

A Formal Analysis of Apple’s iMessage PQ3 Protocol

Felix Linker

Department of Computer Science, ETH Zurich

Ralf Sasse

Department of Computer Science, ETH Zurich

David Basin

Department of Computer Science, ETH Zurich

Abstract

We present the formal verification of Apple’s iMessage PQ3, a highly performant, device-to-device messaging protocol offering strong security guarantees even against an adversary with quantum computing capabilities. PQ3 leverages Apple’s identity services together with a custom, post-quantum secure initialization phase and afterwards it employs a double ratchet construction in the style of Signal, extended to provide post-quantum, post-compromise security.

We present a detailed formal model of PQ3, a precise specification of its fine-grained security properties, and machine-checked security proofs using the TAMARIN prover. Particularly novel is the integration of post-quantum secure key encapsulation into the relevant protocol phases and the detailed security claims along with their complete formal analysis. Our analysis covers both key ratchets, including unbounded loops, which was believed by some to be out of scope of symbolic provers like TAMARIN (it is not!).

1 Introduction

Research on secure instant messaging goes back over two decades, with early proposals including Off-the-Record Messaging [9], the Silent Circle Instant Messaging Protocol [31], iMessage, and Signal [32, 24, 27]. Over time, the security community’s understanding of the threat models and security claims for secure messaging evolved. Modern messaging protocols now offer strong guarantees and can communicate messages secretly even in the presence of adversaries who corrupt different parties in different ways during the protocol’s execution. This is befitting given that strong adversaries, like nation states, are capable of compromising both messaging servers and the end points sending and receiving messages. More recently, security against adversaries with quantum computing capabilities has also become an important concern. This requires protection against adversaries who can “harvest now and decrypt later,” namely adversaries who leverage the decreasing cost of mass storage to store the encrypted data

they intercept and to decrypt it in the future when quantum computers become sufficiently powerful [30].

In this paper, we present our formal analysis of Apple’s advanced, widely deployed iMessage PQ3 Messaging Protocol, or PQ3 for short. PQ3 is used across all of Apple’s devices for device-to-device messaging and underlies many other Apple services, e.g., iMessage, FaceTime, HomeKit, and HomePod hand-off. PQ3 is designed to be performant and to offer strong guarantees against powerful adversaries, including those who later possess quantum computers.

PQ3 employs a double-ratchet construction similar to Signal [32]. The protocol takes a hybrid approach to security and combines classical cryptographic primitives, like elliptic curve Diffie-Hellman, and post-quantum primitives, namely ML-KEM [29], a module-lattice-based key-encapsulation mechanism. The hybrid construction means that PQ3’s security does not solely depend on the security of post-quantum primitives, which are less well understood than their classic counterparts. Moreover, PQ3’s integration of hybrid cryptography into the double ratchet provides stronger guarantees than Signal, where a post-quantum Key Encapsulation Mechanism (KEM) is just integrated into the protocol’s setup phase, but not into its ratcheting (see Section 2).

We analyzed PQ3’s security in detail using the TAMARIN prover [33, 28], a state-of-the-art security protocol model checker. Our formal models and proofs are accessible on Zenodo [25]. We report on our model of PQ3, the adversary assumptions, and the protocol’s desired properties. We use TAMARIN’s specification language to specify the messaging protocol and its use of classical and post-quantum cryptography. We also specify all forms of adversary compromise, including the event in which the attacker obtains a sufficiently powerful quantum computer, allowing them to break all non-post-quantum-secure cryptographic primitives. Essentially, the adversary can compromise any key at any time, either through dedicated key-reveal rules or because they obtained a quantum computer. Using TAMARIN’s property language, we formalize and prove both secrecy and authenticity theorems. These theorems precisely express the protocol’s security guar-

antees capturing fine-grained notions of key compromise.

Our analysis establishes that PQ3 provides strong security guarantees against an active network adversary that can compromise any secret key, unless explicitly stated otherwise. For example, PQ3 provides forward secrecy, post-compromise security, and post-quantum security with respect to a “harvest now, decrypt later” adversary. In contrast to Signal, PQ3 provides post-compromise security also against active classical and “harvest now, decrypt later” adversaries and not only against passive, classical adversaries. Moreover, the fine-grained analysis of compromise possibilities and their effects is useful for guiding secure implementations of PQ3. For example, the compromise of a participant’s long-term identity key impacts all security guarantees and thus should be stored with extra care, for example, in a device’s secure enclave.

Contributions Our first contribution is the formalization and machine-checked verification of PQ3 to prove all our security claims. Namely, we use TAMARIN to prove that PQ3 offers strong security guarantees against a powerful adversary with quantum computing capabilities. These guarantees are fine-grained and comprehensive in that omitting any of the many adversary compromise cases leads to attacks. Our verification thereby provides a formal, machine-checked proof that PQ3 meets the high expectations for a modern device-to-device messaging protocol. This high assurance is important given the prominent role of this protocol, which is used in billions of devices worldwide, and its limited prior analysis.

Our second contribution is to show that symbolic security protocol model checkers, in particular TAMARIN, can verify substantial, real-world protocols with nested loops, in their full complexity. This is non-trivial as it entails reasoning about unboundedly many parallel instances of the protocol, where the runs (two devices sending messages) are themselves unbounded. In fact, it was commonly believed that “unbounded (looping) protocols like Signal, and protocols with mutable recursive data structures [...] are also out of scope for symbolic provers, without introducing artificial restrictions” [3]. Our work shows that this is not the case and provides a general methodology for carrying out such proofs.

Organization In Section 2 we survey related work on messaging protocols and their verification. Afterwards, in Sections 3 and 4 we describe PQ3’s threat model, requirements, and the protocol itself. In Section 5 we present our TAMARIN model of PQ3, the adversary, the protocol’s properties, and details on our proofs. We draw conclusions in Section 6.

2 Related work

2.1 Messaging Protocols

Over the past decades, hundreds of secure messaging systems have been proposed [35]. The underlying protocols differ

in how they bootstrap trust to set up initial keys, the properties they achieve, the adversaries they consider, whether bilateral or group communication is supported, and usability. The strongest protocols support message secrecy and authenticity against very strong adversaries. As servers cannot be trusted, encryption must be carried out end-to-end. Moreover, it is common to consider adversaries who can compromise agents’ long-term secrets, and even their session states.

The security bar is now quite high. Modern protocols like Signal, which is used for example in the Signal app, WhatsApp, and Facebook Secret Conversations, offer security guarantees, even when adversaries compromise the devices of the agents running the protocol. In particular, Signal supports both forward secrecy and post-compromise security [14, 8], also called self-healing or backward secrecy. The former protects the protocol’s participants against the future compromise of past sessions, for example, the loss of a long-term secret should not jeopardize the secrecy of previously exchanged messages. The latter helps the participants to recover or “self-heal” from a past compromise to communicate secretly again in the present and future.

Messaging protocols achieve these strong properties by using *ratcheting*, an approach to continually generate new keys. Ratcheting was first proposed in the Off-the-Record Messaging [9] protocol where, with each message round trip, users establish a fresh ephemeral Diffie-Hellman shared secret. Signal further developed this idea with their *double-ratchet algorithm* [32], which nests two ratchets: an outer public-key ratchet and an inner symmetric-key ratchet. This mechanism ensures that the symmetric keys used for encryption and decryption are updated with every message sent, as opposed to just on every round trip. The protocol can recover from past compromises on every round-trip due to a new Diffie-Hellman secret. Forward secrecy is achieved for the symmetric keys as the ratchet chain does not allow one to compute the previous keys from the current message encryption key, but only the future ones.

More recently, researchers have investigated improvements offering guarantees against adversaries with quantum computing capabilities. The Signal protocol uses the Extended Triple Diffie-Hellman (X3DH) Key Agreement Protocol [27] to negotiate the session key used as the ratchet’s initial root key. The recently developed PQXDH Key Agreement Protocol [24] strengthens X3DH by additionally incorporating a post-quantum KEM like Crystals-Kyber [10], and has been verified using both ProVerif and CryptoVerif [4], as well as with a pen-and-paper game-based reduction proof [19]. It has been proven (see Section 2.2) that PQXDH provides forward secrecy even in the presence of an adversary with quantum computing capabilities, provided all KEM private keys remain uncompromised. However, as the post-quantum KEM is only used in the setup phase, the subsequent use of Signal’s double ratchet does not provide post-compromise security against an adversary with quantum computing capabilities, which PQ3

does. Note that both X3DH and PQXDH additionally provide cryptographic deniability [36], which is not provided by PQ3 and hence out of scope for our work.

2.2 Verification of Messaging Protocols

There has been considerable research on verifying messaging protocols using sophisticated constructions like the double ratchet to achieve strong security guarantees. Researchers have studied Signal and variants of it from both a computational and a symbolic perspective, using both pen-and-paper and machine-checked proofs.

Computational proofs A number of pen-and-paper proofs of messaging protocols involving double ratchets have been constructed in the computational setting. This means, in contrast to the symbolic model (introduced shortly), that protocols are analyzed with respect to computational definitions of security. Agents manipulate bit strings, the adversary’s capabilities are modeled by probabilistic polynomial-time Turing machines, and security definitions are thus probabilistic. These models support a more detailed analysis of cryptography than symbolic abstractions. However, the proofs can be quite complex and hence they typically involve their own abstractions or protocol simplifications. Moreover, given that the proofs are traditional pen-and-paper arguments, they are more error-prone than proofs checked by computers. There are exceptions, namely computational proofs constructed with tools like CryptoVerif [6], but these are usually limited to the study of relatively simple combinations of primitives, not complex protocols like the full Signal or PQ3 double ratchet.

In [5], the authors analyze variants of the double ratchet protocol in the Universal Composability framework. As part of their analysis, they consider when keys must be deleted for different properties to hold. Their proofs are game-based with detailed security definitions. Game-based proofs are also given by [11, 13, 1]. In particular, [13] presents a formal analysis of Signal in the random oracle model. Their focus is on Signal’s key agreement and they reason about loops using induction. [1] carries out game-based proofs for a Signal-like protocol; they provide a rational reconstruction of a generalized protocol that modularly achieves the different kinds of properties one wants from Signal and the use of double ratchets. In all these works, security is shown using pen-and-paper proofs, which are not machine checked, and post-quantum security is not considered.

Concomitantly to our work, Stebila carried out a computational analysis of PQ3 [34], providing a reduction argument for its security. He also formalizes the hybrid cryptography integrated into both PQ3’s initialization and double ratchet, and establishes that this provides both forward secrecy and post-compromise security against both classical and “harvest now, decrypt later” adversaries. The modeling of cryptography is, as is standard for computational formalizations, more

concrete and detailed than in our approach. In contrast, the security model, and the proofs (which are game-based, focused on deriving a bound on the adversary’s advantages) are considerably more complex, and proofs are pen-and-paper based, rather than machine checked.

We believe, as Apple researchers also do, that there is substantial benefit to having both kinds of proofs, as they both have their relative strengths. Computational proofs capture the detailed cryptographic assumptions on the operators used. They can also capture the adversary’s advantage in attacking a protocol, by bounding the probability of success for an adversary with given computational resources. In contrast, symbolic proofs better support machine-checked proofs, using different computer-supported proof techniques, like constraint solving and mathematical induction. This supports giving detailed models of protocols’ and adversary’s operational semantics, considering unboundedly many protocol participants and interleaved parallel sessions, and verifying these against detailed, fine-grained security properties.

This value of symbolic proofs is exemplified by our analysis of *injective agreement* [26] (Section 5.3.2), which formalizes that a protocol provides replay protection. [34] did not consider replay in its analysis, and during our TAMARIN proofs, we uncovered that injective agreement can only be provided under additional assumptions (not present in [34]) on the session-handling layer.

Symbolic proofs In terms of verification, the works closest to ours use the symbolic model of cryptography. In this model, messages are represented as terms in a term algebra (rather than bit-strings) and one uses possibilistic rather than probabilistic definitions of security. TAMARIN [33, 28] and ProVerif [7] are examples of tools constructing proofs in this setting. For example, to show that a key is a secret, one would use these tools to prove that, no matter how arbitrarily many protocol runs are interleaved, including runs where the adversary is active, the adversary cannot possibly learn the intended secret. Such proofs may be constructed automatically or interactively, and attempts to prove false statements generally yield attacks on the specified properties.

[18] analyzes Signal’s session-handling layer Sesame. They use TAMARIN to show that, when sessions are accounted for, Signal does not achieve post-compromise security despite the double ratchet having this property. In this work, we do not consider PQ3’s session-handling layer as its specification was not made available to us. Analyzing PQ3 in conjunction with session handling is an interesting line of future work.

[22] use ProVerif and CryptoVerif [6] to analyze a variant of Signal where they extract the models they analyze from an implementation in a JavaScript dialect. Their models are substantially simplified. For example, they lack the inner ratchet based on symmetric cryptography and only consider a fixed, finite number of protocol sessions without loops.

As previously explained, Signal uses the X3DH protocol

to agree on a shared key (the initial root key) prior to the double ratchet’s start. The post-quantum version PQXDH has been analyzed in [4] both symbolically, using ProVerif, and computationally, using CryptoVerif. As the authors explain “Notably, this is the first machine-checked post-quantum security proof of a real-world cryptographic protocol.” While this is indeed the case, they only consider the initialization part of the Signal protocol. They do not reason about the double ratchet construction, which is based on classical cryptography and thus provides no post-quantum security guarantees.

In [3], the authors analyzed Signal based on an F* implementation. They observe: “Notably, Signal has not been mechanically analyzed for an arbitrary number of rounds before. The ProVerif analysis of the Signal protocol in [22] was limited to two messages (three ratcheting rounds), at which point the analysis already took 29 hours. (With CryptoVerif, the analysis of Signal has to be limited to just one ratcheting round).” Their own proof is however also limited and only verifies properties for the outermost ratchet. In contrast, our proof uses induction within TAMARIN to machine check proofs about both ratchets of PQ3. Even in the classical setting, ignoring our post-quantum extensions, verifying the inner ratchet allows us to establish security properties against stronger adversaries who can compromise session state during the inner ratchet’s execution.

3 Requirements and Threat Model

3.1 Security Requirements

Secrecy PQ3 was designed to provide strong secrecy guarantees, namely *message secrecy*, *forward secrecy*, and *post-compromise security*. Message secrecy means that as long as neither participants’ session states are revealed, the adversary cannot learn any of their exchanged messages. Forward secrecy and post-compromise security limit the window in which an adversary can learn exchanged messages after they compromise parts of the session state. We discussed forward secrecy and post-compromise security already in Section 2.1. In short, forward secrecy protects protocol participants against the future compromise of past sessions, and post-compromise security helps to recover or “self-heal” from a past compromise to communicate secretly again in the present and future.

In our security analysis, we define a secrecy lemma that captures all three notions of secrecy and that addresses the precise implications of partial session state compromise. Describing this fine-grained secrecy lemma requires a detailed understanding of the key material used in PQ3, and is thus deferred to Section 5.3.1.

Authentication and Replay Protection A message recipient can identify the message’s sender. We formulate this as an *agreement property*: the recipient and sender agree on their view of the message. For any message received, allegedly

originating from the peer at message counter i , the peer must have actually sent the message using counter i , intending it to go to the receiver. Moreover, this agreement is *injective* [26]. Namely, a given message is only accepted once by the recipient; hence the protocol provides *replay protection*.

3.2 Threat Model

PQ3 seeks to provide the above security properties even when the protocol is run in the presence of a strong active network adversary who may have access to a powerful quantum computer in the future. As an active network adversary, the adversary can read, reorder, intercept, replay, and send any message to any participant. We assume though that devices use strong randomness and that, short of possessing a quantum computer, the adversary cannot factor large numbers or compute discrete logs. Hence, in the pre-quantum era, cryptographic primitives like (elliptic curve) Diffie-Hellman are secure against the adversary.

By default, the adversary can access every participants’ key material unless we explicitly forbid this. We will refine our threat model for each security property and list all the keys that the adversary must not access for the security property to hold. This allows us to focus on which key material the adversary must access to violate a security guarantee and to abstract from whether this compromise is plausible. For example, recently developed cryptographic primitives, designed to provide post-quantum security, may turn out to be flawed. Some keys are stored in a devices’ main memory and relatively easy to compromise, whereas others, like identity keys, are stored in Apple’s Secure Enclave and are thus much harder to compromise. The fact the adversary can access every key by default allows us to consider all of these cases.

In addition, our threat model accounts for the possibility that the adversary may at some point possess a cryptographically-relevant quantum computer. When this happens, the adversary will be able to break all non-post-quantum-secure primitives, such as elliptic curve Diffie-Hellman, and can access all such secret key material, independently of what a refined threat model may state.

We constrain the adversary’s future quantum computing capabilities by assuming that as soon as the adversary possesses a quantum computer, no honest participant runs the protocol. This models an adversary that anticipates future developments in quantum computing and stores all messages sent by the protocol participants. For this reason, the adversary is a passive quantum attacker and is referred to as a “harvest now, decrypt later” adversary.¹

For setup and session establishment, the protocol leverages Apple’s IDentity Services (IDS) key directory. We assume

¹Note that PQ3 only protects past sessions against quantum attackers. To protect active sessions, PQ3’s relies on an elliptic curve signature scheme, which can be broken by a quantum computer.

that this directory is secure in that it only distributes the participants’ authentic public keys. The problem of key authentication is orthogonal to PQ3 and has recently been addressed by Apple with their rollout of “Contact Key Verification” [2].

4 PQ3 Messaging Protocol

PQ3 is a device-to-device messaging protocol where either device can asynchronously exchange messages at any time, independent of the connection status of their peer’s device. We first describe PQ3 at a high-level of abstraction, followed by a more detailed account. We provide a full pseudocode specification of PQ3 in [25].

4.1 High-level Account

In PQ3, communication between two parties, say Alice and Bob, works roughly as follows. Suppose that Alice wants to initiate messaging with Bob.

1. Alice queries Apple’s IDentity Service (IDS) for Bob’s *pre-key material* and a *long-term identity public key*.
2. Alice derives an initial *root key*, *chain key*, and *message key*. Alice encrypts her first message for Bob using the message key and sends Bob the ciphertext along with a signature and the key material necessary to derive the initial root key.
3. Upon receiving this new message, Bob lacks the key to decrypt the ciphertext, and so he must derive it. Bob first queries the IDS to verify Alice’s long-term identity public key and checks the received signature. He uses the key material received from Alice to derive the initial root, chain, and message key and decrypts the initial message. Alice and Bob have now established a shared session.
4. As long as the session does not change direction (i.e., the current sender keeps sending messages), both parties perform *symmetric ratcheting*. In the symmetric ratchet, participants use the old chain key to derive a new chain key and message key.
5. Whenever the session changes direction (i.e., the current receiver wants to reply), both parties perform *public-key ratcheting*. In the public-key ratchet, participants use the old root key and newly sampled asymmetric key material to derive a new root key.

At this high level of abstraction, Steps 2–5 resemble the standard double-ratchet construction. But there are significant differences in the concrete details on how the ratchets are performed, in particular how a post-quantum KEM is integrated into the ratcheting.

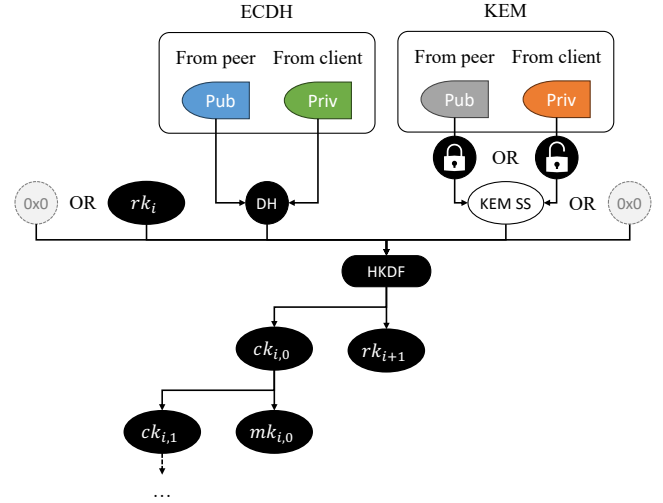


Figure 1: Dependency between the keys used by PQ3. Arrows denote that one value is used to derive another. The lock icons denote KEM encapsulation or decapsulation respectively. Sometimes a zero-byte sequence is used instead of a root key or KEM shared secret.

4.2 More Detailed Account

We now expand on the above account. Although this account is more detailed, we still focus on the essential ideas and we omit some low-level details, like message and key derivation tags. Moreover, we describe some additional features of PQ3 at the end of this section.

Keys PQ3 specifies many keys. Every participant has a *long-term identity key*, a P-256 ECDSA public/private key pair to authenticate messages and other key material. Long-term identity public keys are distributed and authenticated using the IDS. All other keys are used to derive message keys. Figure 1 depicts the dependencies between these keys.

We start by introducing PQ3’s three types of symmetric keys. These symmetric keys are always derived with respect to a given public-key ratchet step (identified by i in Fig. 1). *Message keys* (depicted as $mk_{i,0}$) are the message encryption keys and are derived from *chain keys* (depicted as $ck_{i,0/1}$). Chain keys are derived from either previous chain keys or initially from the same entropy sources as the root keys. *Root keys* (depicted as $rk_{i/i+1}$) are used in every public-key ratchet step and, in particular, maintain the entropy from previous public-key ratchets.

Root and initial chain keys are derived from three entropy sources: the session’s previous root key (or a zero-byte sequence upon session start; rk_i in Fig. 1), an ECDH shared secret (“DH” in Fig. 1), and optionally a KEM shared secret (replaced with a zero-byte sequence when omitted; “KEM SS” in Fig. 1). To establish these shared secrets, every client uses P-256 ECDH public/private key pairs, which we call *ECDH*

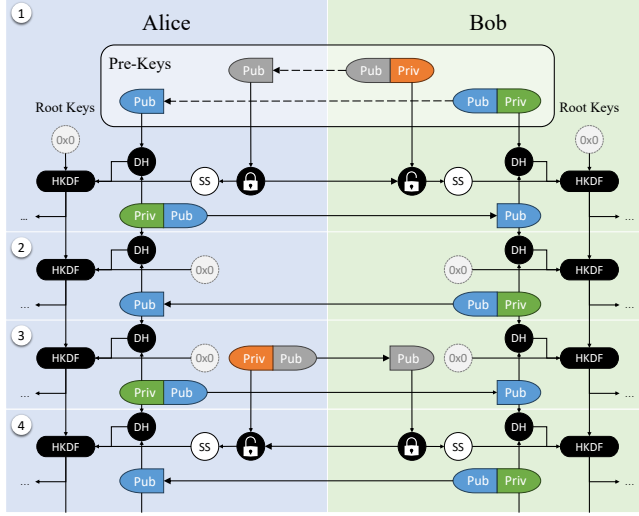


Figure 2: PQ3’s public-key ratchet. Each block 1-4 illustrates a public-key ratchet step. We omit the symmetric ratchet; chain and message keys are derived from the output of the HKDF (denoted by “...”). In Step 1, Alice initiates a session with Bob and uses pre-key material (white box) to derive a root key. Alice sends a freshly encapsulated shared KEM secret (lock icon), and a freshly sampled ECDH public key to Bob that Bob can use to derive session keys. New KEM shared secrets are only encapsulated and shared when a new KEM public key was sent in the previous public-key ratchet (see block 4). Orange/gray key pairs denote ML-KEM keys, green/blue key pairs denote ECDH keys. This figure was inspired by [32].

keys, and ML-KEM 768 or 1024 public/private key pairs, which we call *KEM keys*. Clients establish the ECDH shared secret by combining an ECDH public key from their peer with their own ECDH private key (“ECDH Pub/Priv” in Fig. 1). Clients establish the KEM shared secret either by encapsulating it for their peer using their peer’s KEM public key or by having their peer encapsulate it for them and decapsulating it with their own KEM private key (“KEM Pub/Priv” in Fig. 1).

In general, every client uses distinct, fresh ECDH and KEM keys for every session, the public part of which they send in PQ3 messages to their peer. These session-specific keys are called *ephemeral keys*. Ephemeral keys are short-lived and used only for a specific session. To support asynchronous messaging, clients use ECDH and KEM public *pre-keys* instead of their ephemeral counterparts upon session start (the ECDH and KEM keys depicted in Figure 1 could be either ephemeral or pre-keys). Clients upload their pre-keys to the IDS using timestamped pre-key bundles, which are signed with their long-term identity key. Clients can fetch their peers’ pre-keys from the IDS to start a new session with any of their peers’ clients without requiring that client to be online. Pre-keys can be reused in multiple sessions, but are only used upon session start. PQ3 uses ML-KEM 768 key pairs for ephemeral KEM keys and ML-KEM 1024 key pairs for KEM pre-keys.

Session Establishment In the following, we assume, as before, that Alice wishes to establish a new session with Bob. We depict an example run of PQ3 in Figure 2, specifically showing the key derivations of both parties. The figure shows four public-key ratchet steps (numbered 1-4). Step 1 illustrates session establishment as explained next. Note that all messages sent between parties include a signature by the respective sender for authentication purposes using their long-term identity key. We omit signatures, long-term identity keys, the steps of the symmetric ratchet, and sent messages from

the figure to avoid clutter and to focus on the key material used in root key derivation.

Alice’s actions are depicted in the left, blue half of Figure 2. Alice initiates her session with Bob by performing an IDS query for Bob’s identity. Alice thereby learns three keys from the query’s result: Bob’s long-term identity public key, an ECDH public pre-key, and a KEM public pre-key. Querying and using pre-keys is depicted within the white box in Figure 2. Alice then generates a fresh ECDH ephemeral public/private key pair (“Priv/Pub” in Step 1) and encapsulates a fresh KEM shared secret with Bob’s public pre-key (lock icon in Step 1). The encapsulation algorithm provides Alice with the cleartext KEM shared secret for her use (shown as “SS” in Step 1), and ciphertext to be given to Bob (the lock to the right of “SS”, showing that it used the KEM public pre-key from above). Bob can decapsulate the KEM shared secret with his KEM private pre-key to receive the same KEM shared secret. Alice then combines her ECDH ephemeral private key with Bob’s ECDH public pre-key to obtain the initial ECDH shared secret (depicted as “DH”).

Alice proceeds to derive the initial root key and the associated initial chain key from the ECDH shared secret, the KEM shared secret, and a zero-byte sequence, which stands in for the previous root key. This is depicted on the far left of Figure 2 as “HKDF” in Step 1. She derives a message key from the initial chain key and encrypts her initial message with that message key. She sends Bob the ciphertext, her ECDH ephemeral public key, the KEM encapsulation (with the latter two shown in Figure 2), a hash of Bob’s public pre-keys (the *pre-key hash*), and a signature on all these elements and some additional authenticated data. The exact values of the authenticated data field are unspecified, and the field can be used freely by applications.

Bob uses that message to derive the initial root and chain key. Bob’s actions are depicted in the right, green half of

Figure 2. Bob first performs an IDS query to receive Alice’s long-term identity public key (not depicted in Figure 2), which he uses to verify the message signature. Bob then looks up the private parts of his pre-keys used by Alice, which are identified by the pre-key hash. Bob decapsulates the KEM encapsulation to obtain the KEM shared secret (the open lock symbol in Step 1), and combines Alice’s ECDH public ephemeral key with his ECDH private pre-key to establish the ECDH shared secret (“DH” in Step 1). With these two values (and the zero-byte sequence), Bob computes the initial root and chain key (illustrated by “HKDF” in Step 1) and derives a message key from that chain key to decrypt the ciphertext.

Symmetric Ratchet With a shared root key established, Alice can send any number of additional messages to Bob without the participants updating the root key. Nevertheless, each of these messages will be encrypted with a distinct key derived by symmetric ratcheting. Whenever a participant encrypts a message, they use the current chain key to derive a message key, and then ratchet the chain key forward by deriving a new chain key from the previous one. PQ3 establishes per-message forward secrecy as soon as the previous chain and message keys are deleted, i.e., participants should only store the latest root and chain key. The symmetric ratchet, though, is only executed as long as the conversation’s direction does not change, i.e., as long as the current sender keeps sending. Whenever the current receiver wishes to respond, they perform a public-key ratchet instead.

Public-Key Ratchet Suppose, after receiving some messages from Alice, that Bob wants to reply. This means that the conversation *changes direction*, and whenever this happens clients perform the public-key ratchet. Every public-key ratchet updates the root key and derives a new, initial chain key. The steps taken to derive these new keys are similar to the steps taken during session establishment. Figure 2 illustrates (next to session establishment) three further public-key ratchet steps (numbered 2-4).

To perform the public-key ratchet, Bob first generates a fresh ECDH ephemeral public/private key pair. Depending on the conversation’s state, Bob may additionally perform either of the following two actions: (i) use the encapsulation algorithm to produce a new KEM shared secret and ciphertext (for decapsulation by Alice), or (ii) generate a new KEM ephemeral public/private key pair. Action (i) is performed whenever Bob’s peer, Alice, performed Action (ii) in the previous public-key ratchet. To save bandwidth, Action (ii) need not always be performed. Instead, a custom heuristic determines when a client refreshes its KEM keys. The heuristic accounts for the threat environment, performance, and other requirements. As per iOS 17.4, PQ3 clients send a fresh KEM public key roughly every 50 messages or whenever they have not sent a fresh KEM public key within a week [21].

Bob then derives the next root key and the associated initial chain key. He first combines his freshly generated ECDH ephemeral private key with Alice’s ECDH ephemeral public key to obtain the new ECDH shared secret. He then uses the previous root key, the new ECDH shared secret, and either the new KEM shared secret or a zero-byte sequence (depending on whether Bob performed Action (i)) to derive the next root key and associated initial chain key. He again derives a message key from that chain key to encrypt his message and sends Alice the following values: the ciphertext, his fresh ECDH ephemeral public key, optionally the new KEM encapsulation (Action (i)), optionally his new KEM public key (Action (ii)), and a signature on all the above.

Figure 2 depicts in Step 3 that Alice generates a new ephemeral KEM public/private key pair and sends the corresponding public key to Bob, i.e., Alice executes Action (ii) above. This means that Bob will execute Action (i) in Step 4.

Overall, the cryptographic constructions used are hybrid: all key derivations incorporating a KEM shared secret also involve classical secrets. This design entails (and we establish this formally in our proofs) that PQ3’s security is at least as strong as when using classical cryptography alone. The repeated use of the KEM encapsulation in the protocol therefore strictly strengthens the protocol to provide post-compromise security even against a “harvest now, decrypt later” adversary who managed to access some KEM shared secret.

5 Security Proofs

In this section, we describe how we modeled PQ3 and proved its security using TAMARIN. We briefly introduce TAMARIN (Section 5.1), describe our protocol model (Section 5.2), the formal security properties (Section 5.3), and our proofs (Section 5.4). Our protocol model covers PQ3 in its full complexity, including its nested loops, all its cryptographic primitives, and their combinations. We discuss limitations and proof effort in Section 5.5. All our formal models and proofs are openly accessible on Zenodo [25].

5.1 Background on Tamarin

TAMARIN works in the *symbolic model* of cryptography, which supports a high degree of automation when constructing proofs. TAMARIN uses labeled *multiset rewriting rules* to model setup assumptions and the behavior of protocol participants. The participants play in so-called roles, where the possible actions of each role are given by sets of rules. TAMARIN verifies security properties with respect to an active network adversary who can read, intercept, reorder, replay, and send messages. In addition to this built-in adversary, modelers can give the adversary additional capabilities using explicit rules.

Each rule has a premise and conclusion. These consist of (potentially *persistent*) *facts*, which store the terms that TAMARIN manipulates and reasons about. The rules together

specify an infinite-state transition system. Each state of this transition system includes the protocol-state associated with each role instance, the adversary’s knowledge, all messages being sent on the network, and more. To apply a rule, the facts in its premise must be found in the current global state. When a rule is applied, all non-persistent facts appearing in the premise of the rule are removed from the state and instances of all facts in the conclusion of the rule are added.

All rules are labeled and TAMARIN reasons about traces, which are sequences of the instantiated rules’ labels. For this, TAMARIN supports a subset of first-order logic to specify the properties one then proves. Furthermore, formulas in this logic can also be used to specify *restrictions* on which traces TAMARIN should consider when proving theorems. Restrictions can be used, for example, to state that a participant performs a certain check, e.g., signature verification, in which case traces with failed checks would be excluded.

To model different cryptographic primitives, TAMARIN supports a number of built-in equational theories, for example, for symmetric encryption and message signing. The user can additionally define their own equational theories.

TAMARIN reasons using backwards search. Starting from the protocol’s specification, it negates the property to be verified and searches for a trace representing an attack. If there cannot exist any such trace, then the property is proven. Internally, TAMARIN uses constraint solving, and supports both an automatic mode and an interactive mode. Each step is machine-checked, using sound and complete proof rules. However, as the underlying problem is undecidable, there is no guarantee of termination. Users can help TAMARIN construct proofs in an interactive mode, where again the prover checks each proof step. Users can also help TAMARIN by specifying auxiliary properties that can be proven once and for all and that can be reused in larger proofs.

Finally, TAMARIN also supports a form of induction. This is essentially an induction on the length of a trace with a distinguished special *last* timepoint. Timepoints in general provide an order on the steps in the protocol. For the special last timepoint, the property must be proven, with it being assumed at all previous timepoints. We explain TAMARIN’s induction scheme more detailed in Appendix A.1.1.

5.2 Protocol Model

We used TAMARIN to comprehensively model PQ3 as described in Section 4.2. Our model comprises a set of rules and restrictions, modelling PQ3 as a state transition system, together with an equational theory, modelling cryptographic primitives. In this section, we describe the rules and restrictions and refer to Appendix B for details on our equational theory. The full protocol model is provided on Zenodo [25].

We provide an overview of our model’s protocol rules in Figure 3. Our formal model has three parts. The first part models the generation of long-term signing keys and pre-

keys (rule `UserKeyGen`), and IDS queries (rule `QueryIDS`). These are setup rules, which are the same for all participants, independent of whether they start a session as the sender or receiver. The second and third part model the adversary’s capabilities and PQ3’s protocol flow respectively.

In our model of the adversary’s capabilities, we allow the adversary to compromise every private, root, chain, and message key through dedicated reveal-rules, unless our security lemmas explicitly forbid a certain key to be revealed. Additionally, we model the “harvest now, decrypt later” capability as follows. Whenever participants generate a non-post-quantum-secure key, like a fresh ephemeral ECDH private key, our model saves the key in a persistent state fact (i.e., a fact that is not consumed when it is used in a rule’s premise). The adversary can then access any secrets stored this way after the rule `PQAttackerStart` is applied, but from that point on, no honest participant runs PQ3.

Our model of PQ3’s protocol flow is depicted as the big blue box in Figure 3. The left-hand side depicts all sender-related rules, the right-hand side all receiver-related rules, and in the center is a `Session` fact that stores all information needed to send and receive messages. For example, a `Session` fact stores a participant’s most recently generated ECDH and KEM private keys and the corresponding public keys of their peer, as well as any derived root and chain keys.

A new session is started by applying one of the rules `SessionStartAsSender` or `ReceiverStart`. These are the only two rules that only produce and do not consume a `Session` fact. Most other rules update a session, i.e., they consume and produce a `Session` fact, and they can be applied arbitrarily many times per session. After a new session has started, one of two things can happen. Either the conversation does not change direction and then both participants will apply the symmetric ratchet rules, or the conversation changes direction and the public-key ratchet rules are applied.

When being the receiver, a participant may non-deterministically choose to become sender. When they do, they perform the public-key ratchet. Depending on whether their peer had sent them a new KEM public key previously, they may additionally encapsulate a new KEM shared secret. Also, the new sender may non-deterministically send a new KEM public key themselves to their peer.

A participant changes from the sender to the receiver role when they receive a new message while being in the sender state. When a participant becomes the receiver, they perform the public-key ratchet as well. In one of the two rules, they do so using a decapsulated KEM shared secret, and in the other rule they use a zero-byte sequence instead.

Intuitively, one can consider our model as implementing two nested loops. First, there is the outer, public-key ratchet loop where participants generate new ephemeral ECDH secret keys and derive root and chain keys. Second, there is the inner, symmetric ratchet loop where participants derive message keys and send messages. The symmetric ratchet loop always

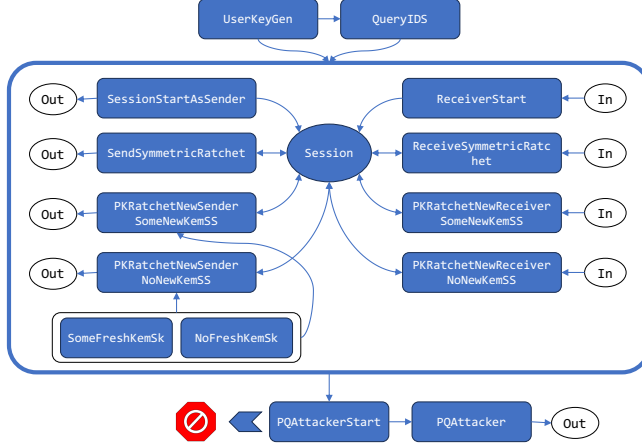


Figure 3: Overview of our formal model. Rectangles denote rules and ellipses denote facts, with their respective name printed inside. Arrows denote fact consumption and generation or rule transition. The white rectangle around `Some/NoFreshKemSk` denotes that either of the rules is applied non-deterministically. The rule `PQAttackerStart` can be applied at any point. When this happens, protocol execution halts (modeling a “harvest now, decrypt later” adversary) and thereafter the rule `PQAttacker` can be applied, which reveals any non-post-quantum-secure secret to the adversary. This figure omits rules that reveal key material.

runs within one iteration of the public-key ratchet loop.

5.3 Properties Specified

5.3.1 Secrecy

PQ3 aims to satisfy three secrecy properties: message secrecy, forward secrecy, and post-compromise security. In our formalization, we combine all three into a single property. This property is formulated as a formula, called a *lemma* in TAMARIN, as one must prove that it holds for the protocol.

Figure 4 contains our secrecy lemma.² It states that the adversary cannot know a message (line 4) that has been previously sent (line 2), unless the adversary succeeds in at least one of four kinds of compromise, listed below. The kinds of compromise are formulated with respect to the keys referenced by the `SessionSecrets` fact. This fact lists all keys and shared secrets used by the sender when sending the respective message, e.g., their most recently sampled ephemeral ECDH public key (`myEcdhPk`) and the most recently encapsulated KEM shared secret (`kemSS`). We sketch a possible attack for each kind of compromise to show that dropping any but the first disjunct yields a counterexample. To learn a message sent with PQ3, the adversary must compromise at least one of:

- The message key used during encryption from either the receiver or sender (line 5 in Figure 4). Should the adversary learn the message key, they could simply decrypt the message themselves.
- One of the chain keys used in the symmetric ratchet to derive the message key from either the receiver or sender (lines 6-7; $a \ll b$ denotes that a is a *subterm* of b [17]). Should the adversary learn one of these chain keys, they could simply derive the message key themselves.

²In the following, we will sometimes shorten the names of facts in lemmas compared to the source files, e.g., `RevealIdentityKey` may become `RevealIDKey`.

- The recipient’s long-term identity key before the message `msg` was sent (line 8). In this case, the adversary could generate a fresh ECDH ephemeral and KEM encapsulation key and send them to the messaging partner in question. This attack allows the adversary to carry out all communication in their victim’s stead.
- One of the ephemeral ECDH secret keys, used to derive the most recently established ECDH shared secret, *and* the KEM shared secret (lines 9-15). This allows the adversary to perform a public-key ratchet step themselves.

The adversary can learn an ECDH secret key either through direct compromise (lines 10-11) or using a quantum computing attack should a sufficiently powerful quantum computer be available (line 9). The compromise of the sender’s ECDH pre-key has no effect because a sender will always sample a fresh ECDH ephemeral key upon session start.

The KEM shared secret can be effectively compromised in two ways. First, the adversary can compromise the KEM secret key used for encapsulation (lines 12-13). Second, the adversary can circumvent the need to compromise the KEM shared secret by compromising a root key derived after that KEM shared secret was established (lines 14-15). In the latter case, if the adversary additionally learns an ECDH secret key used in a subsequent public-key ratchet step, they can derive the respective initial chain key themselves.

In addition to the ECDH and KEM shared secret, the adversary also requires the root key from the previous public-key ratchet to perform the current public-key ratchet themselves. Our threat model, however, permits this root key to be revealed to the adversary anyway.

Recall that our threat model assumes that the adversary can access all key material unless explicitly forbidden. Our secrecy lemma only forbids the adversary to access key material related to sending the message in question. All key-reveal assumptions in lines 5-15 use the key material introduced

```

1 All id me them msg ad myEcdhPk theirEcdhPk kemSS encapPk rk chainKey msgKey #t.
2   ( Sent(id,_,me,them,msg,ad) @ t
3   & SessionSecrets(myEcdhPk,theirEcdhPk,kemSS,encapPk,rk,chainKey,msgKey) @ t)
4 ==> (not Ex #x. K(msg) @ x)
5   | (Ex #x. RevealMessageKey(me,msgKey) @ x) | (Ex #x. RevealMessageKey(them,msgKey) @ x)
6   | (Ex ckC #x. RevealChainKey(me,ckC) @ x & (ckC << chainKey | ckC = chainKey))
7   | (Ex ckC #x. RevealChainKey(them,ckC) @ x & (ckC << chainKey | ckC = chainKey))
8   | (Ex #x. RevealIDKey(them) @ x & x < t)
9   | ( ( Ex #x. PQAttack() @ x)
10      | (Ex #x. RevealECDHPreKey(them,theirEcdhPk) @ x)
11      | (Ex #x. RevealECDHKey(id,me,myEcdhPk) @ x) | (Ex #x. RevealECDHKey(_,them,theirEcdhPk) @ x))
12      & ( (Ex #x. RevealKemKey(me,encapPk) @ x) | (Ex #x. RevealKemKey(them,encapPk) @ x)
13          | (Ex #x. RevealKemPreKey(me,encapPk) @ x) | (Ex #x. RevealKemPreKey(them,encapPk) @ x)
14          | (Ex k #x. RevealRootKey(me,kemSS,k) @ x & k << rk)
15          | (Ex k #x. RevealRootKey(them,kemSS,k) @ x & k << rk)))

```

Figure 4: Secrecy lemma. The lemma formalizes that if a message `msg` was sent using the secret values referenced by `SessionSecrets`, then either the message cannot be known by the adversary (line 4) or the adversary compromised a specific combination of keys (lines 5ff.). Section 5.3.1 explains this lemma, line-by-line, in further details.

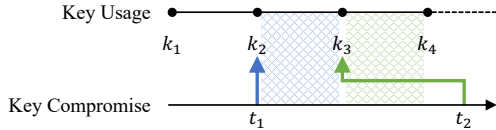


Figure 5: A participant derives four initial chain keys (k_1 - k_4) over time and the adversary compromises k_2 and k_3 at times t_1 and t_2 respectively. Independent of when it occurs (compare t_1 with t_2), key compromise has a similar and limited effect: The adversary can only learn messages sent before the next initial chain key is derived (the shaded areas).

in line 3, which in turn is bound to the `Sent` event in line 2 by the variable `t`. Thus, proving secrecy establishes forward secrecy and post-compromise security as we explain next.

For long-term identity keys, we show that PQ3 provides forward secrecy in that all messages exchanged prior to the compromise of such a key remain secure (see line 8). For most encryption keys (exceptions and details below), we establish forward secrecy and post-compromise security in that to compromise a given message, the adversary must learn the respective key used for that message and the compromise of past or future keys has no effect. Note that “past” and “future” here refer to the points in time when a key was used, not when it was compromised. In particular, this allows us to establish post-compromise security guarantees even *after* the adversary obtained a quantum computer and participants stopped running the protocol. Although participants will no longer rotate keys, as they no longer run the protocol, they will have self-healed from the compromise of any other key than the most recently used one. For an illustration, see Figure 5.

To provide an example for why our secrecy lemma entails forward secrecy and post-compromise security, consider the lemma modelling ECDH key forward secrecy in Figure 6. The lemma resembles our secrecy lemma in Figure 4, but addition-

```

1 All id ecdhKey1 ecdhKey2 m #t1 #t2 #t3.
2   ( SessionSecrets(ecdhKey1, ...) @ t1
3   & SessionInfo(id, ...) @ t1
4   & K(ecdhKey1) @ t2
5   & SessionSecrets(ecdhKey2, ...) @ t3
6   & Sent(id, ..., m, ...) @ t3
7   & t3 < t1)
8 ==> (not Ex #x. K(m)@#x)
9   | (...) // as in secrecy lemma

```

Figure 6: Sketch of a potential formalization of ECDH key forward secrecy. Observe that this property is strictly weaker than our secrecy lemma in Figure 4 because we only add conjuncts to the implication’s left-hand side. Thus, this formalization of forward secrecy is implied by our secrecy lemma.

ally assumes that the adversary learned a relevant ECDH key derived before the current message was sent. This modified lemma accurately models ECDH key forward secrecy, but it is strictly weaker than our secrecy lemma. Formally, this is the case because we only strengthen the implication’s left-hand side. We also cannot drop any disjunct on the implication’s right side because, if we could, our secrecy lemma would not be provable (we sketched attacks on the previous page). Intuitively speaking, the adversary does not gain more power when we explicitly add the event that they learn a respective ECDH key to the trace because we assume that the adversary can access all keys by default anyway.

In general, we establish per-key forward secrecy and post-compromise security upon key rotation. For some keys, forward secrecy and post-compromise security are only established under further constraints. In these cases, our secrecy lemma precisely defines the point in time at which forward secrecy or post-compromise security are established. We list all forward secrecy and post-compromise security guarantees entailed by our secrecy lemma below and, wherever necessary,

describe the constraints on these guarantees. When the adversary does not possess a quantum computer, PQ3 provides:

- Long-term identity key forward secrecy.
- ECDH ephemeral key forward secrecy and post-compromise security.
- ECDH pre-key post-compromise security as soon as a new ECDH ephemeral key is generated by a session's initial recipient.

In practice, PQ3 also provides forward secrecy for ECDH pre-keys as it requires that participants update their pre-keys registered at the IDS every 2 weeks. As soon as a client registers a new pre-key, they establish forward secrecy for all previous session-start messages sent to them.

Should the adversary at some point break all non-ML-KEM keys using a quantum computer, PQ3 still provides:

- ML-KEM key post-quantum forward secrecy and post-compromise security.
- Chain and message key forward secrecy and post-compromise security. These properties are established unconditionally except for chain key post-compromise security, which is established upon the next public-key ratchet. PQ3 establishes these properties even when the adversary possesses a quantum-computer because these keys depend on KEM-encapsulated secrets.

Note that working out and rigorously proving such fine-grained notions of secrecy is nontrivial and one strongly benefits here from a proof assistant. Overall, our TAMARIN proof of secrecy establishes that, in the absence of the sender or recipient being compromised, all keys and messages transmitted are secret. The secrecy property is fine-grained in that compromises can be tolerated in a well-defined sense where the effect of the compromise on the secrecy of data is limited in time and effect as described above. Moreover, we show that PQ3 combines the security of both classical and post-quantum-secure cryptographic primitives. Hence, to break PQ3 one must break both.

5.3.2 Agreement

In contrast to secrecy, formalizing agreement is much simpler. This is because PQ3 relies on the participants' long-term identity keys' security to provide agreement. Compromise of a participant's long-term identity key is both necessary and sufficient to break agreement. It is necessary because an attacker must generate a message signature when trying to spoof a sender, and it is sufficient because a sender need not compromise the sender's encryption keys to send an inauthentic message; they can simply generate their own and send them alongside the faked message.

```

1 All id i s r m ad #t.
2   Received(id, i, s, r, m, ad) @ t
3 ==> ( (Ex #x. Sent(_, i, s, r, m, ad) @ x & x < t)
4       | (Ex #x. RevealIDKey(s) @ x & x < t))

```

Figure 7: Agreement lemma. For every message-receive event, there must be a corresponding message-send event for which the participants agree on the authenticated data, sender, receiver, and message counter, unless the sender's long-term identity key was previously compromised.

```

1 All s1 s2 r1 r2 m ad ecdhPk1 mk1 ecdhPk2 mk2 #t1
2   #t2.
3   ( Received(_, _, s1, r1, m, ad) @ t1
4       & SessionSecrets(ecdhPk1, _, _, _, mk1) @ t1
5       & Received(_, _, s2, r2, m, ad) @ t2
6       & SessionSecrets(ecdhPk2, _, _, _, mk2) @ t2)
7 ==> ( (t1 = t2)
8       | ( ecdhPk1 = ecdhPk2 & mk1 = mk2
9           & s1 = s2 & r1 = r2
10          & Ex #x. ECDHPreKeyGen(r1, ecdhPk1) @ x)
11       | (Ex #x. RevealIDKey(s1) @ x & x < t1 )
12       | (Ex #x. RevealIDKey(s2) @ x & x < t2))

```

Figure 8: Injective agreement lemma. It formalizes that for two message-receive events with the same message m and authenticated data ad , these events must be the same (line 7), or they were sent using the recipients pre-key (lines 8f.), or one sender's identity key was compromised (lines 11ff.).

Our formalization of agreement is split into two TAMARIN lemmas (Figures 7 and 8). The first lemma formalizes agreement: Whenever a participant r receives a message m and authenticated data ad , apparently from s and with message counter i , then either s had previously sent m to r with counter i or that senders' long-term identity has been compromised in the past.

The second lemma formalizes that the agreement is injective [26], meaning that there is a one-to-one mapping from receive-events to send-events. This lemma states that for every two honest message-receive events with the same message and authenticated data, these events must either be identical ($\#t1 = \#t2$), or a recipient's ECDH pre-key rather than an ephemeral key was used to derive the message key (lines 8-10), or either of the senders were compromised. Compromise of one sender suffices to violate injective agreement because agreement does not entail secrecy. The adversary could learn a message by compromising the ECDH and KEM keys of the session. They could then send the message again, which requires the compromise of a long-term identity key, however, to produce the necessary signature.

During our proof efforts, we noticed a trivial violation of injective agreement, which is covered by lines 8-10. PQ3 cannot provide injective agreement for session-start messages (and messages sent as part of the symmetric ratchet directly

thereafter) as pre-keys can be reused for session starts. Thus, recipients will accept session-start messages multiple times. In practice, this case must be addressed by an application’s session-handling layer, which defines under which conditions clients will accept session-start messages from devices they already have an existing session with. We shared this finding with Apple researchers who confirmed that the iMessage session-handling layer indeed addresses this case. Put differently, our formal proofs highlight precisely the assumptions on session-handling needed to securely deploy PQ3.

5.4 Proofs & Proof Methodology

We describe here our proofs and proof methodology for PQ3. Our proof methodology applies to theories that include (possibly nested) loops and for which trace formulas like secrecy or authentication are to be proven. We present our methodology more generally and with further details in Appendix A.

We encountered two challenges when verifying PQ3. First, PQ3 employs a nested loop. If not carefully handled, loops result in prover non-termination as they are unrolled infinitely often. TAMARIN provides induction to address this problem, but using induction correctly, especially when loops are nested, requires postulating nontrivial auxiliary lemmas.

Second, our threat model considers the leakage of “synthetic” key material, derived using a KDF, and our lemmas naturally must refer to this key material. When proving secrecy, we repeatedly encountered cases similar to the following. TAMARIN would consider an honest session sending a message, claiming that the adversary could get the decryption key for this message (violating secrecy) from a completely unrelated session. We call such unrelated sessions *ghost sessions*. In this case, the non-trivial proof goal was to convince TAMARIN that the ghost session must be the same as the honest session or the peer’s session. Note that other protocol models typically only consider the leakage of “atomic” key material, i.e., key material modelled as a fresh term.

To address these two challenges, our methodology uses three kinds of auxiliary lemmas.

Loop-Jump Lemmas These lemmas allow one to skip unrolling the steps of a (nested) loop and jump to a “relevant” point in a loop, for example, its beginning or where a specific term was introduced.

Variable-Linking Lemmas These lemmas establish that for two instances of the same fact using two variables a and b , if both facts have the same value for a , they must have the same value for b .

Adversary-Construction Lemmas These lemmas formalize how an adversary could construct a term. Typically, the adversary can either construct it or access it using a dedicated reveal rule (which in turn typically implies a contradiction to the threat model). Figure 9 depicts

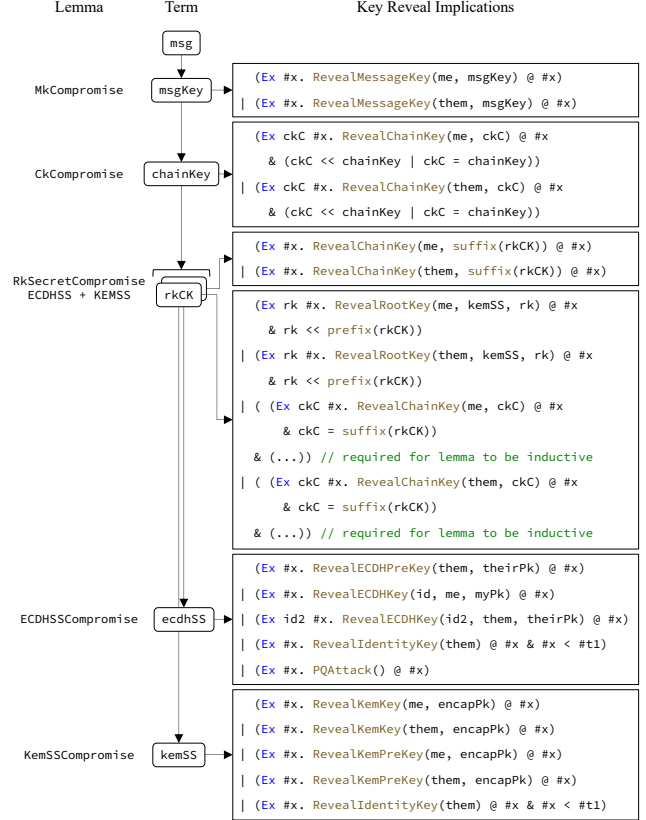


Figure 9: Connection between Adversary-Construction Lemmas for Message Secrecy. Arrows denote logical implication. We omit two conjuncts in `RkSecretCompromiseKEMSS` that are only required to prove the lemma by induction. We provide more details on these conjuncts in our formal model [25].

our model’s adversary-construction lemmas. For example, `CkCompromise` states that the adversary can only know a chain key if they know the value that gets split into the root and chain key (`rkCK`), or they compromised this or a previous chain key.

Loop-jump lemmas are the foundation for proving properties of models including nested loops. Without such lemmas, TAMARIN’s induction fails to prove even the simplest properties of an outer loop. The induction hypothesis will not apply in cases where a step in the outer loop is directly preceded by a step in an inner loop. Moreover, adversary-construction lemmas are required to deal with the complicated terms that are computed in nested loops, and variable-linking lemmas are required to address ghost sessions.

We proved secrecy for PQ3 using a series of adversary-construction lemmas, depicted in Figure 9, which in turn were proven using the loop-jump and variable-linking lemmas in Figures 10 and 11. Concretely, when proving secrecy, TAMARIN first negates the original lemma and tries to construct a trace satisfying the negated lemma, i.e., TAMARIN

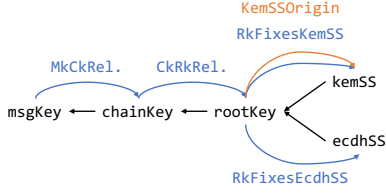


Figure 10: Loop-Jump (orange) and Variable-Linking Lemmas (blue) Related to Key Derivation. Black arrows indicate which variables are used to construct other variables, e.g., a message key is derived from a chain key.

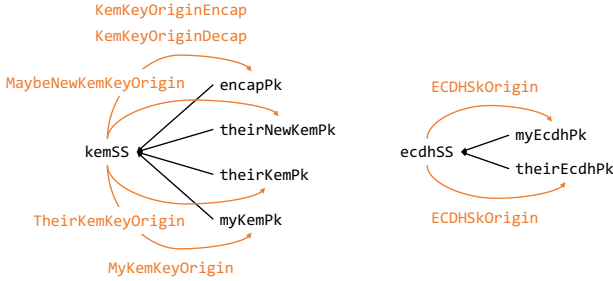


Figure 11: Loop-Jump Lemmas (orange) Related to Establishing Shared Secrets. Black arrows indicate which key material can be used to establish which shared secret.

tries to construct a trace where a message has been sent and the adversary knows it. By solving for how the adversary could learn the message, TAMARIN deduces that the adversary must know the message key used for encryption. This allows us to apply the first adversary-construction lemma `MkCompromise`. This lemma expresses that the adversary can only know the message key if they either know the respective chain key (allowing us to apply the next adversary-construction lemma) or if they access a reveal rule (contradicting our threat model assumptions directly). In the case where the adversary knows a respective smaller term, we can apply the next adversary-construction lemma, etc. Finally, the lemmas `ECDHSSCompromise` and `KemSSCompromise` directly contradict the threat model.

We proved these adversary-construction lemma using the loop-jump and variable-linking lemmas depicted in Figures 10 and 11. A sequence of variable-linking lemmas (depicted in blue) connect message to chain to root keys and to the respective KEM shared secret and ECDH shared secret (Figure 10). Loop-jump lemmas (depicted in orange) then connect the shared secrets to the asymmetric key material used to establish them (Figure 11). This allows TAMARIN to deduce that access to the shared secret requires access to the respective private key material. Beyond the lemmas depicted in Figures 10 and 11, we only use the three loop-jump lemmas `RootKeyConnectionReceive`, `RootKeyConnectionSend`, and `SessionStart`, which jump from an instance of the symmetric ratchet to the most re-

cent public key ratchet (switching from sender to receiver or receiver to sender respectively) and the session start.

We proved both agreement lemmas much like we proved secrecy, but proving agreement was much simpler. PQ3 provides agreement by signing every message. When trying to prove non-injective agreement, TAMARIN immediately finds that to violate agreement, the adversary must generate this signature themselves, which in turn requires access to the signing key. The rule that introduces the signing key, however, can directly be established using the `SessionStart` lemma as signing keys are queried only upon session start.

When attempting to prove injective agreement, TAMARIN will start by constructing a trace with two honest receive events for the same message. Using variable-linking lemmas, we can establish that these two sessions must use the same ECDH shared secret, and using the respective loop-jump lemmas, we can jump to the rule instantiation where the receiver generated their latest ECDH ephemeral key. This allows TAMARIN to derive that the two receive events must have happened in the same session (unless a pre-key was used; but this case is addressed in the lemma directly).

Finally, we only use six auxiliary lemmas not fitting the categories defined above. These lemmas simply limit TAMARIN’s search space to reduce proof construction time. For example, they show that certain events (like session start) can only occur once, or establish well-formedness conditions (for example, that the root key is a subterm of the chain key).

5.5 Discussion

5.5.1 Scope of Analysis

We do not consider session handling, long-term identity or pre-key rollover, and only consider group messaging implicitly. Our analysis covers the protocol design as described in the documentations we received from Apple. PQ3’s implementation is not part of our analysis. Furthermore, as our analysis is based on symbolic models, it abstracts away some details of the concrete implementation, like message lengths and some algorithmic details of the ciphers used.

We did not model session handling as a specification of iMessage’s session handling was not available to us. Moreover, PQ3 is not limited in its use to iMessage. Different applications may have different requirements on their session handling. Studying PQ3 in isolation is therefore desirable in its own right.

A security analysis of group aspects, such as members joining or leaving groups, is not part of PQ3 as it is a device-to-device messaging protocol. In practice, group messaging can be implemented using PQ3 by sending messages via pairwise runs of PQ3 to all group members. Such functionality is provided by an application’s session-handling layer and is thus outside of our analysis. iMessage implements group messaging using multiple, individual device-to-device sessions,

and our analysis establishes the security of each such session.

Beyond the limitations just mentioned, our formal model incorporates all details that were part of the documentation provided to us by Apple. In particular, we did not abstract away any protocol steps that participants may take.

5.5.2 Proof Effort

Our TAMARIN model comprises 32 lemmas in total. Next to the auxiliary lemmas used to prove secrecy and agreement (Section 5.4), our model includes a *sources lemma*, which aids TAMARIN in precomputation steps, and two *executability lemmas*. Executability lemmas effectively “sanity check” a protocol specification by establishing that the participants can run the protocol without adversary involvement. This enhances our confidence that the protocol model faithfully represents the protocol and that its properties do not hold trivially.

All proofs are guided by custom proof heuristics, implemented in Python, and finding the right heuristics to successfully construct proofs required substantial efforts. For example, checking the proof for the lemma formalizing injective agreement (Section 5.3.2) takes around 7 hours and requires 20 GB of RAM on a server using two Intel Xeon CPU E5-2650 v4 @ 2.20GHz. The proofs of other lemmas require up to 100 GB of RAM to be checked. Overall, we estimate that proving PQ3 took around 2.5 person-months of work.

6 Conclusions

We have used TAMARIN to formally verify the device-to-device messaging protocol PQ3. Our analysis is based on machine-checked proofs of fine-grained secrecy and authentication properties. This provides a high degree of assurance that PQ3 functions securely against an active network adversary who can selectively compromise parties, even when sufficiently powerful quantum computers become available. Additionally, the properties we prove give a detailed account of the impact that the compromise of every individual key has. Lastly, we show that TAMARIN is up to the task of reasoning about complex protocols with nested loops, and we have given a general methodology for doing this.

Future work Of particular interest would be the formal analysis of PQ3 in conjunction with session handling, as implemented for iMessage. Whether PQ3’s security guarantees as established here fully transfer to iMessage remains an open question. For example, [18] established that the Signal application may not provide post-compromise security although the protocol does due to the implementation of session handling (see Section 2.2). Furthermore, our formal model could be extended to account for IDS key roll-over, i.e., of long-term identity and pre-keys, and it could be extended to incorporate

enhanced models of cryptographic primitives, such as those suggested by [20, 16, 12, 15].

Acknowledgments

This work was supported by Apple Inc. The Werner Siemens-Stiftung (WSS) provided funding for Felix Linker’s Doctorate as part of the Centre for Cyber Trust (CECYT). We thank both for their support.

Ethics Discussion

PQ3 will be used by billions of users and proving its security benefits these users. We have found no attacks that would warrant responsible disclosure or could put users at risk. Moreover, the description of PQ3 facilitates future research on its security.

Open Science Policy

Our formal model, proofs, and a pseudocode specification of PQ3 Messaging Protocol as well as case studies illustrating our proof methodology (see Appendix A) are available at [25].

References

- [1] Joël Alwen, Sandro Coretti, and Yevgeniy Dodis. “The Double Ratchet: Security Notions, Proofs, and Modularization for the Signal Protocol”. In: *Advances in Cryptology – EUROCRYPT 2019*. 2019. DOI: [10.1007/978-3-030-17653-2_5](https://doi.org/10.1007/978-3-030-17653-2_5).
- [2] Apple Security Engineering and Architecture (SEAR). *Advancing iMessage Security: iMessage Contact Key Verification*. Oct. 2023. URL: <https://security.apple.com/blog/imessage-contact-key-verification/> (visited on 01/12/2024).
- [3] Karthikeyan Bhargavan, Abhishek Bichhawat, Quoc Do, Pedram Hosseini, Ralf Küsters, Guido Schmitz, and Tim Würtele. “DY*: A Modular Symbolic Verification Framework for Executable Cryptographic Protocol Code”. In: *6th IEEE European Symposium on Security and Privacy (EuroS&P)*. Sept. 2021. DOI: [10.1109/EuroSP51992.2021.00042](https://doi.org/10.1109/EuroSP51992.2021.00042).
- [4] Karthikeyan Bhargavan, Charlie Jacomme, Franziskus Kiefer, and Rolfe Schmidt. “Formal Verification of the PQXDH Post-Quantum Key Agreement Protocol for End-to-End Secure Messaging”. In: *33rd USENIX Security Symposium*. 2024. URL: <https://www.usenix.org/conference/usenixsecurity24/presentation/bhargavan>.

- [5] Alexander Bienstock, Jaiden Fairuze, Sanjam Garg, Pratyay Mukherjee, and Srinivasan Raghuraman. “A More Complete Analysis of the Signal Double Ratchet Algorithm”. In: *Advances in Cryptology – CRYPTO 2022*. 2022. DOI: [10.1007/978-3-031-15802-5_27](https://doi.org/10.1007/978-3-031-15802-5_27).
- [6] B. Blanchet. “A Computationally Sound Mechanized Prover for Security Protocols”. In: *2006 IEEE Symposium on Security and Privacy (S&P)*. May 2006. DOI: [10.1109/SP.2006.1](https://doi.org/10.1109/SP.2006.1).
- [7] Bruno Blanchet. “Automatic Verification of Security Protocols in the Symbolic Model: The Verifier ProVerif”. In: *Foundations of Security Analysis and Design VII*. 2014. DOI: [10.1007/978-3-319-10082-1_3](https://doi.org/10.1007/978-3-319-10082-1_3).
- [8] Olivier Blazy, Ioana Boureanu, Pascal Lafourcade, Cristina Onete, and Léo Robert. “How Fast Do You Heal? A Taxonomy for Post-Compromise Security in Secure-Channel Establishment”. In: *32nd USENIX Security Symposium*. 2023. URL: <https://www.usenix.org/conference/usenixsecurity23/presentation/blazy>.
- [9] Nikita Borisov, Ian Goldberg, and Eric Brewer. “Off-the-Record Communication, or, Why Not to Use PGP”. In: *Proceedings of the 2004 ACM Workshop on Privacy in the Electronic Society*. Oct. 2004. DOI: [10.1145/1029179.1029200](https://doi.org/10.1145/1029179.1029200).
- [10] Joppe Bos, Leo Ducas, Eike Kiltz, T Lepoint, Vadim Lyubashevsky, John M. Schanck, Peter Schwabe, Gregor Seiler, and Damien Stehle. “CRYSTALS - Kyber: A CCA-Secure Module-Lattice-Based KEM”. In: *2018 IEEE European Symposium on Security and Privacy (EuroS&P)*. Apr. 2018. DOI: [10.1109/EuroSP.2018.00032](https://doi.org/10.1109/EuroSP.2018.00032).
- [11] Ran Canetti, Palak Jain, Marika Swanberg, and Mayank Varia. “Universally Composable End-to-End Secure Messaging”. In: *Advances in Cryptology – CRYPTO 2022*. 2022. DOI: [10.1007/978-3-031-15979-4_1](https://doi.org/10.1007/978-3-031-15979-4_1).
- [12] Vincent Cheval, Cas Cremers, Alexander Dax, Lucca Hirschi, Charlie Jacomme, and Steve Kremer. “Hash Gone Bad: Automated Discovery of Protocol Attacks That Exploit Hash Function Weaknesses”. In: *32nd USENIX Security Symposium*. Aug. 2023. URL: <https://www.usenix.org/conference/usenixsecurity23/presentation/cheval>.
- [13] Katriel Cohn-Gordon, Cas Cremers, Benjamin Dowling, Luke Garratt, and Douglas Stebila. “A Formal Security Analysis of the Signal Messaging Protocol”. In: *2017 IEEE European Symposium on Security and Privacy (EuroS&P)*. Apr. 2017. DOI: [10.1109/EuroSP.2017.27](https://doi.org/10.1109/EuroSP.2017.27).
- [14] Katriel Cohn-Gordon, Cas Cremers, and Luke Garratt. “On Post-compromise Security”. In: *29th Computer Security Foundations Symposium (CSF)*. June 2016. DOI: [10.1109/CSF.2016.19](https://doi.org/10.1109/CSF.2016.19).
- [15] Cas Cremers, Alexander Dax, Charlie Jacomme, and Mang Zhao. “Automated Analysis of Protocols That Use Authenticated Encryption: How Subtle AEAD Differences Can Impact Protocol Security”. In: *32nd USENIX Security Symposium*. Aug. 2023. URL: <https://www.usenix.org/conference/usenixsecurity23/presentation/cremers-protocols>.
- [16] Cas Cremers, Alexander Dax, and Niklas Medinger. “Keeping Up with the KEMs: Stronger Security Notions for KEMs and Automated Analysis of KEM-based Protocols”. In: *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security (CCS)*. Dec. 2024. DOI: [10.1145/3658644.3670283](https://doi.org/10.1145/3658644.3670283).
- [17] Cas Cremers, Charlie Jacomme, and Philip Lukert. “Subterm-Based Proof Techniques for Improving the Automation and Scope of Security Protocol Analysis”. In: *36th Computer Security Foundations Symposium (CSF)*. July 2023. DOI: [10.1109/CSF57540.2023.00001](https://doi.org/10.1109/CSF57540.2023.00001).
- [18] Cas Cremers, Charlie Jacomme, and Aurora Naska. “Formal Analysis of Session-Handling in Secure Messaging: Lifting Security from Sessions to Conversations”. In: *32nd USENIX Security Symposium*. 2023. URL: <https://www.usenix.org/conference/usenixsecurity23/presentation/cremers-session-handling>.
- [19] Rune Fiedler and Felix Günther. *Security Analysis of Signal’s PQXDH Handshake*. 2024. URL: <https://eprint.iacr.org/2024/702>.
- [20] Dennis Jackson, Cas Cremers, Katriel Cohn-Gordon, and Ralf Sasse. “Seems Legit: Automated Analysis of Subtle Attacks on Protocols That Use Signatures”. In: *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security (CCS)*. Nov. 2019. DOI: [10.1145/3319535.3339813](https://doi.org/10.1145/3319535.3339813).
- [21] Frederic Jacobs. *Invited Talk: Designing iMessage PQ3: Quantum-Secure Messaging at Scale*. Mar. 2024. URL: <https://www.youtube.com/watch?v=RVbHElGe518> (visited on 05/29/2024).
- [22] Nadim Kobeissi, Karthikeyan Bhargavan, and Bruno Blanchet. “Automated Verification for Secure Messaging Protocols and Their Implementations: A Symbolic and Computational Approach”. In: *2017 IEEE European Symposium on Security and Privacy (EuroS&P)*. Apr. 2017. DOI: [10.1109/EuroSP.2017.38](https://doi.org/10.1109/EuroSP.2017.38).

- [23] Hugo Krawczyk. “Cryptographic Extraction and Key Derivation: The HKDF Scheme”. In: *Advances in Cryptology – CRYPTO 2010*. 2010. DOI: [10.1007/978-3-642-14623-7_34](https://doi.org/10.1007/978-3-642-14623-7_34).
- [24] Ehren Kret and Rolfe Schmidt. *The PQXDH Key Agreement Protocol*. Tech. rep. Revision 3. May 2023. URL: <https://signal.org/docs/specifications/pqxdh/pqxdh.pdf>.
- [25] Felix Linker, Ralf Sasse, and David Basin. *A Formal Analysis of Apple’s iMessage PQ3 Protocol*. Zenodo. Jan. 2025. DOI: [10.5281/zenodo.14710688](https://doi.org/10.5281/zenodo.14710688). (Visited on 01/23/2025).
- [26] G. Lowe. “A Hierarchy of Authentication Specifications”. In: *Proceedings 10th Computer Security Foundations Workshop*. June 1997. DOI: [10.1109/CSFW.1997.596782](https://doi.org/10.1109/CSFW.1997.596782).
- [27] Moxie Marlinspike and Trevor Perrin. *The X3DH Key Agreement Protocol*. Tech. rep. Revision 1. Nov. 2016. URL: <https://signal.org/docs/specifications/x3dh/x3dh.pdf>.
- [28] Simon Meier, Benedikt Schmidt, Cas Cremers, and David Basin. “The TAMARIN Prover for the Symbolic Analysis of Security Protocols”. In: *Computer Aided Verification (CAV)*. 2013. DOI: [10.1007/978-3-642-39799-8_48](https://doi.org/10.1007/978-3-642-39799-8_48).
- [29] *Module-Lattice-based Key-Encapsulation Mechanism Standard*. Aug. 2023. DOI: [10.6028/NIST.FIPS.203.ipd](https://doi.org/10.6028/NIST.FIPS.203.ipd).
- [30] Michele Mosca and Marco Piani. *Quantum Threat Timeline*. Tech. rep. Global Risk Institute, Dec. 2023. URL: <https://globalriskinstitute.org/publication/2023-quantum-threat-timeline-report/>.
- [31] Vinnie Moscaritolo, Gary Belvin, and Phil Zimmermann. *Silent Circle Instant Messaging Protocol*. Tech. rep. Dec. 2012. URL: <https://netzpolitik.org/wp-upload/SCIMP-paper.pdf>.
- [32] Trevor Perrin and Moxie Marlinspike. *The Double Ratchet Algorithm*. Tech. rep. Revision 1. Nov. 2016. URL: <https://signal.org/docs/specifications/doubleratchet/doubleratchet.pdf>.
- [33] Benedikt Schmidt, Simon Meier, Cas Cremers, and David Basin. “Automated Analysis of Diffie-Hellman Protocols and Advanced Security Properties”. In: *25th Computer Security Foundations Symposium (CSF)*. June 2012. DOI: [10.1109/CSF.2012.25](https://doi.org/10.1109/CSF.2012.25).
- [34] Douglas Stebila. *Security Analysis of the iMessage PQ3 Protocol*. 2024. URL: <https://eprint.iacr.org/2024/357>.
- [35] Nik Unger, Sergej Dechand, Joseph Bonneau, Sascha Fahl, Henning Perl, Ian Goldberg, and Matthew Smith. “SoK: Secure Messaging”. In: *2015 IEEE Symposium on Security and Privacy (S&P)*. May 2015. DOI: [10.1109/SP.2015.22](https://doi.org/10.1109/SP.2015.22).
- [36] Nik Unger and Ian Goldberg. “Deniable Key Exchanges for Secure Messaging”. In: *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security (CCS)*. Oct. 2015. DOI: [10.1145/2810103.2813616](https://doi.org/10.1145/2810103.2813616).

A Proof Methodology

In Section 5.4, we presented our proof methodology, specialized to how we applied it to the PQ3 Messaging Protocol. In this section, we will describe this methodology in more detail and in its generality: i.e., how one could apply it to other protocols with similar structure.

TAMARIN provides general support for handling loops, based on induction and injective facts, and we begin our account by explaining them. We afterwards introduce two minimal TAMARIN theories that illustrate the issues of nested loops and ghost sessions (see Section 5.4), but on a smaller and simpler scale. These theories will also help illustrate the reasoning behind the loop-jump, adversary-construction and variable-linking lemmas that we have seen. Finally, we present the resulting proof methodology.

A.1 Handling Loops in TAMARIN

A.1.1 Induction

TAMARIN analyzes formulas directly by backward search, as explained in Section 5.1, or by induction. When TAMARIN attempts to prove a formula φ by induction, it rewrites it into the form

$$\text{BC}(\varphi) \wedge (\text{IH}(\varphi) \implies \varphi).$$

The first conjunct, $\text{BC}(\varphi)$, is the base case, and it requires proving φ on the empty trace. The second conjunct, $\text{IH}(\varphi) \implies \varphi$, is the induction step, which requires proving φ on the last element of the trace, where φ is assumed on all previous steps of the trace. $\text{BC}(\varphi)$ is defined as φ , where every formula of the form $f@i$ is replaced with \perp . For example, for a formulation of secrecy such as

$$\forall m, t. \text{Sent}(m)@t \implies \neg(\exists x. K(m)@x),$$

this replacement results in

$$\forall m, t. \perp \implies \neg(\exists x. \perp).$$

$\text{IH}(\varphi)$ is defined as φ but every quantified temporal variable is asserted to not be the last time point. This is done using the

special predicate *last*, which is true if and only if it is provided the last time point as argument. For example, the induction hypothesis of secrecy as defined above would become

$$\forall m, t. \text{Sent}(m)@t \implies \neg(\exists x. K(m)@x \wedge \neg \text{last}(x)) \vee \text{last}(t).$$

After translating ϕ into its inductive form, TAMARIN attempts to prove it as any other formula. Effectively, it attempts to prove the base case and induction step separately, and the induction hypothesis IH is made available (like an auxiliary lemma) in the branch proving the induction step.

In practice, induction is used to prove properties of protocols with loops. However, one can only prove properties of loops by induction when the loops are expressed in terms of facts that appear repeatedly in the protocol's trace. Take the above translation of secrecy as an example. In the induction step, the induction hypothesis becomes effectively vacuous as long as the fact *Sent*(*m*) only occurs at a last time point *t*. In that case, the second disjunct on the right-hand side of the implication will apply, whereas one usually requires the first disjunct to apply to make progress on a proof. Only when we can introduce a new *Sent* fact in the trace that does not occur at the last time point can we use the induction hypothesis.

In particular, and applied to loops, this means two things. (1) Induction can only be applied to formulas that express invariants of loops, but not to formulas that express something that holds after a loop has stopped. A loop will only end once, which makes it impossible to introduce a second end of the loop not occurring at the last time point. (2) Induction cannot be used to prove properties for outer loops without further auxiliary lemmas. When we attempt to prove properties of outer loops, TAMARIN will always also consider the case that a step in the outer loop was preceded by an inner loop of unbounded length. Also in these cases, the fact referenced in the induction hypothesis (the outer loop step) will only occur at the last time point. For both of these reasons, induction must be applied with care and cannot be blindly applied to prove arbitrary properties of protocols with loops.

A.1.2 Injective Facts

Injective facts are commonly used to model loops in TAMARIN. They are defined as facts that (for a fixed first argument) can occur only once in the global state. We call an injective fact's first argument its *loop identifier*. If a fact satisfies the following constraints, it is automatically detected as injective by TAMARIN:

- It is not a persistent fact.
- Its loop identifier is a fresh term.
- It never occurs more than once in a rule's conclusion with the same loop identifier.

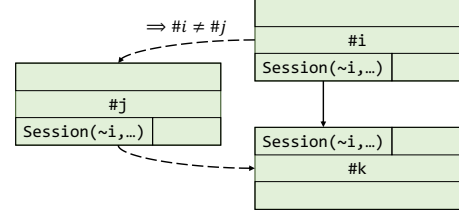


Figure 12: Illustration of a contradiction from an injective instance. The solid arrow indicates premise consumption. Dashed arrows indicate time ordering, i.e., the rule at #j must be applied after #i but before #k. The order of time points requires that #j and #i must be unified as #k consumes a session fact with the matching ID $\sim i$. However, #j must occur strictly after #i, which contradicts this unification occurrence. There is a symmetric case where #j and #k must be unified because #j and #i share a premise.

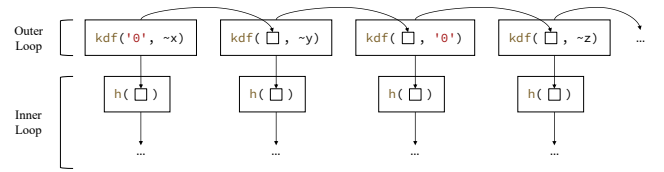


Figure 13: Nested loop example.

- Whenever it occurs in a rule's conclusion, either (a) its loop identifier is freshly generated (the loop has started), or (b) it occurs in the rule's premise with the same loop identifier (a loop step is taken).

Injective facts allow TAMARIN to derive contradictions by exploiting that all injective facts with a shared loop identifier must be linearizable. In particular, a loop step can never occur between two directly connected loop steps (as illustrated in Figure 12). It is possible to prove properties of loops without using injective facts, but using injective facts can drastically simplify proofs, so it is generally advisable to make use of this heuristic.

A.2 Proof Methodology by Example

In what follows, we illustrate the challenges encountered when constructing proofs about nested loops using three simple, minimal theories. We will introduce these theories and our proof methodology on an intuitive level. For full details, see our artifact that provides all the theory files [25].

Nested Loop Example Figure 13 provides an illustration of a nested loop from the nested-loop theory in [25]. The loop models a participant's key derivation similar to the key derivation used in PQ3. The inner loop applies a hash function repeatedly, using non-determinism to leave open how often, to a value derived from a KDF, which we call the *seed*. The outer loop updates the seed. When the outer loop starts, the

seed is derived from a zero-byte sequence and a fresh value. At every outer-loop iteration, the seed is derived from the previous seed and either a zero-byte sequence or a fresh value (determined non-deterministically). The adversary can access all fresh values used in this loop using a reveal oracle.

The key derivation in this theory is similar to the key derivation in PQ3 when focussing on KEM shared secrets. The inner loop abstracts from the chain and message key derivation, while the outer loop abstracts from establishing new KEM shared secrets. At the end of this section, we show how a simpler version of this theory captures the essence of the double ratchet construction, i.e., repeatedly establishing fresh Diffie-Hellman shared secrets.

Now consider proving a simple key secrecy lemma: Every key established in the inner loop either remains confidential, or the most recently used fresh value in the outer loop was revealed to the adversary. To prove this lemma, we establish four auxiliary lemmas:

Outer Loop Step Every step in the inner loop must be preceded by a step in the outer loop. This lemma can be proven straightforwardly by induction.

Fresh Seed Source For every step in the outer loop ratchet, there must be a step in that ratchet deriving the seed that most recently was derived from a fresh value. We can prove this by induction using the *outer loop step* lemma. When proving this lemma, there is only one case that does not immediately lead to a contradiction. When the outer-loop step was immediately preceded by an inner-loop step, there is neither a contradiction nor does the induction hypothesis apply. In that case, we can apply the *outer loop step* to jump to the previous outer-loop step, which will either have used a fresh value to derive its seed (direct contradiction) or a zero-byte sequence (but since it is an outer-loop step, the induction hypothesis applies).

Seed Construction When the adversary derived a seed, they must have used the seed most recently constructed using a fresh value in that derivation. Similarly to the previous lemma, we can only prove this lemma using *outer loop step* as an auxiliary lemma.

Key Construction When the adversary derived a key established in the inner loop, they must have used the most recently generated seed. This lemma can also be proven by induction straightforwardly.

Using these four auxiliary lemmas, Tamarin automatically proves key secrecy. Note that all auxiliary lemmas above are proven using induction but the key secrecy lemma is not.

We can describe above lemmas in more general terms.

Outer Loop Step This lemma “jumps to” the most recent step of the outer loop and allows one to skip unrolling the inner loop infinitely.

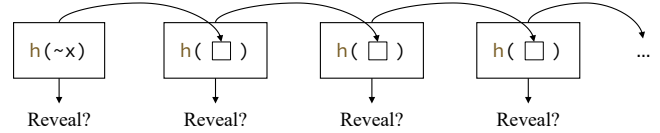


Figure 14: Revealing loop example.

Fresh Seed Source This lemma “jumps to” the step in the outer loop that introduces the “relevant term” (in our case, the fresh term used instead of the zero-byte sequence). It allows one to skip unrolling the outer loop infinitely.

Seed/Key Construction These lemmas link the adversary’s knowledge of the key and the seed to the adversary’s knowledge of the respectively “next term” (the previous seed and the seed established using a fresh value).

Note that lemmas that “jump to” the relevant term in the outer loop (here *fresh seed source*) and that connect terms from the outer loop (here *seed construction*) are only required when there can be unboundedly many outer-loop steps until the relevant step is reached (the relevant step being the one where the fresh term is introduced). To illustrate this point, we also provide a second nested loop theory (*nested-loop-simple* in [25]) that always uses a fresh value to establish the respective seed in the outer loop. In this theory, the lemmas *fresh seed source* and *seed construction* are not needed and unrolling the outer loop is sufficient.

This simplified theory is similar to the key derivation of PQ3 when focussing on the Diffie-Hellman shared secrets and ignoring the KEM shared secrets. It is also similar to the double ratchet as used in Signal [32]. This simplified theory suggests that although the double ratchet construction used in Signal employs a nested loop, no inductive properties must be proven about the outer loop.

Revealing Loop Example The *revealing-loop* theory provided in [25] illustrates the challenges of proving properties of PQ3 when considering the reveal of root, chain, and message keys explicitly and thus illustrates how we addressed the problem of ghost sessions. This example is much simpler than the previous ones and is illustrated in Figure 14. A participant starts a loop in which they repeatedly apply a hash function h to some initial seed $\sim x$. Critically, the model allows the adversary to access any derived value and $\sim x$ using a reveal oracle.

Again, we show how to prove a simple key secrecy lemma: Every hash value derived is secret, i.e., not known by the adversary, unless the adversary compromised any of the previous, intermediate values. To prove this lemma, we require two auxiliary lemmas:

Loop Start Every loop that computes a hash based on some initial seed $\sim x$ is started, sampling $\sim x$. This lemma is straightforward to prove by induction.

Seeds Match If the computed hashes of two loops are identical, their seed must be identical too. Conceptually, this lemma is again simple and can be proven straightforwardly using induction.

With both these lemmas, TAMARIN can prove key secrecy of this example theory using induction. The two auxiliary lemmas help TAMARIN address the case that the adversary learns the hash in question from a ghost session. Using the second lemma, TAMARIN can connect the two sessions using a shared, fresh term. Then, using the first lemma, TAMARIN can instantiate the start of the loop where this shared, fresh term was sampled. From that, TAMARIN can deduce that both sessions must be the same; this enables it to apply the induction hypothesis and to prove the key secrecy lemma.

Again, we can generalize these auxiliary lemmas.

Loop Start This lemma is conceptually similar to the *outer loop step* lemma from the previous example and introduces no new kinds of lemmas.

Seeds Match This lemma links the computation of a value in a loop to the inputs to this computation, not determined by a loop (here, the seed).

Summary With the previous two examples, we showed how to handle nested loops and ghost sessions in TAMARIN. All auxiliary lemmas in these two theories match the three types of lemmas introduced in Section 5.4, which we briefly recapitulate:

Loop-Jump Lemmas These lemmas allow one to skip unrolling the steps of a (nested) loop. Examples: *outer loop step*, *fresh seed source*, *loop start*.

Adversary-Construction Lemmas These lemmas establish that for the adversary to construct one term, they must use another term. Examples: *key construction*, *seed construction*.

Variable-Linking Lemmas These lemmas establish that for two instances of the same fact using two variables a and b , if both facts have the same value for a , they must have the same value for b . Example: *seeds match*.

A.3 Proof Methodology in General

Our proof methodology applies to theories that (i) use asymmetric cryptography to establish shared secrets, which in turn are used to derive symmetric encryption keys, (ii) include a (nested) loop computing these symmetric encryption keys, and (iii) for which trace formulas are to be proven of the following form:

$$\forall \vec{x}. C(\vec{x}) \implies (\neg) \exists \vec{y}. P(\vec{x}; \vec{y}) \vee T(\vec{x}).$$

For all traces that satisfy some context C , there exists (or does not exist) an instance of P bound to that context ($;$ denotes vector concatenation), unless T (which specifies the threat model) applies. For example, for secrecy, C could be “a message was sent”, P could be “the adversary learned that message” (in this case, non-existence would be proven), and T could be “the message encryption keys were revealed to the adversary.”

We require that the protocol is modelled such that there is a single fact that models the protocol’s (nested) loop, which we call the *loop fact*. This allows us to do two things: (a) exploit TAMARIN’s heuristics for injective facts (see Section A.1.2), (b) clearly identify and relate looping variables, which will be critical to our proof methodology. We identify the loop fact’s variables by their position in the fact, and there will generally be two kinds of variables: shared and derived secrets and key material used for establishing the shared and derived secrets. We relate the shared and derived secrets by a strict partial order. That order is defined as the smallest order closed under transitivity for which one variable a is smaller than another variable b if there is a state-transition rule that updates b using a . We say that a loop fact variable can *grow unboundedly* if there is no bound on the size of the terms that the variable can be unified with, for all ground-instantiated traces.

For example, our PQ3 model uses a `Session` fact to model the double-ratchet steps performed by a participant. The order on the shared and derived secrets is depicted in Figure 10 (the black arrows). For PQ3, the variables `msgKey`, `chainKey`, and `rootKey` can grow unboundedly.

Overall, our proof methodology has four steps.

1. For each of the loop fact’s (possibly nested) loops, write a loop-jump lemma that connects a loop instance to its beginning, i.e., to the beginning of the loop overall or to the transition from an outer loop to the respective next inner loop.
2. Identify all the loop facts’ shared and derived secret variables that can grow unboundedly (e.g., `msgKey` for PQ3). For each of these variables, write a variable-linking lemma that connects them to the variables that are directly smaller than them (e.g., a message to a chain key). For some of these unboundedly growing variables, it might additionally be necessary to write a loop-jump lemma that jumps to the rule application assigning a new value to the respective smaller variable. In our experience, this is the case for variables that are updated non-deterministically (i.e., `kemSS` for PQ3).
3. Identify all the loop fact’s shared and derived secret variables that do not grow unboundedly. Typically, these variables will store shared secrets established using asymmetric cryptography. For each of these variables, write loop-jump lemmas that link the usage of that variable to the instantiation of the respective asymmetric key material

```

1 functions: hkdf/2, suffix/1, prefix/1, concat/2, h/1
2
3 equations: concat(prefix(x), suffix(x)) = x
4
5 functions: pqpk/1, encap/2, decap/2
6 equations: decap(encap(k, pqpk(sk)), sk) = k
7
8 functions: default/2, Just/1, None/0, unjust/1
9 equations: default(Just(v), t) = v,
10             default(None, v) = v,
11             unjust(Just(t)) = t

```

Figure 15: Custom functions and equations defined in our formal model.

used to establish the secret. The details of these loop-jump lemmas depend on the protocol specification. For example, the lemmas that link the ECDH shared secret to the respective ECDH keys substantially differ from the lemmas that link the KEM shared secret to the respective encapsulation key.

4. Finally, for all variables connected by variable-linking lemmas, write an adversary-construction lemma that states that in order for the adversary to know the contents of the respective larger variable, they must have either violated the threat model or know the respective smaller variables.

Following these steps, one would write the lemmas `RootKeyConnectionReceive`, `RootKeyConnectionSend`, and `SessionStart` in Step 1, the lemmas depicted in Figure 10 in Step 2, the lemmas depicted in Figure 11 in Step 3, and the lemmas depicted in Figure 9 in Step 4.

B Equational Theory for Protocol Model

We use TAMARIN’s built-in equational theories for signing, symmetric encryption, and Diffie-Hellman key exchange. These respectively model digital signatures, symmetric encryption under message keys, and ECDH key exchanges. We additionally use TAMARIN’s natural numbers theory to model message counters.

In addition to these built-in theories, we specify some custom functions and equations, shown in Figure 15. First, we specify the functions `hkdf`, `suffix`, and `prefix` for key derivation. The function `hkdf` models an HMAC-based key derivation function [23] and takes two arguments: the first is the source of entropy and the second is a domain-separating tag or salt. The `prefix` and `suffix` functions are used for chain and root key derivations, which are derived by splitting a bit-string into a prefix and suffix of equal length. The function `concat` allows one to recover a value given its prefix and suffix. We do not need to use `concat` in the rules modeling the protocol roles of regular parties in our model, but the adversary can use it to reconstruct a value from the prefix and suffix. Additionally, we specify the unary function `h` to

model the pre-key hash used during session establishment, see Section 4.

The functions `pqpk`, `encap`, and `decap` model KEM encapsulation and follow the standard symbolic model for asymmetric encryption. Finally, we use the wrapper function `Just` and the constant `None` to model optional values. The function `default` (together with the accompanying equations) unpacks an optional value or replaces it with a default. For example, we use `Just` and `None` to wrap values that are only sent optionally, e.g., the pre-key hash. The function `unjust` allows the adversary to access the contents of any `Just` value they intercept.