

# Integrating Automated and Interactive Protocol Verification

Achim D. Brucker<sup>1</sup> and Sebastian A. Mödersheim<sup>2</sup>

<sup>1</sup> SAP Research, Vincenz-Priessnitz-Str. 1, 76131 Karlsruhe, Germany  
[achim.brucker@sap.com](mailto:achim.brucker@sap.com)

<sup>2</sup> IBM Research, Säumerstrasse 4, 8803 Rüschlikon, Switzerland  
[smo@zurich.ibm.com](mailto:smo@zurich.ibm.com)

**Abstract.** A number of current automated protocol verification tools are based on abstract interpretation techniques and other over-approximations of the set of reachable states or traces. The protocol models that these tools employ are shaped by the needs of automated verification and require subtle assumptions. Also, a complex verification tool may suffer from implementation bugs so that in the worst case the tool could accept some incorrect protocols as being correct. These risks of errors are also present, but considerably smaller, when using an LCF-style theorem prover like Isabelle. The interactive security proof, however, requires a lot of expertise and time.

We combine the advantages of both worlds by using the representation of the over-approximated search space computed by the automated tools as a “proof idea” in Isabelle. Thus, we devise proof tactics for Isabelle that generate the correctness proof of the protocol from the output of the automated tools. In the worst case, these tactics fail to construct a proof, namely when the representation of the search space is for some reason incorrect. However, when they succeed, the correctness only relies on the basic model and the Isabelle core.

## 1 Introduction

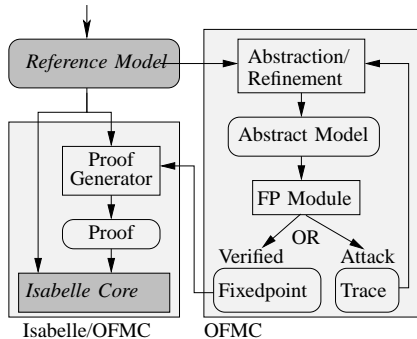
Over the last decade, a number of automated tools for security protocol verification have been developed such as AVISPA [1] and ProVerif [4]. They allow engineers to find problems in their security protocols before deployment. Indeed, several attacks to security protocols have been detected using automated tools. The focus of this work is the positive case—when no attack is found: to obtain a proof of security.

Many automated tools employ over-approximation and abstraction techniques to cope with the infinite search spaces that are caused, e.g., by an unbounded number of protocol sessions. This means to check the protocol in a finite abstract model that subsumes the original model. Thus, if the protocol is correct in the abstract model, then so it is in the original model. However, the soundness of such abstractions depends on subtle assumptions, and it is often hard to keep track of them, even for experts. Moreover, it is often hard to formalize protocols correctly in such over-approximated models. Finally, tools may

also have bugs. For all these reasons, it is not unlikely that insecure protocols are accidentally verified by automated verification tools.

There are semi-automated methods such as the Isabelle theorem prover which offer a high reliability: if we trust in a small core (the proof checking and some basic logical axioms), we can rely on the correctness of proved statements. However, conducting proofs in Isabelle requires considerable experience in both formal logic and proof tactics, as well as a proof idea for the statement to show.

In this work, we combine the best of both worlds: reliability and full automation. The idea is that abstraction-based tools supposedly compute a finite representation of an over-approximated search space, i.e., of what can happen in a given protocol, and that this representation can be used as the basis to automatically generate a security proof in Isabelle. This proof is w.r.t. a clean standard protocol model without over-approximation. If anything goes wrong, e.g., the abstraction is not sound, this proof generation fails. However, if we succeed in generating the proof, we only need to trust in the standard protocol model and the Isabelle core. Our vision is that such automatically generated verifiable proofs can be the basis for reaching the highest assurance level EAL7 of a common criteria certification at a low cost.



**Fig. 1.** The workflow of a protocol verification approach combining automated (e.g., OFMC) and interactive (e.g., Isabelle) techniques.

abstracted version of the protocol description as a set of Horn clauses. From this, the fixed-point (FP) module computes a least fixed-point of derivable events. If this fixed-point contains an attack, this can either be a real attack (that similarly works in the reference model) or a false attack that was caused by the abstraction. By default, OFMC will assume that the attack is false, refine the abstraction based on the attack details and restart the verification. If the computed fixed-point does not contain an attack, it is handed to the proof generator of the Isabelle/OFCM, our extension of the interactive theorem prover Isabelle [25]. The proof generator translates the fixed-point into the terms of the reference model (using annotations about the abstraction OFMC considered)

We have realized the integration of automated and interactive protocol verification in a prototype tool that is summarized in Fig. 1. The protocol and the proof goals (e.g., secrecy, authentication) are specified in the *reference model* (Sect. 2 and Sect. 3). This reference model is not driven by technical needs of the automated verification and is close to other high-level protocol models. The description is fed into the automated tool, where we consider (a novel module of) the Open-source Fixed-point Model-Checker OFMC [24], formerly called On-the-Fly Model-Checker. OFMC first chooses an initial abstraction and produces an abstracted version of the protocol description as a set of Horn clauses. From this, the fixed-point (FP) module computes a least fixed-point of derivable events. If this fixed-point contains an attack, this can either be a real attack (that similarly works in the reference model) or a false attack that was caused by the abstraction. By default, OFMC will assume that the attack is false, refine the abstraction based on the attack details and restart the verification. If the computed fixed-point does not contain an attack, it is handed to the proof generator of the Isabelle/OFCM, our extension of the interactive theorem prover Isabelle [25]. The proof generator translates the fixed-point into the terms of the reference model (using annotations about the abstraction OFMC considered)

and generates an Isabelle proof with respect to the protocol and goal description in the reference model. This proof is fed into the Isabelle core to check the proof.

We emphasize two points. First, the entire approach is completely automatic: after the specification of the protocol and goals in the reference model, no further user interaction is required. Second, we need to trust in only two points—marked with a dark-gray background in Fig. 1: the reference model and the Isabelle core. Bugs in any other part, namely OFMC or the proof generation, can in the worst case result in a failure to verify a (possibly correct) protocol, but they cannot make us falsely accept a flawed protocol as being correct.

Our approach currently has two limitations: we consider a typed model and, based on this, we limit the intruder composition to well-typed messages of the protocol. The first limitation can be justified by implementation discipline (see [19]). The second limitation is not a restriction as we show in Theorem 1.

*Contributions.* An increasing number of works considers the combination of automated methods with interactive theorem proving to obtain both highly reliable and fully automated verification. In this paper, we contribute to this line of work with a novel approach for security protocols. The first novel aspect is that our approach automatically generates a proof from the representation of an over-approximation of the search space computed by an automated protocol verifier. Techniques based on over-approximation, similar to the ones we consider, have turned out to be very successful in protocol verification [4,6,7,10], and our approach is thus the first step towards employing a whole class of established tools for automated proof generation. The second novel aspect is that the proof is entirely based on a standard protocol model without over-approximation close to the model employed for instance in [26]. Our approach thus relates over-approximated representations with standard protocol models.

Practically, we have implemented the integration between Isabelle on the interactive side and the novel FP-module of OFMC on the automated side. The result is a completely automated protocol verifier for an unbounded number of sessions and agents that produces Isabelle-verifiable proofs with respect to a standard protocol model.

## 2 The Reference Protocol Model

We begin with a reference protocol model, which is used in the Isabelle theorem prover and is thus the basis of this work. The model is inspired by the formalization of several security protocols in Isabelle by Paulson and others in [26,3] and is close to the persistent IF model in [23].

*Messages.* We follow the common black-box cryptography model where messages are modeled as symbolic terms. We define the set of all messages in style of an inductive datatype in a functional programming language. Here, all **sans-serif** symbols and the symbols of  $\mathcal{F}$  are *constructors*. The inductively defined datatype is interpreted like a term in the free term algebra, i.e. syntactically different terms are interpreted as being different.

**Definition 1.** Let  $\mathcal{F}$ ,  $\mathcal{L}_A$ ,  $\mathcal{L}_N$ ,  $\mathcal{L}_S$ ,  $\mathcal{L}_P$ ,  $\mathcal{V}_A$ ,  $\mathcal{V}_N$ ,  $\mathcal{V}_S$ ,  $\mathcal{V}_P$ , and  $\mathcal{V}_U$  be pairwise disjoint sets of symbols where  $\mathcal{F}$ ,  $\mathcal{L}_A$ ,  $\mathcal{L}_N$ ,  $\mathcal{L}_S$ , and  $\mathcal{L}_P$  are finite, and  $\mathcal{V}_A$ ,  $\mathcal{V}_N$ ,  $\mathcal{V}_S$ ,  $\mathcal{V}_P$ , and  $\mathcal{V}_U$  are countable. We define the sets of messages, agents, nonces, symmetric keys, and public keys, respectively, to be the following sets of terms:

$$\begin{aligned} \mathcal{M} &= \text{agent } \mathcal{A} \mid \text{nonce } \mathcal{N} \mid \text{symkey } \mathcal{S} \mid \text{pubkey } \mathcal{P} \mid \mathcal{V}_U \\ &\quad \mid \text{crypt } \mathcal{M} \ \mathcal{M} \mid \text{inv } \mathcal{M} \mid \text{script } \mathcal{M} \ \mathcal{M} \mid \text{cat } [\mathcal{M}] \mid \mathcal{F} \ \mathcal{M} \\ \mathcal{A} &= \mathcal{L}_A \times \mathbb{N} \mid \mathcal{V}_A \quad \mathcal{N} = \mathcal{L}_N \times \mathbb{N} \mid \mathcal{V}_N \quad \mathcal{S} = \mathcal{L}_S \times \mathbb{N} \mid \mathcal{V}_S \quad \mathcal{P} = \mathcal{L}_P \times \mathbb{N} \mid \mathcal{V}_P \end{aligned}$$

The set  $\mathcal{A}$  contains both the *concrete* agent names ( $\mathcal{L}_A \times \mathbb{N}$ ) and the *variables* for agent names ( $\mathcal{V}_A$ ). The concrete agent names consist of a label and a natural number. The labels are for the interplay of Isabelle and the automated methods, for now it is sufficient to think just of an infinite set of agents, indexed by natural numbers. Similarly,  $\mathcal{N}$ ,  $\mathcal{S}$ , and  $\mathcal{P}$  define infinite reservoirs of concrete constants and variables for nonces, symmetric keys, and public keys.

For convenience, we write in the examples of this paper simply  $a$ ,  $b$ , or  $i$  for concrete agent names,  $A$ ,  $B$  for agent variables,  $n$ ,  $n_1$ , etc. for concrete nonces,  $N_A$ ,  $N_B$  etc. nonce variables. In general, we use lower-case letters for constants and function symbols and upper-case letters for variables.

We distinguish atomic messages and composed messages (first and second line in the definition of  $\mathcal{M}$ ). Except for the *untyped variables* of  $\mathcal{V}_U$ , all atomic messages are of a particular type, namely from one of the sets  $\mathcal{A}$ ,  $\mathcal{N}$ ,  $\mathcal{S}$ , or  $\mathcal{P}$ . The constructors like `agent` ensure that the respective subsets of the message space are disjoint, for instance, no agent name can be a nonce. We discuss the details of typing below. In examples, we will omit the constructors for convenience, when the type is clear from the context.

Messages can be composed with one of the following operations: `crypt` and `script` represent asymmetric and symmetric encryption, respectively. We also simply write  $\{m\}_k$  for `crypt`  $k$   $m$ , and  $\{\{m\}\}_k$  for `script`  $k$   $m$ . `inv`( $M$ ) represents the private key belonging to a public key. `cat` denotes concatenation. For readability, we omit the `cat` and the list brackets. For instance the term  $\{na_3, b\}_{pk(a)}$  is convenient notation for the following more technical message (for given labeling and numbering): `crypt (pk (agent (honest, 1))) (cat [nonce (na, 3), agent (dishonest, 2)])`. Here we have used two labels, `honest` and `dishonest`, for agents. This represents the default abstraction for agents in the abstract model. We require that the abstraction for the agents is a refinement of the default abstraction. We use this labeling to distinguish honest and dishonest agents also in the reference model.

We use the standard notion of matching messages. The constructors like `agent` here enforce a typing regime: typed variables can only be matched with atomic messages of the same type. Only untyped variables can be matched with composed messages. In rules for honest agents, we only use typed variables. Such a typed model which is standard in protocol verification, even in interactive verification with Isabelle [26,3], considerably simplifies the verification task. The typing can be justified by tagging as in [19].

*Events and Traces.* We define the set of events also as an inductive datatype, based on messages:

$$\begin{array}{c}
\frac{t \in \mathbb{T} \quad \text{iknows } m \in [t] \quad \text{iknows } k \in [t]}{\text{iknows } \{m\}_k \# t \in \mathbb{T}} \quad \frac{t \in \mathbb{T} \quad \text{iknows } m \in [t] \quad \text{iknows } k \in [t]}{\text{iknows } \{m\}_k \# t \in \mathbb{T}} \\
\frac{t \in \mathbb{T} \quad \text{iknows } m_1 \in [t] \quad \dots \quad \text{iknows } m_n \in [t]}{\text{iknows } [m_1, \dots, m_n] \# t \in \mathbb{T}} \quad \frac{t \in \mathbb{T} \quad \text{iknows } m \in [t]}{\text{iknows } f(m) \# t \in \mathbb{T}} \quad f \in \mathcal{F}_{pub} \\
\frac{t \in \mathbb{T} \quad \text{iknows } \{m\}_k \in [t] \quad \text{iknows } \text{inv}(k) \in [t]}{\text{iknows } m \# t \in \mathbb{T}} \\
\frac{t \in \mathbb{T} \quad \text{iknows } \{m\}_{\text{inv}(k)} \in [t] \quad \text{iknows } k \in [t]}{\text{iknows } m \# t \in \mathbb{T}} \\
\frac{t \in \mathbb{T} \quad \text{iknows } \{m\}_k \in [t] \quad \text{iknows } k \in [t]}{\text{iknows } m \# t \in \mathbb{T}} \\
\frac{t \in \mathbb{T} \quad \text{iknows } m_1 \in [t] \quad \dots \quad \text{iknows } m_n \in [t]}{\text{iknows } m_1 \#, \dots \# \text{iknows } m_n \# t \in \mathbb{T}} \\
\frac{t \in \mathbb{T} \quad \text{secret } A \ M \in [t] \quad \text{iknows } M \in [t] \quad \text{honest } A}{\text{attack } M \# t \in \mathbb{T}} \\
\frac{t \in \mathbb{T} \quad \text{request } B \ A \ \text{id } M \in [t] \quad \text{witness } A \ B \ \text{id } M \notin [t] \quad \text{honest } A}{\text{attack } M \# t \in \mathbb{T}}
\end{array}$$

**Fig. 2.** The protocol independent rules of the reference model: the first four are *composition rules (C)*, the next four are *decomposition rules (D)* and the last two are *attack rules (A)*.

**Definition 2.** *The set of events is defined as follows:*

$$\begin{aligned}
\mathcal{E} ::= & \text{iknows } \mathcal{M} \mid \text{state } \mathcal{R} \ [ \mathcal{M} ] \mid \\
& \text{secret } \mathcal{A} \ \mathcal{M} \mid \text{witness } \mathcal{A} \ \mathcal{A} \ \mathcal{I} \ \mathcal{M} \mid \text{request } \mathcal{A} \ \mathcal{A} \ \mathcal{I} \ \mathcal{M} \mid \text{attack } \mathcal{M}
\end{aligned}$$

where  $\mathcal{I}$  is a finite set of identifiers disjoint from all other symbols so far and  $\mathcal{R} \subset \mathcal{I}$ . A trace is a finite sequence  $\langle e_1 \# \dots \# e_n \rangle$  of events  $e_i$ .

The identifier set  $\mathcal{I}$  contains constant symbols for the protocol variables, allowing us to describe as which protocol variable an agent interprets a particular message. We use **Gothic fonts** for identifiers, e.g.  $\mathfrak{A}$  and  $\mathfrak{B}$ .

The event  $\text{iknows } m$  means that the intruder just learned the message  $m$ . The event  $\text{state } \mathfrak{R} \ \text{msgs}$  means that an honest agent playing role  $\mathfrak{R}$  has reached a state of its protocol execution that is characterized by the list  $\text{msgs}$  of messages. We need the other four events for expressing the goals in a protocol independent form when we introduce attack rules below.

*Rules and Protocols.* Based on these definitions, we formalize protocols by a set of inductive rules on traces that have the following form:

$$\frac{t \in \mathbb{T} \quad \phi(t, e_1, \dots, e_n)}{e_1 \# \dots \# e_n \# t \in \mathbb{T}}$$

i.e. whenever  $t$  is a valid trace of the set  $\mathbb{T}$  of traces and  $e_1, \dots, e_n$  are events that fulfill a certain condition  $\phi$  with  $t$ , then also the extension of  $t$  with these events is also part of  $\mathbb{T}$ . Also, we have the rule that the empty trace is part of  $\mathbb{T}$ . Note that we require that all transition rules of honest agents contain only typed variables, i.e. no variables of  $\mathcal{V}_{\mathcal{U}}$ .

Fig. 2 shows the protocol independent rules for the intruder (C) and (D), following the standard Dolev-Yao style intruder deduction, as well as the attack rules (A), which follow the standard definitions of attacks in AVISPA [1]; here  $[t]$  denotes the set of events in the trace  $t$  and  $\mathcal{F}_{pub} \subseteq \mathcal{F}$  is the set of functions that is accessible to the intruder. Before we explain the attack rules, we describe the transition rules of role Alice of the standard example protocol, NSL [20] (more interesting examples are found in Sect. 6):

$$\begin{array}{c} \text{Example 1.} \\ \frac{t \in \mathbb{T} \quad NA \notin \text{used}(t)}{\text{iknows } \{NA, A\}_{pk(B)} \# \text{state } \mathfrak{A} [A, B, NA] \# \\ \text{witness } A B \mathfrak{N}_{\mathfrak{A}} NA \# \text{secret } B NA \# t \in \mathbb{T}} \\ \\ \frac{t \in \mathbb{T} \quad \text{state } \mathfrak{A} [A, B, NA] \in [t] \quad \text{iknows } \{NA, NB, B\}_{pk(A)} \in [t]}{\text{iknows } \{NB\}_{pk(B)} \# \text{request } A B \mathfrak{N}_{\mathfrak{B}} NB \# t \in \mathbb{T}} \end{array}$$

Here,  $\text{used}(t)$  is the set of all atomic messages that occur in  $t$  to allow for the fresh generation of nonces.  $\square$

The goals of a protocol are described negatively by what counts as an attack. This is done by *attack rules* that have the event **attack** on the right-hand side. We now explain the events that we use in attack rules. First, **secret**  $A M$  means that some honest agent (not specified) requires that the message  $M$  is a secret with agent  $A$ . Thus, it counts as an attack, if a trace contains both **iknows**  $M$  and **secret**  $A M$  for an honest agent  $A$  (first attack rule in Fig. 2).

For authentication, the event **witness**  $A B id m$  means that for a particular purpose  $id$ , the honest agent  $A$  wants to transmit the message  $M$  to agent  $B$ . Correspondingly, when  $B$  believes to have received message  $M$  for purpose  $id$  from agent  $A$ , the event **request**  $B A id M$  occurs. It thus counts as an attack, if **request**  $B A id M$  occurs in a trace for an honest agent  $A$  and the trace does not contain the corresponding event **witness**  $A B id M$  (second attack rule in Fig. 2). We call a trace an *attack trace* if it contains the **attack** event.

**Definition 3.** *Let a protocol be described by an inductive set  $R$  of rules. The protocol is said to be safe, if the least set  $\mathbb{T}$  of traces that is closed under  $R$  contains no attack trace.*

Despite some differences, our model is similar to the one of Paulson [26], as discussed in detail in the extended version of this paper [8].

### 3 Limiting Intruder Composition

The closure of the intruder knowledge under the composition rules of the intruder is generally infinite, e.g. the intruder can concatenate known messages arbitrarily. However, many of these messages are useless to the intruder since, due to typing, no honest agent accepts them. It is therefore intuitive that we do not lose any attacks if we limit intruder composition to terms, and subterms thereof, that some honest agent can actually receive. (A similar idea has been considered e.g. in [28].) This limitation on intruder composition makes our approach significantly simpler, and we have therefore chosen to integrate this simplification into our reference model for this first version, and leave the generalization to an unlimited intruder for future versions. We formally define the transformation that limits intruder composition as follows:

**Definition 4.** *For a set of rules  $R$  that contain no untyped variables  $\mathcal{V}_U$ , let  $M_R$  be the set of all messages that occur in an `iknows` or `secret` event of any rule of  $R$  along with as their sub-messages. We say that an intruder composition rule  $r$  can compose terms for  $R$ , if the resulting term of  $r$  can be unified with a term in  $M_R$ . In this case we call  $r\sigma$  an  $r$ -instance for compositions of  $M_R$  if  $\sigma$  is the most general unifier between the resulting term of  $r$  and a message in  $M_R$ . We say that  $R$  is saturated if it contains all  $r$ -instances for composition in  $M_R$ .*

Since we excluded untyped variables, atomic messages in  $M_R$  are typed, i.e. of the form `type(·)`. Also, due to typing, every finite  $R$  has a finite saturation.

**Theorem 1.** *Given an attack against a protocol described by a set of rules  $R \cup C \cup D \cup A$  where  $C$  and  $D$  are the intruder composition and decomposition rules and  $A$  are the attack rules. Let  $R'$  be a saturated superset of  $R$ . Then there is an attack against  $R' \cup D \cup A$ .*

The proof is found in the extended version of this paper [8].

### 4 The Abstract Protocol Model

We now summarize two kinds of over-approximations of our model that are used in our automated analysis tool to cope with the infinite set of traces induced by the reference model. These techniques are quite common in protocol verification and a more detailed description can be found in [23]; we discuss them here only as far as they are relevant for our generation of Isabelle proofs. The first technique is a data abstraction that maps the infinite set of ground atomic messages (that can be created in an unbounded number of sessions) to a finite set of equivalence classes in the style of abstract interpretation approaches. The second is a control abstraction: we forget about the structure of traces and just consider *reachable events* in the sense that they are contained in some trace of the reference model. Neither of these abstractions is safe: each may introduce false attacks (that are not possible in the reference model). Also, it is not guaranteed in general that the model allows only for a finite number of reachable events, i.e. the approach may still run into non-termination.

*Data Abstraction.* In the style of abstract interpretation, we first partition the set of all ground atomic messages into finitely many equivalence classes and then work on the basis of these equivalence classes. Recall that atomic ground messages are defined as a pair  $(l, n)$  where  $l$  is a label and  $n$  is a natural number. We use the label to denote the equivalence class of the message in the abstraction. The abstract model thus identifies different atoms with the same label, and hence we just omit the second component in all messages of the abstract model.

There is a large variety of such data abstractions. For the proof generation the concrete way of abstraction is actually irrelevant, and we just give one example for illustration:

*Example 2.* The initial abstraction that OFMC uses for freshly created data is the following. If agent  $a$  creates a nonce  $n$  for agent  $b$ , we characterize the equivalence class for that nonce in the abstract model by the triple  $(n, a, b)$ . This abstraction can be rephrased as “in all sessions where  $a$  wants to talk to  $b$ ,  $a$  uses the same constant for  $n$  (instead of a really fresh one)”. For a large number of cases, this simple abstraction is sufficient; in the experiments of Sect. 6, only two examples (NSL and Non-reversible Functions) require a more fine-grained abstraction. Note that OFMC automatically refines abstractions when the verification fails, but this mechanism is irrelevant for the proof generation. The protocol from Example 1, rules for  $A$ , then look as follows:

$$\frac{t \in \mathbb{T}}{\text{iknows } \{(\mathfrak{N}_{\mathfrak{A}}, A, B), A\}_{\text{pk}(B)} \# \text{state } \mathfrak{A} [A, B, (\mathfrak{N}_{\mathfrak{A}}, A, B)] \# \text{witness } A B \mathfrak{N}_{\mathfrak{A}} (\mathfrak{N}_{\mathfrak{A}}, A, B) \# \text{secret } B (\mathfrak{N}_{\mathfrak{A}}, A, B), \# t \in \mathbb{T}}$$

$$\frac{t \in \mathbb{T} \quad \text{state } \mathfrak{A} [A, B, NA] \in [t] \quad \text{iknows } \{NA, NB, B\}_{\text{pk}(A)} \in [t]}{\text{iknows } \{NB\}_{\text{pk}(B)} \# \text{request } A B \mathfrak{N}_{\mathfrak{B}} NB \# t \in \mathbb{T}}$$

Here, we only abstract  $NA$  in the first rule where it is created by the  $A$ . It is crucial that the nonce  $NB$  is not abstracted in the second rule: since  $NB$  is not generated by  $A$ ,  $A$  cannot be sure a priori that was indeed generated by  $B$ . In fact, if we also abstract  $NB$  here, the proof generation fails, because the resulting fixed-point does no longer over-approximate the traces of the reference model. More generally, fresh data are abstracted only in a rule where they are created. Finally, observe that the condition is gone that tells that the freshly created  $NA$  never occurred in the trace before, because now agents may actually use the same value several times in place of the fresh nonce.

The key idea to relate our reference model with the abstract one is the use of labels in the definition of concrete data. Recall that each concrete atomic message in the reference model is a pair of a label and a natural number. The finite set of labels is determined by the abstraction we use in the abstracted model; in the above example, we use  $\mathcal{L}_{\mathcal{N}} = \{\mathfrak{N}_{\mathfrak{A}}, \mathfrak{N}_{\mathfrak{B}}\} \times \mathcal{L}_{\mathcal{A}} \times \mathcal{L}_{\mathcal{A}}$  where  $\mathcal{L}_{\mathcal{A}}$  is the abstraction of the agents (for instance  $\mathcal{L}_{\mathcal{A}} = \{(\text{honest}, \text{dishonest})\}$ ). As the atomic messages consist of both such a label and a natural number, the reference model is thus endowed with an infinite supply of constants for each equivalence



class of the abstract model. The relationship between data in the reference and abstract models is straightforward: the abstraction of the concrete constant  $(l, n)$  is simply  $l$ . Vice-versa, each equivalence class  $l$  in the abstract model represents the set of data  $\{(l, n) \mid n \in \mathbb{N}\}$  in the reference model.

It is crucial that in the reference model, the labels are merely annotations to the data and the rules do not care about these annotations, except for the distinction of honest and dishonest agents as discussed before. The labels however later allow us to form a security proof in the reference model based on the reachable events in the concrete model.

We need to take the abstraction into account in the reference model when creating fresh data. In particular, we need to enforce the labeling that reflects exactly the abstraction. We extend the assumptions of the first rule from Example 1 by a condition on the label of the freshly created  $NA$ :

$$\frac{t \in \mathbb{T} \quad NA \notin \text{used}(t) \quad \text{label}(NA) = (\mathfrak{N}_{\mathfrak{A}}, A, B)}{\text{iknows } \{NA, A\}_{\text{pk}(B)}, \text{state } \mathfrak{A} [A, B, NA] \# \text{witness } A B \mathfrak{N}_{\mathfrak{A}} NA \# \text{secret } B NA \# t \in \mathbb{T}}$$

where  $\text{label}(l, n) = l$ . (Recall that every concrete value is a pair of label and a natural number.)

*Control Abstraction.* We now come to the second part of the abstraction. Even with the first abstraction on data, the model gives us an infinite number of traces (that are of finite but of unbounded length). The idea for simplification is that under the data-abstraction, the trace structure is usually not relevant anymore. In the reference model, we need the trace structure for the creation of fresh data and for distinguishing potentially different handling of the same constant in different traces. Under the data-abstraction, however, all these occurrences fall together. In the abstract model we thus abandon the notion of traces and consider only the set  $\mathbb{E}$  of events that can ever occur.

*Example 3.* Our running example has now the following form:

$$\frac{\text{iknows } \{(\mathfrak{N}_{\mathfrak{A}}, A, B), A\}_{\text{pk}(B)}, \text{state } \mathfrak{A} [A, B, (\mathfrak{N}_{\mathfrak{A}}, A, B)], \text{witness } A B \mathfrak{N}_{\mathfrak{A}} (\mathfrak{N}_{\mathfrak{A}}, A, B), \text{secret } B (\mathfrak{N}_{\mathfrak{A}}, A, B) \in \mathbb{E}}{\text{state } \mathfrak{A} [A, B, NA] \in \mathbb{E} \quad \text{iknows } \{NA, NB, B\}_{\text{pk}(A)} \in \mathbb{E}} \frac{}{\text{iknows } \{NB\}_{\text{pk}(B)}, \text{request } A B \mathfrak{N}_{\mathfrak{B}} NB \in \mathbb{E}}$$

## 5 Turning Fixed-Points into Proofs

We now turn to the proof generator itself (see Fig. 1), putting the pieces together to obtain the security proof with respect to the reference model.

Let  $R_G$  denote in the following the given set of the reference model that describes the protocol, its goals, and the intruder behavior, as described in Sect. 2

and 3. Recall that OFMC chooses an abstraction for the data that honest agents freshly create, and refines these abstractions if the verification fails. As described in Sect. 4, the connection between the data abstraction and the reference model is made by annotating each freshly created message with a label expressing the abstraction. Since this annotation is never referred to in the conditions of any rule, the set of traces remains the same modulo the annotation. We denote by  $R_M$  the variant of the rules with the annotation of the freshly created data.

The next step is the inductive definition of the set of traces  $\mathbb{T}$  in Isabelle, representing the least fixed-point of  $R_M$ . For such inductive definitions, Isabelle proves automatically various properties (e.g., monotonicity) and derives an induction scheme, i.e. if a property holds for the empty trace and is preserved by every rule of  $R_M$ , then it holds for all traces of  $\mathbb{T}$ . This induction scheme is fundamental for the security proof.

We define in Isabelle a set of traces  $\mathbb{T}'$  that represents the over-approximated fixed-point  $FP$  computed by OFMC, expanding all abstractions. We define this via a *concretization function*  $\llbracket \cdot \rrbracket$ :

$$\begin{aligned} \llbracket l \rrbracket &= \{(l, n) \mid n \in \mathbb{N}\} & \llbracket f \ t_1 \ \dots \ t_n \rrbracket &= \{f \ s_1 \ \dots \ s_n \mid s_i \in \llbracket t_i \rrbracket\} \\ \llbracket F \rrbracket &= \cup_{f \in F} \llbracket f \rrbracket & \mathbb{T}' &= \{e_1 \# \dots \# e_n \mid e_i \in \llbracket FP \rrbracket\} \end{aligned}$$

This replaces each occurrence of an abstract datum with an element of the equivalence class it represents, and then builds all traces composed of events from the fixed-point. Note that while  $FP$  is finite,  $\mathbb{T}'$  is infinite.

As the next step, the proof generation module proves several auxiliary theorems using a set of specialized tactics we designed for this purpose. Each proved auxiliary theorem can be used as a proof rule in subsequent proofs. We thus use Isabelle as a framework for constructing a formal tool in a logically sound way (see also [31]).

The first auxiliary theorem is that  $\mathbb{T}'$  does not contain any attacks. The theorem is proved by unfolding the definition of the fixed-point and applying Isabelle's simplifier.

The main part of the proof generation is an auxiliary theorem for each rule  $r \in R_M$  that  $\mathbb{T}'$  is closed under  $r$ . In a nutshell, these theorems are also shown by unfolding the definition of  $\mathbb{T}$  and applying the simplifier. In this case, however, on a more technical level, we need to convert the set comprehensions of the definition into a predicate notation so that the simplifier can recognize the necessary proof steps. This is actually also the point where the labels that annotate the abstraction silently fulfill their purpose: the rules are closed under *any* concretization of the abstract data with the elements from the equivalence class they represent. By our construction, the proof generation does not need to take care of the abstraction at all.

Using the theorems that all rules are closed under  $\mathbb{T}'$ , we can now show the last auxiliary theorem, namely that  $\mathbb{T} \subseteq \mathbb{T}'$ , i.e., that OFMC indeed computed an over-approximation of what can happen according to the reference model. This theorem is proved by induction, using the induction scheme we have auto-

Protocol	FP	time [s]	Protocol	FP	time [s]
ISO 1-pass (sk)	40	21	ISO 2-pass mutual (sk)	130	266
ISO 2-pass (sk)	56	60	NSCK	137	6756
NSL	75	64	TLS (simplified)	166	14202
DenningSacco	76	3984	ISO 2-pass (pk)	168	2965
ISO 1-pass (pk)	82	177	ISO 3-pass mutual (sk)	229	6992
Bilateral Key Exchange	87	182	ISO 2-pass mutual (ccf)	300	1808
Andrew Secure RPC	104	617	ISO 2-pass mutual (pk)	322	5295
DenningSacco	117	1472	ISO 1-pass (ccf)	418	1731
ISO 2-pass (ccf)	124	269	ISO 3-pass mutual (ccf)	664	14434
NSL (w. key server)	127	242			

**Table 1.** Analyzing security protocols using Isabelle/OFMC

matically obtained from the definition of  $\mathbb{T}$  above, i.e. we show that the subset relation is preserved for each rule of  $R_M$ , using the set of the auxiliary theorems.

Finally, we derive our main theorem that  $\mathbb{T}$  contains no attack, which immediately follows from  $\mathbb{T} \subseteq \mathbb{T}'$  and  $\mathbb{T}'$  containing no attack. We have thus automatically derived the proof that the protocol is safe from the given reference model description and OFMC’s output.

## 6 Experimental Results

For first experiments, we have considered several protocols from the Clark-Jacob library [11] and a simplified version of TLS. Tab. 1 shows the results in detail, namely the size of the fixed-point and the time to generate and check the Isabelle proof. Here, we have considered for each protocol all those secrecy and authentication goals that do actually hold (we do not report on the well-known attacks on some of these protocols that can be detected with OFMC). The runtime for generating the fixed-point in OFMC is negligible ( $< 10$  s for each example), while the runtime for Isabelle/OFMC is significantly larger. We hope to improve on the proof generator performance by fine-tuning and specializing the low-level proof tactics where we currently use generic ones of Isabelle.

We suggest, however, that the proof generation time does not affect the experimentation with OFMC such as testing different designs and variants of a newly designed protocol, because the proof generation is meant only as a final step when the protocol design has been fixed and verified with OFMC.

## 7 Related and Future Work

There is a large number of automated tools for protocol verification. [4,6,7,10] in particular are close to the method that is implemented in the new fixed-point module of OFMC: they are all based on an over-approximation of the search

space as described in Sect. 4: the first over-approximation concerns the fresh data, following the abstract interpretation approach of [14] and the second concerns the control structure, i.e., considering a set of reachable events rather than traces. We propose that the other over-approximation-based tools can similarly be connected to Isabelle as we did it for OFMC.

The work most closely related to ours is a recent paper by Goubault-Larrecq who similarly considers generating proofs for an interactive theorem prover from the output of automated tools [18]. He considers a setting where the protocol and goal are given as a set  $S$  of Horn clauses; the tool output is a set  $S_\infty$  of Horn clauses that are in some sense saturated and such that the protocol has an attack iff a contradiction is derivable. He briefly discusses two principal approaches to the task of generating a proof from  $S_\infty$ . First, showing that the notion of saturation implies consistency of the formula and that the formula is indeed saturated. Second, showing the consistency by finding a finite model. He suggests that the first approach is unlikely to give a practically feasible procedure, and rather follows the second approach using tools for finding models of a formula.

In contrast, our work, which is closer to the first kind of approach, shows that this proof generation procedure does indeed work in practice for many protocols, comparable to the results of [18]. We see the main benefit of our approach in the fact that we can indeed use the output of established verification tools dedicated to the domain of security protocols. However, note that the work of [18] and ours have some major differences which makes results hard to compare. First, we consider a reference model where the protocol is modeled as a set of traces; the generated proofs are with respect to this reference model and all abstractions are merely part of the automatic tools. In contrast, [18] considers only one protocol model based on Horn clauses, close to the abstract model in our paper. Taking the soundness of all these abstractions for granted is a weakness of [18]. However, also our approach takes some things for granted, namely a strictly typed model and, based on this, specialized composition rules for the intruder. The typing is common in protocol verification and can be justified by a reasonable protocol implementation discipline [19]. The second assumption is justified by Theorem 1. Our next steps are concerned with lifting these two assumptions from our reference model, i.e. allowing for an untyped reference model with unbounded intruder composition. First experiments suggest that at least the unbounded intruder composition can be feasibly integrated into the proof generation procedure.

Similarly, [32] performs static analysis of security protocols using the tool *Rewrite* which can generate proofs for the theorem prover *Coq*. Like in the case of [18], the resulting proof is with respect to an over-approximated model only and takes the soundness of all abstractions for granted.

A completely different approach to achieve the same goal is currently followed by Meier. Based on [21], he considers an embedding of the *Scyther* tool [15] into Isabelle in order to generate proofs automatically when they fall into the scope of the Scyther method. The advantage is here that one does not rely on a typed

model, while at the stage of this writing the proof generation is not in all cases completely automated.

Several automated verification tools are based on, or related to, automated theorem provers, e.g. SATMC [2] generates Boolean formulae that are fed into a SAT-solver, and ProVerif [4] can generate formulae for the first-order theorem prover SPASS [30]. While there is some similarity with our approach, namely connecting to other tools including the subtle modeling issues, this goes into a different direction. In fact, these approaches additionally rely on both the correctness of the translation to formulae, and the correctness of the automated theorem prover that proves them. In contrast, we generate the proof ourselves and let Isabelle check that proof.

Several papers such as [9,5,23,13] have studied the relationships between protocol models and the soundness of certain abstractions and simplifications in particular. For instance, [13] shows that for a large class of protocols, two agents (an honest and a dishonest one) is sufficient. Recall that we have used this as a standard abstraction of honest agents. While such arguments have thus played an important role in the design of the automated tool and its connection to the reference model, the correctness of our approach does not rely on such arguments and the question whether a given protocol indeed satisfies the assumptions. Rather, it is part of the Isabelle proof we automatically construct that the abstract model indeed covers everything that can happen in the reference model. The automated verifier may thus *try out* whatever abstraction it wants, even if it is not sound. Once again, in the worst case, the proof in Isabelle simply fails, if the abstraction is indeed unsound for the given protocol, e.g., when some separation of duty constraints invalidate the assumptions of the two-agents abstraction.

This was indeed one of the main motivations of our work: our system does not rely on the subtle assumptions and tricks of automated verification. This also allows for some heuristic technique that extends the classical abstraction refinement approaches such as [12]. There, the idea is to start with a simple most abstraction, and when the automatic verification fails, to refine the abstraction based on the counter-example obtained. This accounts for the effect that the abstraction may lead to incompleteness (i.e., failure to verify a correct system), but it is essential that one uses only sound abstractions (i.e., if the abstract model is flawless then so is the concrete model). With our approach, we are now even able to try out potentially unsound abstractions in a heuristic way, i.e. start with abstractions that *usually* work (like the two-agent abstraction). If they are unsound, i.e. the Isabelle proof generation fails, then we repeat the verification with a more refined abstraction.

Isabelle has been successfully used for the interactive verification in various areas, including protocol verification [26,3]. These works are based on a protocol model that is quite close to our reference model (see Sect. 2). There are several works on increasing the degree of automation in interactive theorem provers by integrating external automated tools [17,29,16,27,22]. These works have in

common that they integrate generic tools like SAT or SMT tools. In contrast, we integrate a domain-specific tool, OFMC, into Isabelle.

As further future work, we plan the development of additional Isabelle tactics improving the performance of the verification in Isabelle. Currently, the main bottleneck is a proof step in which a large number of existential quantified variables need to be instantiated with witnesses. While, in general, such satisfying candidates cannot be found efficiently, we plan to provide domain-specific tactics that should be able infer witnesses based on domain-specific knowledge about the protocol model. Finally, we plan to eliminate the limitations on the intruder inductions from the model explained in Sect. 3. While this limitation can be reasonably justified in many cases, the fact that our approach relies on it is a drawback, both in terms of efficiency and also theoretically. In fact, tools like ProVerif instead employ a more advanced approach of rule saturation that allow to work without the limitation. We plan to extend our approach to such representations of the fixed-point.

*Acknowledgments.* The work presented in this paper was partially supported by the FP7-ICT-2007-1 Project no. 216471, “AVANTSSAR: Automated Validation of Trust and Security of Service-oriented Architectures” ([www.avantssar.eu](http://www.avantssar.eu)). We thank Luca Viganò for helpful comments.

## References

1. Armando, A., Basin, D., Boichut, Y., Chevalier, Y., Compagna, L., Cuellar, J., Hankes Drielsma, P., Héam, P.C., Mantovani, J., Mödersheim, S., von Oheimb, D., Rusinowitch, M., Santiago, J., Turuani, M., Viganò, L., Vigneron, L.: The AVISPA Tool for the Automated Validation of Internet Security Protocols and Applications. In: CAV’05. Springer-Verlag (2005). URL <http://www.avispa-project.org>
2. Armando, A., Compagna, L.: SAT-based Model-Checking for Security Protocols Analysis. *Int. J. of Information Security* **6**(1), 3–32 (2007)
3. Bella, G.: *Formal Correctness of Security Protocols*. Springer-Verlag (2007)
4. Blanchet, B.: An efficient cryptographic protocol verifier based on prolog rules. In: CSFW’01, pp. 82–96. IEEE Computer Society Press (2001)
5. Blanchet, B.: Security protocols: from linear to classical logic by abstract interpretation. *Information Processing Letters* **95**(5), 473–479 (2005)
6. Boichut, Y., Héam, P.C., Kouchnarenko, O., Oehl, F.: Improvements on the Genet and Klay technique to automatically verify security protocols. In: AVIS’04, pp. 1–11 (2004)
7. Bozga, L., Lakhnech, Y., Perin, M.: Pattern-based abstraction for verifying secrecy in protocols. *Int. J. on Software Tools for Technology Transfer* **8**(1), 57–76 (2006)
8. Brucker, A., Mödersheim, S.: Integrating Automated and Interactive Protocol Verification (extended version). Tech. Rep. RZ3750, IBM Zurich Research Lab (2009). [domino.research.ibm.com/library/cyberdig.nsf](http://domino.research.ibm.com/library/cyberdig.nsf)
9. Cervesato, I., Durgin, N., Lincoln, P.D., Mitchell, J.C., Scedrov, A.: A Comparison between Strand Spaces and Multiset Rewriting for Security Protocol Analysis. In: ISSS 2002, LNCS 2609, pp. 356–383. Springer-Verlag (2003)
10. Chevalier, Y., Vigneron, L.: Automated Unbounded Verification of Security Protocols. In: CAV’02, LNCS 2404, pp. 324–337 (2002)

11. Clark, J., Jacob, J.: A survey of authentication protocol: Literature: Version 1.0 (1997). [www.cs.york.ac.uk/~jac/papers/drareview.ps.gz](http://www.cs.york.ac.uk/~jac/papers/drareview.ps.gz)
12. Clarke, E., Fehnker, A., Han, Z., Krogh, B., Ouaknine, J., Stursberg, O., Theobald, M.: Abstraction and counterexample-guided refinement in model checking of hybrid systems. *Int. J. of Foundations of Computer Science* **14**(4), 583–604 (2003)
13. Comon-Lundh, H., Cortier, V.: Security properties: two agents are sufficient. In: ESOP'2003, LNCS 2618, pp. 99–113. Springer-Verlag (2003)
14. Cousot, P.: Abstract interpretation. *Symposium on Models of Programming Languages and Computation, ACM Computing Surveys* **28**(2), 324–328 (1996)
15. Cremers, C.: Scyther. Semantics and Verification of Security Protocols. Phd-thesis, University Eindhoven (2006)
16. Erkök, L., Matthews, J.: Using Yices as an automated solver in Isabelle/HOL. In: AFM'08 (2008)
17. Fontaine, P., Marion, J.Y., Merz, S., Nieto, L.P., Tiu, A.F.: Expressiveness + automation + soundness: Towards combining SMT solvers and interactive proof assistants. In: H. Hermanns, J. Palsberg (eds.) TACAS, LNCS 3920, pp. 167–181. Springer-Verlag (2006)
18. Goubault-Larrecq, J.: Towards producing formally checkable security proofs, automatically. In: CSF'08, pp. 224–238. IEEE Computer Society (2008)
19. Heather, J., Lowe, G., Schneider, S.: How to prevent type flaw attacks on security protocols. In: CSFW'00. IEEE Computer Society Press (2000)
20. Lowe, G.: Breaking and fixing the Needham-Schroeder public-key protocol using FDR. In: TACAS'96, LNCS 1055, pp. 147–166. Springer-Verlag (1996)
21. Meier, S.: A formalization of an operational semantics of security protocols. Diploma thesis, ETH Zurich (2007). [people.inf.ethz.ch/meiersi/fossp](http://people.inf.ethz.ch/meiersi/fossp)
22. Meng, J., Quigley, C., Paulson, L.C.: Automation for interactive proof: First prototype. *Information and Computation* **204**(10), 1575–1596 (2006)
23. Mödersheim, S.: On the Relationships between Models in Protocol Verification. *J. of Information and Computation* **206**(2–4), 291–311 (2008)
24. Mödersheim, S., Viganò, L.: The open-source fixed-point model checker for symbolic analysis of security protocols. In: Fosad 2007–2008–2009, LNCS, vol. 5705, pp. 166–194. Springer-Verlag (2009)
25. Nipkow, T., Paulson, L.C., Wenzel, M.: Isabelle/HOL: A Proof Assistant for Higher-Order Logic. LNCS 2283. Springer-Verlag (2002)
26. Paulson, L.C.: The inductive approach to verifying cryptographic protocols. *J. of Computer Security* **6**(1-2), 85–128 (1998)
27. Paulson, L.C., Susanto, K.W.: Source-level proof reconstruction for interactive theorem proving. In: K. Schneider, J. Brandt (eds.) TPHOLs, LNCS 4732, pp. 232–245. Springer-Verlag (2007)
28. Roscoe, A.W., Goldsmith, M.: The perfect spy for model-checking crypto-protocols. In: DIMACS'97 (1997)
29. Weber, T., Amjad, H.: Efficiently checking propositional refutations in HOL theorem provers. *J. of Applied Logic* **7**(1), 26 – 40 (2009)
30. Weidenbach, C., Schmidt, R.A., Hillenbrand, T., Rusev, R., Topic, D.: System description: Spass version 3.0. In: CADE-21, pp. 514–520. Springer-Verlag (2007)
31. Wenzel, M., Wolff, B.: Building formal method tools in the Isabelle/Isar framework. In: K. Schneider, J. Brandt (eds.) TPHOLs 2007, LNCS 4732, pp. 352–367. Springer-Verlag (2007)
32. Zunino, R., Degano, P.: Handling exp, × (and Timestamps) in Protocol Analysis. In: FoSSaCS 2006, LNCS 3921, pp. 413–427. Springer-Verlag (2006)