a message is an “attachment”. Definitions of such attachments also lack security considerations of the mentioned optional fields, instead building on the VC recommendation that, as we just explained, is incomplete in this sense. Thus, security controls are not clarified anywhere in this bundle of specifications which SSI is intended to rely upon.

While this layered thinking is commonplace in software engineering, problems are known to arise from a security perspective [11]. For example the DIDComm guide claims there are no possible interception attacks on DIDComm by a man-in-the-middle<sup>4</sup>; but that claim assumes an unrealistic threat model where communication is always between two honest agents. Yet the Dolev-Yao threat model is realistic for systems deployed on the Web, where dishonest or compromised agents may actively assume roles in an ecosystem [19]. That being said, the specifications including DIDComm [16] and the W3C VC data model [41] provide tools that may be assembled securely, as demonstrated in this work.

A critical threat detectable in a Dolev-Yao threat model is that, if the domain is omitted in a VP, then SSI protocols are prone to the replay attacks in Row 6 Tab. 2. The relevance of this threat is explained next, building on the example in Sec. 3.3<sup>5</sup>. A holder, e. g., a student wants to authenticate themselves towards Eve, e. g., for some student discount at some online shop. During that authentication session the holder presents to Eve a Verifiable Presentation (VP) of the student credential. This process may even complete successfully and the student may even receive their student discount. However, during this authentication process, Eve initiates a simultaneous authentication session with another verifier, e. g., a university to prove that she is the student despite not being the student. Eve uses information transmitted in messages from the university, when communicating with the student in the other session. Specifically, Eve is able to replay the challenge nonce  $n_c$  from the honest verifier (the university) to the honest holder (the student). Subsequently, this nonce is included in the VP of the student. After receiving this manipulated VP, Eve replays the VP to the university. With the matching nonce and the signature of the holder on this VP, the university believes that they were communicating with the holder the whole time, except they were communicating with

Eve masquerading as the actual holder. This fully compromises the student’s account, and hence would be a critical vulnerability.

Preliminary work applying related symbolic methods to SSI [18] analyses a 2-party DID Exchange handshake. That work considers a weak formulation of agreement and does not fully account for threats such as dishonest verifiers exploiting honest verifiers, that are fully accounted for in our Dolev-Yao network threat model.

### A.2 Novel attacks on multi-party authentication

Most definitions explored, such as *forward secrecy* and *2-party (injective) agreement* are formulated in a reasonably standard way in the repository. For agreement, the invariant that must hold in every trace representing the history of an execution, is: an occurrence of an event listing messages used by the agent performing the authentication implies the existence of an event listing all the messages sent by the agent being authenticated, and, furthermore, messages match pairwise. Forward secrecy is modelled as two phases: phase 1, before a breach where sessions run as normal, and, phase 2, after a breach when the private keys of agents are revealed and where sessions continue to run. Secrecy is only asserted about sessions that completed during phase 1.

We explain in more detail *multi-party agreement*. As explained in Sec. 5.2, a novel insight of this work is that multi-party authentication helps to explain some SSI design decisions that secrecy and two-party authentication properties miss. In particular, if the holder does not check the signature of a VC before the holder uses the VC, then there are reissuing attacks (row 8 Tab. 2). Besides the ProVerif code defining the threat model under which this attack exists, the attack vector is illustrated as an MSC in the repository<sup>6</sup>.

Multi-party agreement properties are modelled by inserting events in both the authenticator and the other parties being authenticated as follows. When the verifier is the authenticator, event `auth_VerifierCompletesProtocol`, appears after the last action of the verifier, parameterised on  $m_4, m_5, m_6, m_7, m_8$  in Tab. 4. Notice the verifier only makes assertions about the messages that it sends and receives. Also the message  $m_9$  is excluded, since  $m_9$  is sent by the verifier without expecting a response, and hence the verifier cannot check whether that message was received. The corresponding events `auth_IssuerSendsLastMessageToHolder` and `auth_ProverSendsLastMessageToVerifierInProtocolFull` appear immediately before the last message sent by the issuer and prover respectively to ensure they occur when a final message is sent (any subsequent inputs after the last output can be ignored by the same argument for excluding  $m_9$  in the event above). The event associated with the prover is parameterised on messages labelled  $m_0, m_1, m_2, m_3, m_4, m_5, m_6, m_7, m_8$  in Tab. 3 and Tab. 4. The definition of prover is extended for this property such that messages  $m_0, m_1, m_2, m_4$  are passed as additional parameters to “Prover”, in order to remember the messages that were exchanged by the “Holder” process during the issuance phase of the protocol, so they may be asserted in the relevant event. The event associated with the issuer is parameterised on its messages  $m_0, m_1, m_2, m_3$  in Tab. 3.

The authentication query (an invariant that must hold along any execution path of the protocol) is expressed in Fig. 4. Since the query

<sup>4</sup><https://didcomm.org/book/v2/mitm>

<sup>5</sup><https://github.com/uvdsl/ssi-protocol-verify/blob/main/PlainVCs/doc/msc-mitm-attack.pdf> illustrates the attack vector.

<sup>6</sup><https://github.com/uvdsl/ssi-protocol-verify/blob/main/PlainVCs/doc/msc-njagreement-attack.pdf> illustrates the attack vector.

**Table 3: Processes in the  $\pi$ -calculus for the issuance phase.**

Holder ( $P, sk_P, I, pk_I, V, pk_V$ )	Issuer ( $I, sk_I, attr, P, pk_P$ )
$new\ ssk_{PI}, n_p, n_h;$ $let\ m'_0 := (n_p, pk(ssk_{PI}))\ in$ $let\ m_0 := \{(m'_0, sig(m'_0, ssk_{PI}))\}_{pk_I}\ in$ $\overline{ch}(m_0);$ $ch(m_1);$ $let\ ((n'_p, n_i, spk_I), s_1) := adec(m_1, ssk_{PI})\ in$ $if\ check((n'_p, n_i, spk_I), s_1, pk_I)\ then$ $if\ n'_p = n_p\ then$ $let\ m'_2 := ((n_i, P, I, n_h), sig((n_i, P, I, n_h), sk_P))\ in$ $let\ m_2 := \{(m'_2, sig(m'_2, ssk_{PI}))\}_{spk_I}\ in$ $\overline{ch}(m_2);$ $ch(m_3);$ $let\ (((P', attr, I'), s_I), P'', n'_h, s_H), s_3) := adec(m_3, ssk_{PI})\ in$ $if\ check((((P', attr, I'), s_I), P'', n'_h, s_H), s_3, spk_I)\ then$ $if\ check((((P', attr, I'), s_I), P'', n'_h, s_H), spk_I)\ then$ $if\ check((P', attr, I'), s_I, pk_I)\ then$ $if\ (P', I', P'' n'_h) = (P, I, P, n_h)\ then$ $!Prover(P, sk_P, VC, V, pk_V)$	$new\ ssk_I, n_i;$ $ch(m_0);$ $let\ ((n_p, spk_{PI}), s_0) := adec(m_0, sk_I)\ in$ $if\ check((n_p, spk_{PI}), s_0, spk_{PI})\ then$ $let\ m'_1 := (n_p, n_i, pk(ssk_I))\ in$ $let\ m_1 := \{(m'_1, sig(m'_1, sk_I))\}_{spk_{PI}}\ in$ $\overline{ch}(m_1);$ $ch(m_2);$ $let\ ((n'_i, P', I', n_h), sp), s_2) := adec(m_2, ssk_I)\ in$ $if\ check(((n'_i, P', I', n_h), sp), s_2, spk_{PI})\ then$ $if\ check((n'_i, P', I'), sp, pk_P)\ then$ $if\ (n'_i, P', I') = (n_i, P, I)\ then$ $let\ claims := (P, attr, I)\ in$ $let\ VC := (claims, sig(claims, sk_I))\ in$ $let\ m'_3 := ((VC, P, n_H), sig((VC, P, n_H), sk_I))\ in$ $let\ m_3 := \{(m'_3, sig(m'_3, ssk_I))\}_{spk_{PI}}\ in$ $\overline{ch}(m_3);$

**Table 4: Processes in the  $\pi$ -calculus for the provenance phase.**

Prover ( $P, sk_P, VC, V, pk_V$ )	Verifier ( $V, sk_V, RULE, pk_P, pk_I, URI$ )
$new\ ssk_{PV}, n_p;$ $let\ m'_4 := (n_p, pk(ssk_{PV}))\ in$ $let\ m_4 := \{(m'_4, sig(m'_4, ssk_{PV}))\}_{pk_V}\ in$ $\overline{ch}(m_4);$ $ch(m_5);$ $let\ ((n'_p, n_v, spk_V), s_5) := adec(m_5, ssk_{PV})\ in$ $if\ check((n'_p, n_v, spk_V), s_5, pk_V)\ then$ $if\ n'_p := n_p\ then$ $let\ m'_6 := (n_v, URI)\ in$ $let\ m_6 := \{(m'_6, sig(m'_6, ssk_{PV}))\}_{spk_V}\ in$ $\overline{ch}(m_6);$ $ch(m_7);$ $let\ ((n_c, RULE), s_7) := adec(m_7, ssk_{PV})\ in$ $if\ check((n_c, RULE), s_7, spk_V)\ then$ $let\ (claims, s_I) := VC\ in$ $if\ claims = RULE\ then$ $let\ VP := ((VC, n_c, V), sig((VC, n_c, V), sk_P))\ in$ $let\ m_8 := \{VP, sig(VP, ssk_{PV})\}_{spk_V}\ in$ $\overline{ch}(m_8);$ $ch(m_9);$ $let\ ((tkn, s_{tkn}), s_9) := (adec(m_9, ssk), spk_V)\ in$ $if\ check((tkn, s_{tkn}), s_9, spk_V)\ then$ $if\ check(tkn, s_{tkn}, pk_V)\ then$	$new\ ssk_V, n_i, n_c, tkn;$ $ch(m_4);$ $let\ ((n_p, spk_{PV}), s_4) := adec(m_4, sk_V)\ in$ $if\ check((n_p, spk_{PV}), s_4, spk_{PV})\ in$ $let\ m'_5 := (n_p, n_v, pk(ssk_V))\ in$ $let\ m_5 := \{(m'_5, sig(m'_5, sk_V))\}_{spk_{PV}}\ in$ $\overline{ch}(m_5);$ $ch(m_6);$ $let\ ((n'_v, uri'), s_6) := adec(m_6, sk_V)\ in$ $if\ check((n'_v, uri'), s_6, spk_{PV})\ then$ $if\ (n'_v, URI') = (n_v, URI)\ then$ $let\ m'_7 := (n_c, RULE)\ in$ $let\ m_7 := \{(m'_7, sig(m'_7, ssk_V))\}_{spk_{PV}}\ in$ $\overline{ch}(m_7);$ $ch(m_8);$ $let\ (((P', attr', I'), s_I), n'_c, V'), sp), s_8) := adec(m_8, ssk_V)\ in$ $if\ check((((P', attr', I'), s_I), n'_c, V'), sp), s_8, spk_{PV})\ then$ $if\ check((((P', attr', I'), s_I), n'_c, V'), sp, pk_P)\ then$ $if\ check((P', attr', I'), s_I, pk_I)\ then$ $if\ ((P', attr', I'), n'_c, V') = ((P, attr, I), n_c, V)\ then$ $let\ m'_9 := (tkn, sig(tkn, sk_V))\ in$ $let\ m_9 := \{sig(m'_9, ssk_V)\}_{spk_{PV}}\ in$ $\overline{ch}(m_9);$

$\forall m_4, m_5, m_6, m_7, m_8. \ event(auth\_VerifierCompletesProtocol(m_4, m_5, m_6, m_7, m_8))$   
 $\Rightarrow \exists m_0, m_1, m_2, m_3. \ event(auth\_IssuerSendsLastMessageToHolder(m_0, m_1, m_2, m_3))$   
 $\wedge \ event(auth\_ProverSendsLastMessageToVerifierInProtocolFull(m_0, m_1, m_2, m_3, m_4, m_5, m_6, m_7, m_8))$

**Figure 4: Invariant expressing multi-party authentication by the verifier of both the prover and issuer.**

is about the perspective of the verifier as an authenticator, messages observed by the verifier are universally quantified, while messages unknown to the verifier, which are observed by the holder during issuance interacting with the verifier, are existentially quantified. Notice that, although there is no interaction between the issuer and verifier, the query ensures that the VC appearing inside the message  $m_8$  of the verifier matches the VC inside the message  $m_8$  of the prover, and hence matches the VC inside message  $m_3$  of the Holder preceding the prover and hence matches the VC inside message  $m_3$  of the issuer. Therefore, by transitivity, the verifier and issuer indirectly agree on a specific VC.

### A.3 Formulating unlinkability v.s. the issuer

The novel formulations of unlinkability towards the issuer are also formal contributions of this paper (lines 9, 10, 13 of Tab. 2). It is commonplace when symbolically verifying protocols to express unlinkability as an equivalence problem between a process modelling an idealised system that is trivially unlinkable by definition and another process modelling more realistic behaviours where the same identities are used across multiple sessions [4, 10, 24, 25].

To model the unlinkability of an honest prover and verifier in the presence of a malicious (or overly curious) issuer, our trick in setting up the trust model is to model only honest verifiers and provers who interact with each other and also to expose the secret keys of the issuer. This gives the issuer full power to behave as a Dolev-Yao attacker, attempting to manipulate sessions between honest participants. An open variable for the issuer's private key models the assumption that the attacker has the private keys of the issuer.

A rationale for our trust model is that, if the issuer were able to exploit the protocol to determine whether an honest prover has used the same VC that the attacker issues in two provenance sessions with honest verifiers, then the issuer would be able to exploit its position in the network to track the prover. One may place a counter argument that an issuer will likely be "honest but curious" and can be modelled with less capabilities than a full Dolev-Yao attacker; yet, this argument is irrelevant since proofs go through and hence, no matter how devious the issuer is, they will be unable to track participants in honest sessions of the provenance phase (with the exception of sessions in which a malicious issuer poses as a verifier in the session themselves, of course). Perhaps counter-intuitively, this formal property is advantageous to the issuer: If accused of abusing their knowledge to track the VCs they issue, that claim may be countered by the issuer arguing that the protocol makes such tracking impossible even if a sophisticated devious attacker were to assist the issuer.

In contrast to the secrecy and unlinkability problems, which reason over infinitely many session, we restrict this analysis to two sessions, so that the formulation of the problem is amenable to the bounded equivalence checker DeepSec. An applied  $\pi$ -calculus processes modelling the idealised and real-world scenarios that should be equivalent appears in the relevant DeepSec files in the repository. Both the idealised process and real-world process begin with a preamble defining the secret keys of the honest provers and verifiers as follows, and releasing the public keys (or DIDs

containing a public key) to the network.

```
new sk_prover1, sk_prover2, sk_verifier;
let pk_prover1 = pk(sk_prover1) in key(pk_prover1);
let pk_prover2 = pk(sk_prover2) in key(pk_prover2);
let pk_verifier = pk(sk_verifier) in key(pk_verifier);
```

The processes initiated consist of three parallel threads. In both the specification and real-world scenarios, there are two parallel honest verifiers, who are prepared to engage in a session with one of two honest provers which correspond to the public keys advertised above. These verifiers are parameterised as follows.

```
Verifier( DID_verifier, sk_verifier, attr,
          pk_prover1, pk_prover2, pk_issuer, URI )
```

The above is a mild variant of the verifier processes defined in Tab. 4, where the public keys of two provers  $pk\_prover1$  and  $pk\_prover2$  are both accepted by the verifier when checking the signature on the VP. Parameters such as  $URI$  and  $attr$  are open variables, thus may be publicly known (and even manipulated by the attacker).

The system and real-world processes differ in how the honest holder is modified. Both begin as specified by the "Holder" process in Tab. 3, parameterised on a private key of the holder,  $pk\_prover1$ , and a public key based on the secret key of the issuer known to the attacker,  $pk(sk\_issuer)$ . The holder is therefore prepared to receive a VC issued by an attacker and interact with an honest verifier.

In the real world, once the VC is issued, the holder continues much as in Tab. 3 by starting two prover sessions loaded with the VC that has just been issued and the public keys of the honest verifier. This models the holder using the same VC twice in different sessions, i. e. an expected usage pattern. In contrast, in the idealised process, the holder is modified such that, after having being issued a VC it is prepared to engage in two prover sessions that employ two different fresh VCs with the relevant attributes, as if issued by the attacker, rather than the VC that was just issued to the holder. If the attacker cannot distinguish this setup from the above real-world usage pattern, then not only can the attacker not tell whether or not the same credential was used twice, but it also cannot tell whether a particular credential was used at all. The presence of a fresh prover further verifies that the identity of any prover involved in a provenance session is also not revealed to the attacker.

In order to strengthen the threat model in the setting of anonymous credentials, the verifier is dropped from the process modelling the real and idealised worlds, and the variable  $sk\_verifier$  is turned into an open free variable, indicating that the attacker may know that variable (and also manipulate it, e.g., by correlating the secret key of the issuer and verifier). This strengthens the threat model by assuming that the verifier may also be an attacker, and furthermore the issuer and verifier may attempt to collude to trace the holder of a credential. Our verification of that model in DeepSec shows that anonymous credentials are not vulnerable to attacks on unlinkability in the face of this threat.

The above explanations and discussions highlight that this paper makes a novel contribution to the symbolic verification of security and privacy properties of protocols, as well as applying appropriate established methodologies to evaluate the security and privacy of VC protocols.