

Automated Analysis of Protocols that use Authenticated Encryption: How Subtle AEAD Differences can impact Protocol Security

Cas Cremers¹, Alexander Dax^{1,3}, Charlie Jacomme², and Mang Zhao^{1,3}

¹CISPA Helmholtz Center for Information Security, Germany

²Inria Paris, France

³Saarland University

Abstract

Many modern security protocols such as TLS, WPA2, WireGuard, and Signal use a cryptographic primitive called Authenticated Encryption (optionally with Authenticated Data), also known as an AEAD scheme. AEAD is a variant of symmetric encryption that additionally provides authentication. While authentication may seem to be a straightforward additional requirement, it has in fact turned out to be complex: many different security notions for AEADs are still being proposed, and several recent protocol-level attacks exploit subtle behaviors that differ among real-world AEAD schemes.

We provide the first automated analysis method for protocols that use AEADs that can systematically find attacks that exploit the subtleties of the specific type of AEAD used. This can then be used to analyze specific protocols with a fixed AEAD choice, or to provide guidance on which AEADs might be (in)sufficient to make a protocol design secure. We develop generic symbolic AEAD models, which we instantiate for the Tamarin prover. Our approach can automatically and efficiently discover protocol attacks that could previously only be found using manual inspection, such as the Salamander attack on Facebook’s message franking, and attacks on SFrame and YubiHSM. Furthermore, our analysis reveals undesirable behaviors of several other protocols.

1 Introduction

Authenticated Encryption (AE) and Authenticated Encryption with Associated Data (AEAD) are some of the most commonly used cryptographic building blocks. AEAD primitives are built from symmetric encryption primitives and augmented with authentication mechanisms. Their applications include the vast majority of encrypted internet data, such as in TLS, WPA2 from IEEE 802.11 (WiFi), WireGuard, and by messaging apps such as Signal or WhatsApp. For example, in TLS, TLSCiphertext is constructed from an AEAD applied to a header and payload: both are authenticated, but only the payload is encrypted, and the plaintext header includes the content type and the ciphertext length.

While AEADs are ubiquitous in modern secure communications, there is no commonly agreed “strong” security notion that they should satisfy. In fact, the current landscape of security notions for AEADs is rather chaotic: there are many proposed frameworks and security notion variants [1–3, 7–9, 15, 17, 26, 31, 32, 43, 51, 55]. For some of these notions, their implication relations are known [7], but many of them are hard to compare for technical reasons.

In reality, there are good examples of recent protocol attacks that exploit subtle properties of concrete AEAD schemes, such as [26, 33, 42]. These have all been found through manual inspection of the protocol and knowledge of the particular AEAD scheme, such as exploiting the reuse of a nonce, a number meant to be used only once. We would like to formally prove the absence of such attacks: i.e., that a protocol (e.g., TLS, WhatsApp, WPA2), when instantiated with a specific AEAD (e.g., AES-GCM), satisfies a desired security notion. At a methodological level, these attacks can be hard to model because they require a methodology that is not only precise enough to capture AEAD details (e.g. impact of nonce reuse) but also scales well enough to allow for modeling the often complex possible executions of a protocol that determine whether such attack requirements can be met.

In this work, we develop the first systematic methodology for the automated analysis of security protocols that can find attacks that leverage subtle behaviors of concrete AEAD schemes, or show their absence. For our methodology, we leverage the TAMARIN prover [45], a symbolic protocol analysis tool, which we augment with novel fine-grained models of AEAD primitives. This tool choice enables us to analyze several non-trivial protocols. We also considered using tools in the computational model (e.g., [5, 6, 12]), but these currently do not yet scale to the complexity of the protocols we are interested in, and cannot find attacks.

One of the challenges in developing our models is that they require finding middle ground between theoretical security notions, weaknesses exploited in practice, and suitability for automated analysis. We identify the core theoretical and practical concerns of AEADs, and use these to develop our

generic symbolic AEAD models. We notably identify how the *collision resistance* of AEAD is a central issue of their design: we illustrate how a lack of it leads to multiple attack classes, and how satisfying collision resistance implies that many existing security notions are met. In our protocol case studies, we rediscover previously reported attack classes, such as *accountability* or *authentication*, but also identify a new attack class concerning *content agreement* in group messaging scenarios: can a dishonest group member send a single message that will be interpreted differently by multiple parties?

Contributions Our main contributions are the following:

- We develop the first systematic automated methodology for analyzing security protocols that takes the subtle properties of specific AEAD instantiations into account.
- In case studies, we show our methodology effectively rediscovers known attacks on several protocols, including YubiHSM [42], Facebook’s Message Franking [26], and SFrame [33]. We also rediscover a theoretical attack variant on Facebook’s Message Franking first mentioned in [31]. Moreover, our analysis uncovers unexpected behavior in WebPush [56], Whatsapp Group Messaging [59], and Scuttlebutt [54].
- We formally prove the missing or conjectured relations between existing AEAD security notions w.r.t. collision resistance, completing the picture in the domain.

We provide the full formal TAMARIN models and analysis scripts at [20].

Overview We first give background on AEADs in Section 2. We build the foundations for our methodology by revisiting the AEAD landscape and real-world attack patterns in Section 3, and prove some missing relations between AEAD notions. We then develop our symbolic modeling and analysis approach in Section 4. We evaluate our approach on several protocol case studies and present details on selected case studies in Section 5. We discuss limitations in Section 6 and describe further related work in Section 7. We conclude in Section 8.

We give all proofs and details on the remaining case studies in the full version [19].

2 Background on AEADs and protocol attacks

The modern approach to protecting privacy and authenticity of messages is to use authenticated encryption. This primitive evolved as a variant of symmetric encryption that highly efficiently offers two additional properties that are useful in real-world applications: authentication of the encrypted data, and concurrent authentication of some plaintext data. One could in theory achieve this by using a combination of symmetric encryption and MACs, but AEADs are much more

efficient and ensure the authenticated part is strictly bound to the encrypted-and-authenticated part.

The efficient combination of these constructions has proven to be surprisingly intricate. In standard encryption it is desirable that if Alice encrypts the same message twice, the attacker cannot tell this from the ciphertexts. This property is typically achieved by ensuring that when encrypting the next message, some data x is integrated that was not used before, ensuring uniqueness of the ciphertext. As we will see later, this sometimes fails for various reasons, causing some x to be re-used. Intuitively, we would hope that while the attacker is now able to detect message re-encryption, there should not be any further negative consequences. For some provably secure AEADs, it turns out that the situation is worse.

Furthermore, much like symmetric encryption, AEADs in real-world deployments typically construct and send ciphertext incrementally. To start decrypting partial ciphertexts as soon as possible, the syntax or API of some AEADs allows to decouple decryption from the verification of authentication, which can be helpful but also cause problems elsewhere.

Historically, Authenticated Encryption with Associated Data (AEAD) was introduced by Rogaway [51] to entwine privacy and authenticity for both messages and headers in a single and compact mode. The definition of AEAD is given in the nonce-based pattern, where the nonce is named after the *number* that are supposed to be used only *once*. The nonce-based AEADs are expected to relax the security requirements of the randomized or counter-based pattern – ensuring no reuse of the nonce during the encryption is sufficient for privacy and authenticity. Moreover, Bellare and Hoang [8] initialized the study on binding keys and other optional inputs to the ciphertexts.

In this section, we first recall the formal syntax of AEADs and the canonical privacy and integrity definitions in Section 2.1. In Section 2.2, we review the main attacks on protocols based on subtle AEADs behaviors and weaknesses. In Section 2.3, we summarize and classify the main AEAD frameworks in the literature.

2.1 Formal AEAD syntax and requirements

Notations. We consider that all algorithms defined in this paper are parameterized implicitly by the security parameter. Let $s \leftarrow S$ denote sampling a variable s uniformly at random from a set S . Let $x \leftarrow X$ denote the execution of a probabilistic algorithm X followed by assigning the output to a variable x . We write $x \leftarrow X$ if the algorithm X is deterministic. Let \perp denote an special error symbol that is not included in any set in this paper. We use $_$ to denote a variable that is irrelevant.

In the presentation of this work, we focus on the definition of nonce-based AEADs. The main reason, besides reducing complexity, is that randomness- and counter-based AEADs can be cast as instances of nonce-based AEADs. Looking ahead, all our symbolic models will cover nonce-

$\text{Exp}_{\text{AEAD}}^{\text{IND\$-CPA}}:$ 1 $b \leftarrow \{0, 1\}$ 2 $\mathcal{L}_c \leftarrow \emptyset$ 3 $k \leftarrow \text{KGen}()$ 4 $b' \leftarrow \mathcal{A}^{\text{ENC}}()$ 5 return $\llbracket b = b' \rrbracket$	$\text{Exp}_{\text{AEAD}}^{\text{IND\$-CCA}}:$ 1 $b \leftarrow \{0, 1\}$ 2 $\mathcal{L}_c \leftarrow \emptyset$ 3 $k \leftarrow \text{KGen}()$ 4 $b' \leftarrow \mathcal{A}^{\text{ENC,DEC}}()$ 5 return $\llbracket b = b' \rrbracket$	$\text{ENC}(N, H, m):$ 6 if $(N, H, m, _) \in \mathcal{L}_c$ 7 return \perp 8 if $b = 0$ 9 $c \leftarrow \text{Enc}(k, N, H, m)$ 10 else $c \leftarrow \{0, 1\}^{\ell(m)}$ 11 $\mathcal{L}_c \leftarrow \mathcal{L}_c \cup \{(N, H, m, c)\}$ 12 return c	$\text{DEC}(N, H, c):$ 13 if $(N, H, _, c) \in \mathcal{L}_c$ 14 return \perp 15 $m \leftarrow \text{Dec}(k, N, H, c)$ 16 if $m \neq \perp$ 17 $\mathcal{L}_c \leftarrow \mathcal{L}_c \cup \{(N, H, m, c)\}$ 18 return m
---	---	---	--

Figure 1: IND\\$-CPA and IND\\$-CCA security for an AEAD = (KGen, Enc, Dec) scheme.

$\text{Exp}_{\text{AEAD}}^{\text{CTI-CPA}}:$ 1 $\mathcal{L}_c \leftarrow \emptyset$ 2 $k \leftarrow \text{KGen}()$ 3 $(N, H, c) \leftarrow \mathcal{A}^{\text{ENC}}()$ 4 if $c \in \mathcal{L}_c$ 5 return 0 6 return $\llbracket \text{Dec}(k, N, H, c) \neq \perp \rrbracket$	$\text{Exp}_{\text{AEAD}}^{\text{CTI-CCA}}:$ 1 $\mathcal{L}_c \leftarrow \emptyset$ 2 $k \leftarrow \text{KGen}()$ 3 $(N, H, c) \leftarrow \mathcal{A}^{\text{ENC,DEC}}()$ 4 if $c \in \mathcal{L}_c$ 5 return 0 6 return $\llbracket \text{Dec}(k, N, H, c) \neq \perp \rrbracket$	$\text{ENC}(N, H, m):$ 7 $c \leftarrow \text{Enc}(k, N, H, m)$ 8 $\mathcal{L}_c \leftarrow \mathcal{L}_c \cup \{c\}$ 9 return c $\text{DEC}(N, H, c):$ 10 return $\text{Dec}(N, H, c)$
--	--	---

Figure 2: CTI-CPA and CTI-CCA security for an AEAD = (KGen, Enc, Dec) scheme.

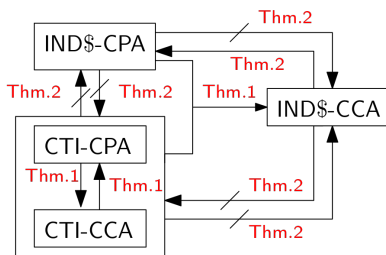


Figure 3: The relation between integrity and privacy for AEAD.

based AEADs and can then trivially be used to model randomized or nonce-based AEADs.

Definition 1 ([51]). Let Key, Nonce, Header, Message, Ciphertext respectively denote the space of keys, nonces, headers (aka. associated data), messages, and ciphertexts. An authenticated encryption with associated data scheme $\text{AEAD} = (\text{KGen}, \text{Enc}, \text{Dec})$ is a tuple of algorithms where

- KGen the key generation algorithm outputs a symmetric key $k \in \text{Key}$, i.e., $k \leftarrow \text{KGen}()$.
- Enc the encryption algorithm inputs a key $k \in \text{Key}$, a nonce $N \in \text{Nonce}$, a header $H \in \text{Header}$, and a message m and (deterministically) outputs a ciphertext c , i.e., $c \leftarrow \text{Enc}(k, N, H, m)$.
- Dec the decryption algorithm inputs a key $k \in \text{Key}$, a nonce $N \in \text{Nonce}$, a header $H \in \text{Header}$, and a ciphertext $c \in \text{Ciphertext}$ and deterministically outputs a message $m \in \text{Message} \cup \{\perp\}$, i.e., $m \leftarrow \text{Dec}(k, N, H, c)$.

Over such schemes, the N, H and ciphertext c need to be sent over the network¹, and the correctness of the scheme

¹We stress that AEAD schemes can be used offline in practice, where nonces and headers both are hidden from attackers' view. However, this paper focuses on a more common case where the attackers might have access to the nonce and headers, e.g., which are transmitted over network.

requires that the decryption of a ciphertext with the same parameters N, H, k indeed returns the plaintext. We assume that the decryption with inputs outside the corresponding spaces must output \perp .

The two core security guarantees are integrity and privacy.

Definition 2 (Privacy [51]). We say an $\text{AEAD} = (\text{KGen}, \text{Enc}, \text{Dec})$ is ϵ -IND\\$-CPA secure, if the below defined advantage of any attacker \mathcal{A} against $\text{Exp}_{\text{AEAD}}^{\text{IND\$-CPA}}$ experiment in Fig. 1 is bounded by:

$$\text{Adv}_{\text{AEAD}}^{\text{IND\$-CPA}} := |\Pr[\text{Exp}_{\text{AEAD}}^{\text{IND\$-CPA}}(\mathcal{A}) = 1] - \frac{1}{2}| \leq \epsilon$$

Definition 3 (Integrity [51]). We say an $\text{AEAD} = (\text{KGen}, \text{Enc}, \text{Dec})$ is ϵ -CTI-CPA secure, if the below defined advantage of any attacker \mathcal{A} against $\text{Exp}_{\text{AEAD}}^{\text{CTI-CPA}}$ experiment in Fig. 2 is bounded by:

$$\text{Adv}_{\text{AEAD}}^{\text{CTI-CPA}} := \Pr[\text{Exp}_{\text{AEAD}}^{\text{CTI-CPA}}(\mathcal{A}) = 1] \leq \epsilon$$

Both for integrity and privacy, we can define two security variants, CTI-CCA and IND\\$-CCA, based on whether the attacker also has access to a decryption oracle during the experiment, see e.g., the definition of the experiment for CTI-CCA in Fig. 2. We summarize the well-known relations in Fig. 3, and give the corresponding theorems in the full version of this paper [19].

2.2 Historical real-world protocol attacks exploiting AEADs

Beside the well-studied privacy and integrity, some recent attacks against practical application protocols suggest that the underlying AEADs need an evolution to meet stronger security guarantees. We identify six main classes of these

protocol attacks that have occurred in the wild and briefly describe their high-level requirements.

A1 Nonce reuse attacks - While nonces are expected to be used only once, this can fail in practice for three main reasons. First, protocol designs might aim to establish nonces, but their complex state machines may hide edge cases in which they are in fact reused, as in e.g., WPA2 [57]. Second, the generation of nonces might involve external sources, which may be unreliable, e.g., YubiKey [42] or Trustzone [55]. Third, implementations may be flawed. For example, the Zerologon attack [62] notably exploited the fact that the nonce underlying the AES-CBF8 mode in Microsoft Netlogon protocol is a constant string of zero bits. The encryption of a block of zero bits equals to 0^{16} with probability $1/256$ for any key k , breaking the authentication of windows servers.

A2 Padding oracle attacks [58] - Many AEADs and symmetric encryption schemes are constructed from block ciphers and require the length of input messages to be multiple of a fixed value. Messages whose length is not a multiple are extended before encryption using a so-called *padding scheme*. These can enable plaintext recovery attacks if the attacker has a way to determine if a ciphertext is correctly padded or not, e.g., through timing leaks or error messages. Padding oracle attacks have found on many protocols, including SSL [16], IPSec [25], and GPG [47].

A3 SSH fragmentation attacks [4] - SSH was designed for securing Internet traffic over the unstable channel, where ciphertext blocks in a packet might get lost. The length of a SSH packet is encrypted in its first block. If the number of delivered blocks is less than the length decrypted from the first ciphertext block, no ciphertext integrity is executed.

If an attacker can inject the first ciphertext block and observe the error message reported by the SSH connection, then the plaintext of the transmitted ciphertext can be recovered.

A4 Partitioning oracle attacks [43] - Some real-world applications do not sample the AEAD symmetric keys randomly but simply pick users' passwords. Thus, attackers might know a set of possible password candidates and perform brute-force attacks. Even worse, if attackers have access to a partitioning oracle, which tells whether the password of a ciphertext belongs to some known sets, then the password can be recovered exponentially faster.

In practice, attackers sometimes can obtain the partitioning oracle by observing the reply messages responding to a selected ciphertext. This causes the vulnerability of applications in the real world, such as Shadowsocks [43].

A5 Salamander attack [26] - The end-to-end secure messaging provides high security against the surveillance of the server but potentially prevents the server from blocking the abusive messages. To mitigate this, Facebook invents a abuse report mechanism that allows each user to report the received abusive messages from a claimed sender.

However, this mechanism turns out to be broken because a malicious sender could send a single encrypted attachment

that would decrypt to both an abusive message and an innocent message under two distinct keys.

A6 Sframe attack [33] - An AEAD scheme authenticates the owners of a symmetric key of a ciphertext rather than the sender's identity. This is especially relevant for group communication, where an AEAD cannot use the shared group key to authenticate the specific sender. To provide sender authenticity in groups, while keeping low bandwidth cost, the IETF SFrame protocol v01 [49] requires senders to sign a portion of the AEAD ciphertext using digital signatures.

Unfortunately, the sender identity authenticity of SFrame mechanism turns out to be broken, since the underlying AEAD schemes, AES-CM-HMAC and AES-GCM, do not provide collision resistance for the unsigned portion. This means, a malicious group member holding the symmetric key can forge the unsigned portion of other group members' ciphertexts.

2.3 Theoretical AEAD frameworks

Apart from the previously outlined classes of real-word attacks linked to AEADs, on the theoretical side, many variants of AEADs have been designed in the past twenty years following the seminal work from [51]. Each of those variants come with their own flavors of properties like, e.g., integrity, confidentiality, nonce-misuse resistance, or robustness, leading to dozens of distinct security definitions. Furthermore, these AEAD variants differ in functionality, with some enabling e.g. ciphertext fragmentation or nonce-hiding. We categorize the main differences between distinct AEAD variants as follows:

F1 does each ciphertext (or a part of it) bind to a set of its encryption inputs? This question motivates the study of a novel (compactly) committing AEAD (ccAEAD [19, Definition 6]) regime as well as various security properties, such as collision resistance [19, Definition 9], commitment [19, Definition 10], sender binding [19, Definition 14], and receiver binding [19, Definition 15] [2, 8, 26, 31]. Roughly speaking, the collision resistance prevents the collisions between AEAD encryption with different inputs. The commitment ensures that each valid AEAD decryption indicates the agreement on a subset of its encryption/decryption inputs. The sender- and receiver binding properties are relevant in the abuse-reporting scenarios. While the sender binding allows every ccAEAD receiver to report abusive messages, the receiver binding prevents malicious receivers from framing honest ccAEAD senders. We give their full definitions in the full version [19] - motivated by A5.

F2 can we find collisions on valid decryption inputs for the same ciphertext? This question motivates the study of a novel property called robustness [1, 43]. Briefly speaking, robustness prevents attackers from having a single ciphertext decrypt to multiple distinct valid messages on different inputs. - motivated by A4, A6.

F3 is the AEAD supporting fragmentation of the ciphertexts? That is, can we start decrypting chunks of data before having

verified the whole ciphertext?

F4 is the decryption atomic, or split into a decryption and an integrity check? [7] - motivated by **A2**.

F5 is the AEAD nonce-hiding? That is, is the nonce explicitly needed for the decryption, or is it included and hidden inside the ciphertexts? [9, 17]

F6 is the AEAD nonce-misuse resistant? [32] Must a nonce be used once strictly, or can repeat? - motivated by **A1**.

3 Generalizing real-world AEAD (in)security for systematic analysis

In this section we develop systematic generalizations of AEAD security and weaknesses, which form the foundation of the symbolic models that we will design in Section 4.

In the previous section we recalled attack classes and security frameworks from the literature. However, these do not exist within a single systematic framework, and basing our symbolic modeling on them would lead to incomparable and ad-hoc models. For our more systematic approach, we first identify the core properties of concrete real life AEAD schemes that lead to the attacks: privacy and integrity, collision-resistance, and nonce-misuse resistance. We show the relevance of these four points by:

- providing in Table 1 the security and weaknesses of many widely deployed AEADs with respect to those core properties;
- illustrating how collision resistance theoretically allows covering notions from **F1** and **F2** in Section 3.2; and
- summarizing the concrete existing collision capabilities for deployed AEADs.

In particular, we identify the collision resistance property as a central concern, which we then investigate first from the theoretical and then the practical point of view.

3.1 Core properties

Stepping back from the many theoretical definitions, we identify three main causes for the protocol attacks:

- **A1** comes from a *misuse of nonces*.
- **A2** and **A3** from a *decryption misuse*, where the decryption is not atomic but performed in two steps, in which case we lose the integrity and privacy guarantees.
- **A4**, **A5**, and **A6** actually all stem from a lack of *collision-resistance*.

This leads us to summarizing the concrete security guarantees for AEADs in three categories:

- **privacy** and **integrity** - the core guarantees that we defined previously, and are expected to be met by all AEADs. This is what is lost under decryption misuse.
- **collision-resistance** - this guarantee hinders attackers from coming up with collisions over the output of Enc, i.e. find two distinct sets of inputs \vec{i}_1 and \vec{i}_2 such that $\text{Enc}(\vec{i}_1) = \text{Enc}(\vec{i}_2)$.

- **nonce-misuse resistance** - this guarantees that using a weak nonce twice or the same nonce for distinct message does not lead to a compromise.

With respect to those core properties, we provide in Table 1 the security and weaknesses of many widely deployed AEADs. In addition to the concrete constructions, we also provide in this table the generic constructions of AEAD such as Encrypt then Mac (EtM), whose security guarantees depend on the concrete encryption and MAC algorithm instantiations. For the generic construction EtM, we distinguish two cases based on whether the encryption and mac keys are related, e.g. derived from k with a key derivation function, or unrelated, e.g. simply the first and the second half of the input k .

Notably, while all of AEADs in the table do provide integrity and privacy (otherwise they would not be used), only some of them tolerate that a single nonce is reused twice for different messages. Moreover, we can also observe that the picture for collision-resistance is very disparate and many deployed schemes do not meet it.

The nonce-misuse, privacy, and integrity properties are now well-understood in the community. In contrast, the collision part is more nascent: there are multiple variants for it on the security notions side, and in practice the concrete weaknesses have not been systematized. In this section, we carry on clarifying the theoretical and practical implications of collision-resistance of AEADs.

3.2 Generalizing AEAD collision resistance and relations

We consider the CMT-4 definition in [8] as a natural definition for full collision resistance and recall it below. Roughly speaking, full collision resistance means that each AEAD ciphertext can only be computed by unique input. While this appears as a strong property, its absence may introduce surprising behaviors for some protocols. Looking forward, this is illustrated by some known attacks or our case-studies, which seem to indicate that despite its strength, full collision resistance is a meaningful and desirable notion.

Definition 4 (Full Collision Resistance). *We say an AEAD = (KGen, Enc, Dec) has ϵ -full collision resistance (or ϵ -full-CR), if the below defined advantage of any attacker \mathcal{A} against the $\text{Expr}_{\text{AEAD}}^{\text{full-CR}}$ experiment in Fig. 4 is bounded by*

$$\text{Adv}_{\text{AEAD}}^{\text{full-CR}} := \Pr[\text{Expr}_{\text{AEAD}}^{\text{full-CR}}(\mathcal{A}) = 1] \leq \epsilon$$

Relationship with existing frameworks It turns out that this notion of collision resistance, while straightforward, is enough to cover in practice multiple notions of the literature from [2, 8, 29, 31, 43]. Informally, these notions are:

- *tidyness* - for a fixed key, is the encryption function the inverse of the decryption one? It implies that collisions over encryptions or decryptions are equivalent.

Concrete AEAD	Integrity and Privacy	Full Collision Resistance	Nonce Misuse Resistance
XSalsa20-Poly1305	•	✗ [2]	✗ Xor of plaintexts
AES-GCM	✓ [34,44]	✗ [26]	✗ Forgeability + xor of plaintexts
ChaCha20-Poly1305	✓ [50]	✗ [2]	✗ Xor of plaintexts
OCB3	✓ [11,41]	✗ [2]	✗ Forgeability + equality of blocks
EtM (unrelated keys)	✓ [51]	✗ [31] ⁴	✗ Encryption dependent
AES-CCM	✓ [30,37]	•	✗ Xor of plaintexts
AES-EAX	✓ [10,46]	•	✗ Xor of plaintexts
EtM (related keys)	✓ [51]	✓ [31]	✗ Encryption dependent
CAU-C4	✓ [8]	✓ [8]	✗ Forgeability + Xor of plaintexts
AES-GCM-SIV	✓ [32,35]	✗ [2]	✓ [32]
CAU-SIV-C4	✓ [8]	✓ [8]	✓ [8]

✓ : proven in the cited work(s). • : we conjecture that this holds, but do not know of a proof.
✗ : does not hold, with reference or explanation of counterexample.

Table 1: AEADs (in)-security guarantees: *Integrity and Privacy* refers to IND \mathcal{S} -CPA and CTI-CPA. *Full Collision Resistance* refers to Definition 4. For *Nonce Misuse Resistance* we indicate the potential impact of reusing nonces if the AEAD scheme does not have this property.

```

ExpAEADfull-CR:
1 ((k1, N1, H1, m1), (k2, N2, H2, m2)) ← $\mathcal{S}$  A()
2 if ⊥ ∈ {k1, N1, H1, m1, k2, N2, H2, m2} or (k1, N1, H1, m1) = (k2, N2, H2, m2)
3   return 0
4 c1 ← Enc(k1, N1, H1, m1), c2 ← Enc(k2, N2, H2, m2)
5 return [c1 = c2]

```

Figure 4: full-CR security for an AEAD = (KGen, Enc, Dec).

- *commitment* (CMT- l and CMTD- l for $l \in \{1, 3, 4\}$ [8]) - can we find collisions either over the encryption or the decryption, with different parts of the inputs being allowed to stay fixed based on l ? In order to capture more variants in this property class that are not included in [8], in this paper we rename CMT- l to *collision resistance* (X-CR) and CMTD- l to *input bound ciphertexts* (X-IBC), where $X \subseteq (k, N, H, m)$ ² denotes the inputs that a AEAD scheme commits to. We recall the definitions of X-CR and X-IBC respectively in [19, Definition 9] and [19, Definition 10] and their relations in [19, Theorem 3]. In particular, the full-CR in Definition 4 is identical to (k, N, H, m) -CR in [19, Definition 9]³.
- *full robustness* (FROB [29]) and *even fuller robustness* (eFROB [31]) - is any attacker able to compute a ciphertext that decrypts correctly under two distinct inputs? This notion was initially defined for randomized AEADs. In this paper, we extend the robustness notions FROB and eFROB for randomized AEADs to a unified notion X-FROB for nonce-based AEADs in [19, Definition 12]

²Here, we slightly abuse notation and use (\cdot) to denote a set. Thus, by $X \subseteq (k, N, H, m)$ we mean that X is a subset of the set (k, N, H, m) . For a single element set, we sometimes also omit the parenthesis and regard it as a single element. For instance, we write $k \in X \Leftrightarrow (k) \subseteq X$.

³In [19, Theorem 4] we will show that (k, N, H, m) -CR implies all variants, which motivated our choice to abbreviate (k, N, H, m) -CR to full-CR.

, where $X \subseteq (k, N, H, m)$ denotes the degree of robustness. Moreover, we prove that X-FROB is equivalent to X-IBC in [19, Theorem 6].

- *key committing* KC security [2] - is any attacker able to compute a ciphertext that decrypts correctly under different keys but same nonce? In this paper, we recall the KC definition in [19, Definition 13] and show that X-FROB with $k \in X$ implies KC, while the reverse does not hold, in [19, Theorem 7].
- *multi-key collision resistance* (MKCR) [43] - is any attacker able to compute a ciphertext that decrypts correctly under multiple keys but same nonce and header? The MKCR is parameterized by the number of distinct keys $\kappa \geq 2$. In this paper, we focus on the simplified case where $\kappa = 2$. We recall the MKCR definition in [19, Definition 13] and show that KC implies the simplified MKCR, while the reverse does not hold, in [19, Theorem 7].
- *receiver binding* (r-BIND) [31] - is any attacker able to compute a ciphertext that can be verified under the different header and message? This notion was initially defined for a variant of compactly committing AEAD (ccAEAD), and showed how it can be instantiated for instance with an Encrypt then Mac construction⁴. Note that [31] also introduces how to transform any AEAD to ccAEAD by a “*traditionally committing encryption*” approach (ccAEAD[AEAD]). In this paper, we recall the r-BIND definition in [19, Definition 15] and show its relations with other security notions (in this list) in [19, Theorem 8] and [19, Theorem 9].

We provide the full relations between the above notions in

⁴ [31] proposes to use HMAC-SHA256 to instantiate a keyed random oracle, which is technically false without an additional assumption, as $HMAC(k, m) = HMAC(H(k), m)$ whenever k is bigger than 256 bits.

the theorem below, which is illustrated in Fig. 5. We give the detailed proofs in the longer version [19]. While some of the relations were conjectured before ([8]), we are the first to provide the full proofs, as well as provide generalizations of some notions to enable a comparison.

Theorem (Informal). *For any AEAD scheme, we have that*

1. X-FROB implies X-CR for any $X \subseteq (k, N, H, m)$. If AEAD is tidy, the reverse also holds. See [19, Theorem 3].
2. X-FROB/X-CR/X-IBC resp. implies X'-FROB/X'-CR/X'-IBC for any $X' \subseteq X \subseteq (k, N, H, m)$. See [19, Theorem 4].
3. X-FROB and X-IBC are equivalent for any $X \subseteq (k, N, H, m)$. See [19, Theorem 5].
4. k -FROB implies KC but not in reverse. See [19, Theorem 6].
5. KC implies MKCR but not in reverse. See [19, Theorem 7].
6. (H, m) -FROB of AEAD implies r-BIND of ccAEAD[AEAD]. r-BIND of ccAEAD[AEAD] implies X-FROB of AEAD for any $X \subseteq (H, m)$. See [19, Theorem 8].
7. Neither KC nor MKCR of AEAD implies r-BIND of ccAEAD[AEAD]. The reverse is same. See [19, Theorem 9].

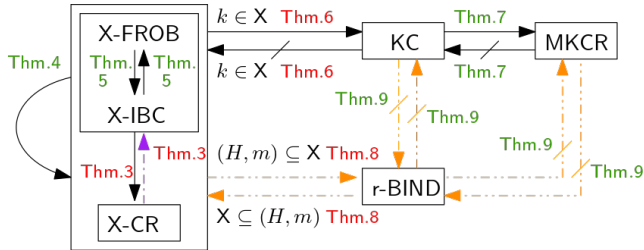


Figure 5: The relation between collision related properties for AEAD with key space Key. The black arrow \rightarrow indicates the general implication. The purple dash-dotted arrow \dashrightarrow indicates the implication for tidy AEAD. The orange dash-dot-dotted arrow \dashrightarrow indicates the implication for ccAEAD[AEAD]. The X in the figure is a subset of (k, N, H, m) , i.e., $X \subseteq (k, N, H, m)$. The theorems highlighted with red color are claimed or proven in other papers. The theorems highlighted with green color are part of our third contribution.

Note that [19, Theorem 3] was proven in [8]. [19, Theorem 6] and [19, Theorem 8] were respectively claimed in [8] and [31] without giving any proofs. Proofs for [19, Theorems 4, 5, 7 and 9] are part of our third contribution. Recall that we have (k, N, H, m) -CR = full-CR in this figure. This theorem indicates that the full collision resistance implies all other existing notions in this figure under the tidyness assumption, which is in fact met by all classical constructions.

3.3 Collision attacks on deployed AEADs

In general, any kind of collision between two ciphertexts can lead to a security issue, and we will advocate that general use AEADs should be fully resistant to collisions. However, many popular deployed AEADs do not meet the full collision resistance, as shown in Table 1. Below, we recall the known attacks against various kinds of collision resistances of different AEAD schemes in the literature.

1. r-BIND: [31] shows a generic attack against any EtM construction with unrelated keys by finding the second key that causes collision by . This attack also applies to real-world modes using Carter-Wegman MACs, e.g., GCM and ChaCha20-Poly1305. [26] shows a concrete attack against AES-GCM and OCB by finding the nonce that causes collision and sketches an faster attack by doing birthday attack on keys. Moreover, at the hand of a corollary of Theorem 1 in [52], [26] claims that this attack also applies to any so-called *rate-1* AEAD, that is, “one blockcipher call per block of message” [26]. This potentially indicates the vulnerability of AES-GCM-SIV and ChaCha20-Poly1305 and any EtM constructions.
2. KC: [2] extends the known attack in [26] against AES-GCM to new proof-of-concept attacks against several commonly used AEAD, including AES-GCM, ChaCha20-Poly1305, AES-GCM-SIV, and OCB3. This attack shows how to create ciphertext collision on two distinct keys. Then, [2] also shows that their new attacks also make impacts in some real-world scenarios, such as the binary polyglots setting.
3. MKCR: [43] shows a novel partitioning oracle attack that feasibly breaks the MKCR security with parameter $\kappa \geq 2$ of widely used AEAD schemes, including AES-GCM, AES-GCM-SIV, ChaCha20-Poly1305, and XSalsa20-Poly1305.
4. X-CR and X-IBC: [8] finds that all above attacks also break the k -CR and -IBC security of respective AEAD schemes. Thus, AES-GCM, AES-GCM-SIV, XSalsa20-Poly1305, and ChaCha20-Poly1305 and OCB are all k -CR insecure, i.e., CMT-1-insecure in [8].

4 Symbolic models for automated verification

We next describe how to, using the generalizations we developed in the previous section, develop symbolic models for AEADs that encompass many of the essential weaknesses from Section 2.2. For each essential weakness, we will develop a specific model that gives the attacker the capability to use the given weakness. Analyzing a protocol with those attacker capabilities enabled will then allow us to automatically find attacks on protocols that may rely on subtle AEAD weaknesses.

Our models cover:

AEAD	nColl	KeysColl
AES-GCM	[26]	[2, 26, 31, 43]
AES-GCM-SIV	[26]	[2, 26, 31, 43]
ChaCha20-Poly1305	[26]	[2, 26, 31, 43]
Encrypt-then-MAC (EtM)	[26]	[26, 31]

Table 2: Effective attacks against collision resistance of several AEADs in the literature. **nColl** describe collisions where, for given keys and a header, the attacker uses brute-force over the nonce to produce colliding ciphertexts. In **KeysColl**, the attacker brute-forces, given a nonce and header, over the keys.

For the generic Encrypt-then-MAC paradigm we give concrete attacks for CTR, OFB, CBC, and CFB modes in the full version [19].

- *collisions Coll*- covering **A4**, **A5**, **A6**, and definitions from **F1** and **F2**.
- *nonce reuse NR*- covering **A1** and **F6**.
- *decryption misuse Forge*- covering **A2**, **A3**, **F3** and **F4**.

Some modern protocols, like [26] or [49], rely on additional features of AEADs that we cover in a modular fashion:

- *explicit tag Tag*- for most AEADs, one can extract a verification tag from the ciphertext, needed to model protocols like [49] for **A6**.
- *explicit commit Com*- to extract a value from the ciphertext committing to the inputs of the encryption. Needed to model protocols like [26] covering **A5** and **F1**.

Collisions can then be lifted to the tag or the commit in a modular fashion, and are essentially only impacting on the complexity of mounting concrete attacks.

We additionally build a model **Leak** that provides an explicit capability to reveal the nonce used for encryption to the attacker. Not sending out the nonce by default but using a dedicated functionality allows accounting for nonce hiding AEADs covering **F5**. While we cannot claim completeness of our models w.r.t. to all possible AEAD weaknesses that may arise in the future, we provide a set of models based on our analysis of the real-world security of AEADs **Section 3** that covers most practical attacks.

We develop and specify the previously enumerated models of AEADs in the *symbolic model* of cryptography, an abstract model used in the formal methods community to express and automate the analysis of cryptographic protocols in **Section 4.1**. We then present symbolic models of the before-mentioned AEAD weaknesses in **Section 4.2**.

4.1 The symbolic model of cryptography

The symbolic model uses function symbols to denote algorithms, and capture their properties through equations. For instance, an encryption is modeled by two binary function symbols enc and dec , with the equation:

$$\text{dec}(\text{enc}(k, m), k) = m$$

Note that the randomness or nonce is not explicit in this classical modeling. And crucially, in the symbolic model, only the equations that are explicitly specified imply equalities. This results in the so-called perfect cryptography assumption: in the previous example, the encryption is perfect, in the sense that given $\text{enc}(k, m)$ and not k , the attacker learns absolutely nothing about m or k as it cannot apply the decryption equation. The attacker cannot change the content of the message, and no collisions will exist.

While the previous assumption may seem too restrictive, it allows for highly automated tools which are one of the strengths of the symbolic model. These tools were already successfully used to automatically find attacks on protocols [23, 60] and aid standardization processes to avoid design-level flaws [21, 24, 48].

In recent years, effort was put into improving the symbolic model with better and more fine-grained support for cryptographic primitives. [22] introduced a stronger version of symbolic Diffie-Hellman group models, while [18] and [36] gave more fine-grained models of cryptographic hash functions and digital signatures, respectively.

4.2 Symbolic AEAD models

We first explicitly model all the input parameters, making the enc and dec having four inputs, $\text{enc}(k, n, h, m)$. Then, we model the multiple weaknesses previously discussed. While we focus on providing models for nonce-based AEADs, as it is the most fine-grained model of AEADs, it is easy to derive from them models for counter-based or randomized AEADs. They can typically be modelled by removing the explicit nonce as $\text{enc}(k, h, m)$, and all equations or capabilities given in the following and that do not directly relate to the nonce can be transposed to this case.

Practical Collision models Coll We start by adding collision capabilities that match the known real-world collision capabilities. When using these models reports an attack on the protocol in one of the automated tools, we can then investigate its feasibility in practice based on the concrete AEAD used and the message encodings by referring to **Section 3**, and in particular **Table 2**.

We start with the capability that can be reasonably computed on many AEADs and was shown to be practical by [26] for Facebook’s Message Franking protocol. As an example, consider the scenario where an attacker tries to produce some colliding ciphertexts given two keys. One option would be to brute-force over the nonce for a fixed header, e.g., an empty header. If successful, the attacker would have a ciphertext that could be decrypted to distinct plaintexts under a common nonce and header using the two keys. We model this nonce finding algorithm in the symbolic model as an additional function symbol c_n , modeling the colliding nonce, which will take as input all the given parameters the collision depends on. We

then add the collision model **nColl**:

$$\begin{aligned} & \text{enc}(k_1, c_n(k_1, k_2, h, m_1, m_2), h, m_1) \\ &= \text{enc}(k_2, c_n(k_1, k_2, h, m_1, m_2), h, m_2) \end{aligned} \quad (\text{nColl})$$

Another widespread collision capability is captured by adding two function symbols c_k^1 and c_k^2 with the collision model **KeysColl**:

$$\begin{aligned} & \text{enc}(c_k^1(n, h, m_1, m_2), n, h, m_1) \\ &= \text{enc}(c_k^2(n, h, m_1, m_2), n, h, m_2) \end{aligned} \quad (\text{KeysColl})$$

To check whether a potential attack found using **KeysColl** might be feasible, refer to [Table 2](#). Whereas for **KeysColl** the attacker needs to produce a collision on the AEAD for a fixed nonce and header, the same kind of collision appears also to be feasible in the case where one of the keys is already fixed. We model this slight variation of **KeysColl** as well (**KeyColl**).

Notice that we, for instance, set that the two colliding encryptions may necessarily use the same nonce and header in those equations. This is caused by many existing protocols implementing that nonces and associated data can be computed independently by the parties, or that they cannot be sent out twice with distinct values.

Generic Collision models FullColl The previous collision models allow to efficiently check for the collisions that are most likely to be practical for existing AEADs and given their use in protocols. While obtaining an attack in those models is instantly interesting, we may miss some future practical attacks. Indeed, as illustrated by [Table 1](#), most AEADs do not meet the full collision resistance property. As we in fact do not know if they meet any kind of collision resistance properties as no proofs exists for many of them, it is possible that in the coming years, new practical ways of building collisions on the existing AEADs are discovered. As such, from a security point of view, for any non collision resistant AEADs, it is prudent to consider that many more collisions are possible than currently practical.

Our approach makes it easy to define such a prudent model, and capture all attacks that may be possible on a protocol if the AEADs is not fully collision resistant. We do this by allowing more collisions and changing which part of the encryption input is fixed on both sides, and which part the attacker is brute-forcing over.

$$\begin{aligned} & \text{enc}(k, n, h, m) \\ &= \text{enc}(\text{gen-c}_k(n, n_2, h, h_2, m, m_2), n_2, h_2, m_2) \end{aligned} \quad (\text{FullKeyColl})$$

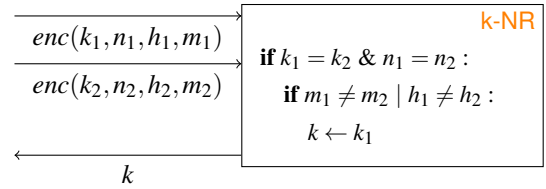
$$\begin{aligned} & \text{enc}(k, n, h, m) \\ &= \text{enc}(k_2, n_2, h_2, \text{gen-c}_m(k, k_2, n, n_2, h, h_2)) \end{aligned} \quad (\text{Full-mColl})$$

$$\begin{aligned} & \text{enc}(k, n, h, m) \\ &= \text{enc}(k_2, \text{gen-c}_n(k, k_2, h, h_2, m, m_2), h_2, m_2) \end{aligned} \quad (\text{Full-nColl})$$

$$\begin{aligned} & \text{enc}(k, n, h, m) \\ &= \text{enc}(k_2, n_2, \text{gen-c}_h(k, k_2, n, n_2, m, m_2), m_2) \end{aligned} \quad (\text{Full-adColl})$$

With **FullKeyColl**, **Full-mColl**, **Full-nColl**, and **Full-adColl**, we capture the capability of an attacker to find collisions by just finding one distinct k , n , h , or m , respectively. These models may cover collisions that are impractical as their main purpose is to check whether the analyzed protocol relies on collision resistance of AEAD schemes. Indeed, if we get an attack in such a model, it implies that a strong collision resistance notion is needed to prove the security of the protocol in the computational model. Further, and as we see later in [Section 5](#), some of these attack may even be practical and could not have been easily discovered in another way.

Nonce-reuse NR The nonce-reuse issue **A1** is slightly more complex to model, as we cannot capture it using an equation. We thus have to use a less classical way to model primitives: we model the attacker capability by providing access to an additional process, or oracle, that does the following:



In this oracle, the attacker can provide two ciphertexts. If those ciphertexts are encrypted under the same key and nonce but differ in either header of message, the attacker learns the secret encryption key. Similar to the collision model, **Coll**, we included a model of this process into our set of AEAD models and call it **k-NR**. This process models the strongest possible leak, namely leakage of the secret key. We can also make it more fine-grained by leaking, e.g., $m_1 \oplus m_2$ instead of k . As not all tools in the symbolic model provide support for exclusive-or like equations, we modeled an over-approximation **m-NR**, which leaks both m_1 and m_2 as an example. Note that with these kinds of oracle-like models, the concrete leaked values can be decided by the capabilities of the chosen tool and the actual weaknesses listed in [Table 1](#).

Decryption Misuses Forge Some protocols, notably SSH, that allow for ciphertext fragmentation, also use AEADs in a non-recommended way by splitting the atomic `dec` operation into verification and decryption. This may also be the case for protocols building their own AEAD based on the EtM construction. In such a case, instead of `dec` that checks integrity, we use a weak decryption function `w-dec` and a verification function `verify`, with the equations:

$$\begin{aligned} & \text{w-dec}(k, n, h, \text{enc}(k, n, h, m)) = m \\ & \text{verify}(\text{enc}(k, n, h, m), k, n, h, m) = \text{true} \end{aligned}$$

We model the fact that the decryption is weak by making decryption succeed on messages forged by the attacker using

the `forge` algorithm:

$$\text{w-dec}(k_2, n_2, h_2, \text{forge}(\text{enc}(k, n, h, m), m_2)) = m_2$$

Remark that a limitation of this **Forge** model is that the attacker cannot compute a valid ciphertext for some function of the message m , which is sometimes possible. Assume that for a given protocol we know that plaintexts are pairs of elements, denoted by $\langle x, y \rangle$, we can also add dedicated forgery rules:

$$\begin{aligned} \text{w-dec}(k_2, n_2, h_2, \text{forge}_1(\text{enc}(k, n, h, \langle m_1, m_2 \rangle))) &= m_1 \\ \text{w-dec}(k_2, n_2, h_2, \text{forge}_2(\text{enc}(k, n, h, \langle m_1, m_2 \rangle))) &= m_2 \end{aligned}$$

If the encryption is XOR based, the attacker should also be able to encrypt at this stage any XOR of a value to the ciphertext. This limitation notably implies that we cannot cover in general premature release of ciphertexts or the SSH fragmentation attacks. While we can do this for particular cases as illustrated, lifting this limitation generically in the symbolic model requires advances in the existing tools and symbolic techniques that we consider out of scope for this work.

Explicit Tag Tag Despite the recommendations, some protocols do not use AEADs only through a decryption and encryption API, but actually rely on some more low-level detail. For instance, schemes rely on the fact that the ciphertext is often a pair (encryption, tag), where the encryption is a basic symmetric encryption of the message and the tag is what provides the integrity. Instead of going to such a low-level, which would be AEAD dependent, we capture this possibility modularly by adding a new function symbol `get_tag`, that inputs ciphertext $\text{enc}(k, n, h, m)$. We can then model collisions over the tags, by adding a variant of each of the previous collision equations over the tag, with, e.g., **nTag** being:

$$\begin{aligned} \text{get_tag}(\text{enc}(k_1, c_n(k_1, k_2, h, m_1, m_2)), h, m_1) \\ = \text{get_tag}(\text{enc}(k_2, c_n(k_1, k_2, h, m_1, m_2)), h, m_2) \end{aligned}$$

Reasoning about explicit tags allows for a top-down approach rather than bottom up. For example, it allows us to ignore implementation details, such as to which side of the encryption the tag is concatenated.

Explicit commitment Com We model compactly committing AEADs by adding a `get-commit` function symbol similar to the `get_tag`. Once again, this allows for a modular model of compactly committing AEADs, where we only specify that a commitment can be extracted, but do not specify how. This extraction can be combined with the collision capabilities to model non-committing AEADs, as a collision on the ciphertexts directly translates to a collision on the commitment. Modeling the fact that the collisions are only on the commitment in a more fine grained way would be possible, but would not yield better attack finding capabilities as they

are covered by the ciphertexts collisions. Further, it appears that collisions on the ciphertext or the commitment only differ in the ease of mounting attacks, the commitment being smaller and easier to manipulate than the whole cipher.

Nonce-Leaking Leak Following **F5**, we capture that an AEAD may not hide the nonce by adding a nonce extraction function symbol `get_nonce` along with the needed equation:

$$\text{get_nonce}(\text{enc}(k, n, h, m)) = n$$

This equation can now be also used instead of sending the nonce to the network explicitly.

Concrete encodings for tools Previously, we presented general equations that can be used to capture multiple AEAD weaknesses in symbolic models. However, symbolic analysis tools often have restrictions on the type of equations that are supported. The most efficient class of equations that are supported by tools satisfy the so called *subterm convergence* property: the right hand side of the equation only contains terms that occur as subterms in the left-hand side. This is not the case for all equations we introduced previously. However, all of them can be expressed in an equivalent fashion using only subterm convergent equations. For instance, for the **nColl** equation the same attacker capability can be captured with the following equations, where we introduce a new function symbol `ct` as an encoding artifact:

$$\begin{aligned} \text{enc}(k_1, c_n(\text{ct}(k_1, k_2, h, m_1, m_2)), h, m_1) \\ = \text{ct}(k_1, k_2, h, m_1, m_2) \\ \text{enc}(k_2, c_n(\text{ct}(k_1, k_2, h, m_1, m_2)), h, m_2) \\ = \text{ct}(k_1, k_2, h, m_1, m_2) \end{aligned} \quad (\text{n-subtermColl})$$

4.3 Automated analysis methodology

We now have a set of models to capture multiple weaknesses of AEADs. To analyze a protocol, the following steps should be followed with the symbolic tool of choice:

1. Verify the protocol in all possible threat models (malicious participants, AEADs weaknesses)
2. If there is an attack based on collisions or nonce misuse, check which AEAD the protocol is using and whether it has the corresponding weakness (**Table 1**);
3. If an attack is from **KeyColl**, **KeysColl**, or **nColl**, use **Section 3.3** to check whether the use AEAD is non collision resistant. If it is not, check **Table 2** to evaluate if the attack is practical or not.
4. If an attack is from one of the over-approximated capabilities **FullKeyColl**, **Full-mColl**, **Full-nColl**, **Full-adColl**, there are two consequences:
 - Collision Resistance is probably needed to prove the protocol computationally secure.

- The attack may however be impractical, and one needs to check the trace to see if the attacker can have enough control over the ciphertext inputs to create a collision.

False attacks In the previously outlined methodology, we say that an attack found in a symbolic model may not be a true attack. To provide more details: during modeling, we sometimes on purpose overapproximate the possible AEAD weaknesses, both to completely rule out classes of attacks or detect subtle attacks. This may indeed lead to us finding “false attacks” that are possible on the design of the protocol but not on all possible implementations for all concrete primitives. In such cases, when by following the methodology we get an attack trace, we have to carefully inspect it, and identify precisely whether it could lead to an attack on the concrete protocol. This typically depends on details of the plaintext encoding used by the protocol (json, concatenation, lengths, etc).

The TAMARIN prover Our methodology is generalized enough to not be bound to a specific tool. The tool of choice needs to support custom equational theories and explicit means to express attacker knowledge. These are criteria fulfilled by various state-of-the-art symbolic tools like [14, 27, 39]. We chose the TAMARIN prover [45] to demonstrate our methodology, as it offers a straightforward way to add custom equational theories and oracle-like processes needed for NR.

Automated analysis setup We split the models from Section 4 into two general classes:

1. collisions and nonce misuse (Coll, NR, Leak)
2. explicit functionalities (Forge, Com, Tag)

Class 1) corresponds to a set of weaknesses that we can check on any protocol using an AEAD scheme. We build a library of those models and a script that verifies the security of a given protocol against those models. Class 2) only makes sense on protocols that do rely on some explicit functionality, like an explicit commitment. Hence, we only model them in the relevant cases where the protocol relies on these explicit functionalities. For our set of case studies, we want to explore their security guarantees against all our AEAD models. To do so, we implemented a Python script that for a given protocol, automatically executes TAMARIN for all possible combinations of threat models, and provides a summary of the secure or insecure scenarios.

When doing this exhaustively, it would mean running TAMARIN 2^{10} times for each case study of class 1) and up to 2^{19} times for class 2). We optimize the script by re-using strict implications of some of our models, e.g. FullKeyColl makes the attacker strictly more powerful than KeyColl. Additionally, some models can be restricted to not be used at the same time

with others. This reduces the possible model combinations to 2^9 (and up to 2^{13} for class 2).) As this is still a huge number of calls, we can use the same implications mentioned before to do some dynamic pruning. The number of prunable model combinations can vary a lot depending on the case study. The total TAMARIN calls that our script automatically made for our case studies can be found in Table 3. Using these implications is useful, both for dynamic pruning and to optimize the search, but also to provide a more compact view of the final results, only displaying non-redundant secure or insecure scenarios. Our script provides us with a summary of the security of a scheme, that can then be formatted in a table as illustrated in Table 4.

A limiting factor in our analysis is the run-time of the protocol models. As the problem of automatically analyzing protocol models is in general undecidable, running TAMARIN could lead to non-termination. We deal with this possibility by introducing timeouts into our experiments such that for each TAMARIN call, we either find a proof, a potential attack trace, or we have a timeout.

Using our technique, which automatically modifies the model for each of the AEAD model combinations, can lead to non-termination more easily, especially on fragile models that were manually tailored toward termination. We selected the value of the timeout depending on the run-time of the protocol model with the classic AEAD model in use.

As an exhaustive search might not be feasible for every future case study, in Appendix A we describe how one can correctly choose the right AEAD model for a certain protocol model.

5 Case studies

We demonstrate our symbolic models for automated verification on a set of eight protocols, classified into four categories depending on the analyzed security property:

- Key Secrecy - rediscovering the attack on *YubiHSM* [61]
- Authentication - rediscovering the attack on *SFrame* [49]
- Accountability - rediscovering the attacks on the accountability of the Facebook Message Franking mechanism [28] and finding that the Web Push [56] standard does not provide server accountability.
- Content Agreement - analysis of multiple group messaging and content delivery protocols, namely SaltPack [53], WhatsApp Groups [59], Scuttlebutt [54], and GPG [38].

We tested our methodology on a computing cluster with Intel® Xeon® Gold 6244 CPUs and 1TB RAM against all possible combinations of the threat models from Section 4. We automate this process using a Python program as described in Section 4.3. For each TAMARIN call within our script we limit TAMARIN to use 4 threads and set the timeout to 60 seconds per TAMARIN call.

For the 8 case studies (plus 3 variants) we had a total evaluation time of 17 hours and 29 minutes with a total of 1404

Protocol	AEAD Scheme	Model	Analysis Results	Time (s)	Novel?	Status
YubiHSM [61]	AES-CCM	NR	Key secrecy attack	2	[42]	Fixed
SFrame [49]	AES-GCM, EtM CTR	Tag	Authentication attack	<1	[33]	Fixed
FB Message Franking [28]	AES-GCM	Coll	Content Agreement attack	8	[26]	Fixed
FB Message Franking [28]	AES-GCM	Coll	Framing attack	3	[26, 31]	Fixed
GPG SED [38]	PGP-CFB	Coll	No Content Agreement	<1	✓	Deprecated
GPG SEIPDv2 [38]	AES-OCB	Coll	No Content Agreement	<1	✓	Infeasible
Saltpack [53]	XSalsa20-Poly1305	Coll	No Content Agreement	8	✓	Infeasible
WebPush [56]	AES-GCM	Coll	Server Accountability	8	✓	Reported
WhatsApp [59]	EtM CBC	Coll	No Content Agreement	3	✓	Reported [‡]
Scuttlebutt [54]	XSalsa20-Poly1305	Coll	No Content Agreement	3	✓	Reported*

* = Feasibility depends on the collision resistance of XSalsa20-Poly1305 (not in Table 2.) See discussion in the full version [19].

‡ = Reported to WhatsApp. Feasibility heavily relies on implementation details, which are not open source.

Table 3: Summary of the main analysis results from our case-studies, illustrating the generality of our models by rediscovering previous attacks and finding new subtleties. In each case, we give the threat model, the used AEAD scheme, the analysis result, as well as the time it took TAMARIN to find it. We also give some additional notes on the status of the observation.

Protocol	Threat Model	Content Agreement
GPG SED	Full-mColl \wedge Full-nColl \wedge Full-adColl \wedge Forge \wedge k-NR \wedge m-NR \wedge Leak	✓
	KeyColl	✗
GPG SEIPDv2	FullKeyColl \wedge Full-nColl \wedge Full-adColl \wedge Forge \wedge k-NR \wedge m-NR \wedge Leak	✓
	Full-mColl	✗
Scuttlebutt	Full-adColl \wedge Forge \wedge k-NR \wedge m-NR \wedge Leak	✓
	KeysColl \vee Full-mColl \vee nColl	✗

Table 4: Example of how our methodology can, given a protocol and a security property, automatically establish the minimal requirements on the AEAD guarantees for the property to hold. We achieve this by analyzing all possible AEAD models, here applied to content agreement for GPG SED, GPG SEIPDv2, and Scuttlebutt. For each analyzed protocol, we obtain the strongest combination of AEAD models under which content agreement holds (✓), which directly yields minimal requirements on the AEAD. The weakest combinations of AEAD models under which a potential violation of the target property (here content agreement) is found is marked with (✗).

TAMARIN calls. We give an overview of the results in Table 3. We show an excerpt of the detailed results in Table 4, while all results are reproducible and can be found in GitHub [20].

Because of space limitations, we only highlight three case studies: The *Facebook Message Franking mechanism* [28], the *Web Push API* [56], and the *Whatsapp group messaging encryption* in Sections 5.1 to 5.3. We refer the reader to the full version [19] for details on the remaining case studies and their detailed attack scenarios.

5.1 Facebook Message Franking

In the setting of End-to-End encryption, reporting the abusive behavior of users seems hard to achieve without weakening security guarantees. In 2016, Facebook introduced *Message Franking* [28] to allow reporting of offensive message attachments. The idea is for a recipient of a malicious message attachment to use a cryptographically sound way to prove that

it was sent by a specific sender.

[26] found an attack against Facebook’s message franking mechanism in 2018. The practical attack they demonstrated involved finding a collision on the used AEAD’s ciphertext. As the sender in this scenario was able to choose the cryptographic keys, messages, and the nonce, they showed how to compute two keys k_1 and k_2 , two message attachments (for which one is the malicious one) m_1 and m_2 , and a nonce n , such that the encryption of m_1 under n and k_1 leads to the same ciphertext as the encryption of m_2 under k_2 and n .

After reporting this attack to Facebook, Facebook immediately patched it. That attack demonstrates the practicality and the impact of collision attacks on real-world schemes.

To show that such attacks can be found on the design level by our analysis, we modeled Facebook’s Message Franking mechanism in TAMARIN. In the initial setting, with the attacker being a malicious sender, we could automatically find the reported attack in a few seconds using the *KeysColl* model.

In addition to analyzing the property violated by the initial attack – can a malicious sender avoid detection? – we also studied the converse property – can a malicious receiver create a fake report? The converse property got first reported as a concern by [31].

We thus additionally model a malicious receiver that tries to report an honest user. In this threat model, we, therefore, look at frameability properties. Being able to frame another party can be severe in practice, for instance, by falsely accusing another person of having sent illegal material. After testing our AEAD models against it, we could re-discover a potential attack [31] on the previously property. However, this attack would require finding a collision on the ciphertext for which one key, the nonce, and the ciphertext itself need to be fixed. Unless further weaknesses of AEADs are found, in this particular case over AES-GCM, this attack is, as of now, impractical.

5.2 Web Push

The Web Push protocol provides means for a server to push notifications to clients by depositing an encrypted notification to the push service that will be fetched by the client when they go online. Web Push is standardized at IETF [56], and, for instance, Apple is planning to integrate it into its ecosystem.

Web Push aims to provide confidential push notifications from a server to its users and to ensure certain privacy properties, like the unlinkability of unique identifiers through the push notification content. Given the wide array of possible applications and concrete use cases, we consider it interesting to check whether the server is accountable or not: can a client prove to a third party that it received a particular push notification from a given server? In contexts where push notifications trigger important actions from a user, protecting users from malicious servers that would try to make the user act and then be able to claim never having done so is critical. The importance of this guarantee would depend on the actual deployment and usage of Web Push; we are currently in an ongoing discussion with IETF on this point. To include this case in our threat model, we thus consider a malicious server controlled by the attacker and verify if it is possible to upload one notification that could be interpreted in two different ways, for instance, offensively or benignly.

Our analysis reports that this guarantee does not hold w.r.t. **Full-mColl**: a single notification can be decrypted validly to two different plaintexts, depending on whether we use the current or deprecated public key of the user. As **Full-mColl** is a strong over-approximation an attack first seemed impractical. After manual inspection of the counterexample trace given by TAMARIN, we could see that this theoretical attack carries over to the real world: WebPush relies on AES-GCM, and we can then reuse similar techniques as for Facebook Message Franking attack: concretely, an attacker can brute-force over the salt used to produce a nonce/key pair to encrypt the

message to find a collision over the unauthenticated part of the ciphertext, and then inject at the end of the ciphertext the needed block to create a collision over the tag. The practicality of the attack depends on the encoding of the plaintexts, and the severity depends on whether the server being not accountable is critical given the use case.

5.3 Content Agreement

We focus now on analyzing the design of multiple messaging mechanisms. We study them in the multiple-recipient setting, trying to answer the following question: *Can a dishonest member of the group send a single message that will be read differently by some recipients?* This question leads us to analyze Content Agreement in the following contexts:

- end-to-end encrypted group messaging applications, like WhatsApp or Signal, or
- dedicated encrypted message mechanism, like GPG, Salt-pack, or Scuttlebutt.

Our study reveals that there is a discrepancy between existing guarantees, which we summarize in [Table 5](#).

We only present the WhatsApp group case study next, and refer the reader to the full version [19] for the others.

WhatsApp groups We model the design of WhatsApp’s group messaging. Because the code of WhatsApp is not available, we constructed our model based on the available information provided in its whitepaper [59, p. 10]. While it relies on the Signal X3DH protocol to establish pairwise channels between the members of the groups, sending a message is slightly different:

- The sender generates a so called *sender key* (also part of the Signal library), and sends this key to each participant over the corresponding pairwise channel;
- To send a message, there is then a single encrypted payload which is uploaded to the server.

While content agreement is trivially broken in Signal itself because of the pairwise channels, it could intuitively be expected within the setting of group messaging. It is however not guaranteed, as reported by our model. Our model captures a group of three people, where one of them is the attacker. We then aim to verify that a given message uploaded to the server will yield the same plaintext for all group members. Our automated analysis reports an attack on this property when enabling ciphertext collisions under **KeyColl**.

The group messaging mechanism relies on an AES-CBC encryption which is then signed with an independent key. This is similar to the Encrypt-then-Mac with unrelated keys scenario. It means that the complexity of mounting an attack in practice is equivalent to the complexity of finding meaningful collisions over AES-CBC. We have seen that with the current capabilities, this strongly depends on the concrete encoding of plaintexts, and whether we can find so-called

Protocol	Content Agreement with CR	Content Agreement without CR	Notes
Whatsapp	✓	✗	Practicality depends on plaintext encodings
Scuttlebutt	✓	✗	Practical
GPG SED (to be deprecated)	✓	✗	Practical
GPG SEIPD v1/v2	✓	✓	Only theoretical attacks
SaltPack	✓	✓	Only theoretical attacks
Signal	✗	✗	Pairwise channels, hence no content agreement

Table 5: Content Agreement summary, with and without Collision Resistance (CR): for a set of group messaging applications and multiple recipients message sending mechanism, we summarize whether a given message can yield to different message for multiple users. In this table, we mention that the Signal application does not meet consistency as a side-remark: as Signal uses pairwise channels to send messages in groups, a different message can be sent to each member of the group.

polyplot plaintexts [2]. As WhatsApp is closed source, verifying the practicality of the attack would require to reverse engineer the full message encoding, which we consider out of scope for this paper. However, given the variety of possible message contents (notification, GIFs, media, React, . . .) it is likely that the encoding would be loose enough to carry out the attack.

5.4 Disclosure

Using our methodology, we detected several undesirable behaviors in the protocol design of Scuttlebutt, Web Push, and WhatsApp Groups [54, 56, 59]. While the behaviors are possible on the protocol design level, their implementation-level feasibility depends on low-level encoding choices.

The behaviors we found did not violate the main specified goals of the respective protocols, and hence we did not mark them as “attacks”. Nevertheless, we contacted the developers of the affected protocols and explained our observations, such that they can assess their implementation-level feasibility, with distinct feedback:

- The developers of the WebPush standard acknowledged the issue, and a discussion is ongoing to determine how to best document the possible behaviors in the standard;
- WhatsApp considered this outside their threat model, and noted that using a different AEAD would still allow a variant of the behavior with the same effect; and
- The Scuttlebutt developers did not respond.

6 Limitations

In an ideal world, we would like to (a) cover all possible AEAD definitions and weaknesses, and (b) have the guarantee that if our method reports an attack, the attack is always feasible in practice. Unfortunately this is not the case yet.

In terms of possible AEAD definitions there are subtle differences that we currently do not capture yet. This includes, for example, properties beyond collisions and nonce-reuse,

such as the “s-way committing security property” [8] that generalizes the CMT notions to the multi-user setting. Our models can also be improved with respect to the *Forge* capability, as discussed in Section 4. On the positive side, we define general models and capture for instance collisions that given the current knowledge are not practical, but could become so in the future, e.g., with new developments on AES. While we do not claim to cover all possible AEAD attacks in the future, this allows to future proof protocols.

With respect to practical feasibility of attacks, the fundamental problem is that our analysis method and standard cryptographic analyses in fact consider protocols *designs* and not their implementation details. For example, this includes abstracting away from encoding details, i.e., how values and compound structures are exactly mapped to bitstrings. Yet such details are critical to determine whether certain collision attacks are possible or not. As a consequence, when we find an attack on the protocol design, this should intuitively be interpreted as: there exists an encoding scheme for which the protocol implementation is insecure. We argue that security of a protocol design should avoid depending on its encoding scheme, and if not, specify the requirements explicitly. The problems we found here using our framework are therefore real concerns for the protocol designs. Still, manual inspection of the implementation is still needed to check whether the encodings allow for generation of the required collisions; but our analysis indicates the types of collisions needed.

7 Further Related Work

We provided background and summary of the related works on AEADs in the first two sections. Here we discuss additional works in formal analysis of security protocols.

Improving the symbolic models of primitives to enable automated attack finding has recently been explored for several types of basic primitives, like cryptographic hashes [18], Diffie-Hellman groups [22], or digital signatures [36]. Our work is centered on AEAD, for which no systematic approach

had been attempted yet.

Ad-hoc approaches include a specific form of nonce-reuse in the Tamarin analysis of WPA [23] and the analysis of Yubikey [42]. In a different approach, [40] modeled the fine-grained block based encryption in the tool ProVerif, but this approach did not scale to protocols of the complexity considered here. Overall, our work is the first to systematically explore weaknesses of concrete algorithms or formal definitions for AEADs and provide models amenable to automation.

ProVerif [13] is, besides Tamarin, the other major tool for automated analysis of protocols in the symbolic model. While we developed our models and case-studies in Tamarin, they could also be used in the ProVerif framework.

As we are focused on automated attack finding due to real life weaknesses of AEADs, our work is orthogonal to tools from the computational model [5, 6, 12] that are all focused on proving security, and cannot find attacks. Our automated models can be useful to establish the assumptions on the AEAD before attempting a computational proof.

With respect to our case studies, based on the absence of collision resistance of AEADs, also referred to as robustness or key-commitment, [26] already reported an attack on the Facebook abuse reporting mechanism, which violates accountability. We are the first to report on behaviors in which is content agreement is not satisfied due to AEAD weaknesses. Other undesirable behaviors have arisen due to collisions, which are linked to oracle partitioning attacks [43], where an attacker can obtain a better than brute force advantage against the OPAQUE protocol. While we can model the relevant AEAD weakness in our framework, modeling the violated security property in the symbolic model is left to future work.

8 Conclusions

We developed the first methodology to analyze the impact of detailed AEAD behaviors on the protocols that use them. Our methodology thus enables detecting protocol weaknesses for a given AEAD, or conversely, determining a protocol's AEAD requirements. The case studies indicate that our methodology is effective and efficient in finding potential weaknesses, notably automatically finding attacks that previously could only be found by manual inspection, thus bringing a new class of attacks within the realm of automated detection.

Acknowledgments

This work received funding from the France 2030 program managed by the French National Research Agency under grant agreement No. ANR-22-PECY-0006.

References

- [1] Michel Abdalla, Mihir Bellare, and Gregory Neven. Robust Encryption. Cryptology ePrint Archive, Report 2008/440, 2008. <https://ia.cr/2008/440>.
- [2] Ange Albertini, Thai Duong, Shay Gueron, Stefan Kölbl, Atul Luykx, and Sophie Schmieg. How to Abuse and Fix Authenticated Encryption Without Key Commitment. In *31st USENIX Security Symposium*, 2020.
- [3] Martin R Albrecht, Jean Paul Degabriele, Torben Brandt Hansen, and Kenneth G Paterson. A surfeit of SSH cipher suites. In *ACM Conference on Computer and Communications Security (CCS)*, pages 1480–1491, 2016.
- [4] Martin R Albrecht, Kenneth G Paterson, and Gaven J Watson. Plaintext recovery attacks against SSH. In *IEEE Symposium on Security and Privacy (S&P)*, pages 16–26, 2009.
- [5] David Baelde, Stéphanie Delaune, Charlie Jacomme, Adrien Koutsos, and Solène Moreau. An interactive prover for protocol verification in the computational model. In *IEEE Symposium on Security and Privacy (S&P)*, pages 537–554. IEEE, 2021.
- [6] Gilles Barthe, Benjamin Grégoire, Sylvain Héraud, and Santiago Zanella Béguelin. Computer-aided security proofs for the working cryptographer. In *Annual Cryptology Conference*, pages 71–90. Springer, 2011.
- [7] Guy Barwell, Daniel Page, and Martijn Stam. Rogue Decryption Failures: Reconciling AE Robustness Notions. In *Proceedings of the 15th IMA International Conference on Cryptography and Coding*, page 94–111, 2015.
- [8] Mihir Bellare and Viet Tung Hoang. Efficient Schemes for Committing Authenticated Encryption. Cryptology ePrint Archive, Report 2022/268, 2022. <https://ia.cr/2022/268>.
- [9] Mihir Bellare, Ruth Ng, and Björn Tackmann. Nonces are Noticed: AEAD Revisited. Cryptology ePrint Archive, Report 2019/624, 2019. <https://ia.cr/2019/624>.
- [10] Mihir Bellare, Phillip Rogaway, and David Wagner. The EAX mode of operation. In *International Workshop on Fast Software Encryption*, pages 389–407, 2004.
- [11] Ritam Bhaumik and Mridul Nandi. Improved security for OCB3. In *International Conference on the Theory and Application of Cryptology and Information Security*, pages 638–666, 2017.

- [12] Bruno Blanchet. CryptoVerif: Computationally sound mechanized prover for cryptographic protocols. In *Dagstuhl seminar "Formal Protocol Verification Applied"*, volume 117, page 156, 2007.
- [13] Bruno Blanchet, Vincent Cheval, and Véronique Cortier. Proverif with lemmas, induction, fast subsumption, and much more. In *IEEE Symposium on Security and Privacy (S&P)*, 2022.
- [14] Bruno Blanchet, Ben Smyth, Vincent Cheval, and Marc Sylvestre. ProVerif 2.00: automatic cryptographic protocol verifier, user manual and tutorial, 2018.
- [15] Alexandra Boldyreva, Jean Paul Degabriele, Kenneth G Paterson, and Martijn Stam. Security of symmetric encryption in the presence of ciphertext fragmentation. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pages 682–699, 2012.
- [16] Brice Canvel, Alain Hiltgen, Serge Vaudenay, and Martin Vuagnoux. Password interception in a SSL/TLS channel. In *Annual International Cryptology Conference*, pages 583–599, 2003.
- [17] John Chan and Phillip Rogaway. Anonymous AE. Cryptology ePrint Archive, Report 2019/1033, 2019. <https://ia.cr/2019/1033>.
- [18] 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*, 2023.
- [19] Cas Cremers, Alexander Dax, Charlie Jacomme, and Mang Zhao. Automated Analysis of Protocols that use Authenticated Encryption: Analysing the Impact of the Subtle Differences between AEADs on Protocol Security (Full version). <https://inria.hal.science/hal-04126116>, August 2023.
- [20] Cas Cremers, Alexander Dax, Charlie Jacomme, and Mang Zhao. Tamarin models and analysis scripts to reproduce the results in this paper, 2023. <https://github.com/AutomatedAnalysisOf/AEADProtocols>.
- [21] Cas Cremers, Marko Horvat, Jonathan Hoyland, Sam Scott, and Thyla van der Merwe. A comprehensive symbolic analysis of TLS 1.3. In *ACM Conference on Computer and Communications Security (CCS)*, pages 1773–1788, 2017.
- [22] Cas Cremers and Dennis Jackson. Prime, Order Please! Revisiting Small Subgroup and Invalid Curve Attacks on Protocols using Diffie-Hellman. In *IEEE 32nd Computer Security Foundations Symposium (CSF)*, pages 78–93, 2019.
- [23] Cas Cremers, Benjamin Kiesel, and Niklas Medinger. A Formal Analysis of IEEE 802.11's WPA2: Countering the Kracks Caused by Cracking the Counters. In *29th USENIX Security Symposium*, pages 1–17, 2020.
- [24] Alexander Dax, Robert Künnemann, Sven Tangermann, and Michael Backes. How to wrap it up—a formally verified proposal for the use of authenticated wrapping in PKCS# 11. In *IEEE 32nd Computer Security Foundations Symposium (CSF)*, pages 62–6215, 2019.
- [25] Jean Paul Degabriele and Kenneth G Paterson. On the (in) security of IPsec in MAC-then-encrypt configurations. In *ACM Conference on Computer and Communications Security (CCS)*, pages 493–504, 2010.
- [26] Yevgeniy Dodis, Paul Grubbs, Thomas Ristenpart, and Joanne Woodage. Fast Message Franking: From Invisible Salamanders to Encryptment. Cryptology ePrint Archive, Report 2019/016, 2019. <https://ia.cr/2019/016>.
- [27] Santiago Escobar, Catherine Meadows, and José Meseguer. Maude-NPA: Cryptographic protocol analysis modulo equational properties. In *Foundations of Security Analysis and Design*, pages 1–50. Springer, 2009.
- [28] Facebook - Messenger Secret Conversations Technical Whitepaper. <https://about.fb.com/wp-content/uploads/2016/07/messenger-secret-conversations-technical-whitepaper.pdf>, 2017. accessed: 2022-08-08.
- [29] Pooya Farshim, Claudio Orlandi, and Răzvan Roşie. Security of Symmetric Primitives under Incorrect Usage of Keys. Cryptology ePrint Archive, Report 2017/288, 2017. <https://ia.cr/2017/288>.
- [30] Pierre-Alain Fouque, Gwenaëlle Martinet, Frédéric Valette, and Sébastien Zimmer. On the Security of the CCM Encryption Mode and of a Slight Variant. In *International Conference on Applied Cryptography and Network Security*, pages 411–428, 2008.
- [31] Paul Grubbs, Jiahui Lu, and Thomas Ristenpart. Message Franking via Committing Authenticated Encryption. Cryptology ePrint Archive, Report 2017/664, 2017. <https://ia.cr/2017/664>.
- [32] Shay Gueron and Yehuda Lindell. GCM-SIV: full nonce misuse-resistant authenticated encryption at under one cycle per byte. In *ACM Conference on Computer and Communications Security (CCS)*, pages 109–119, 2015.
- [33] Takanori Isobe, Ryoma Ito, and Kazuhiko Minematsu. Security Analysis of SFrame. In *European Symposium on Research in Computer Security*, pages 127–146, 2021.

- [34] Tetsu Iwata, Keisuke Ohashi, and Kazuhiko Minematsu. Breaking and repairing GCM security proofs. In *Annual Cryptology Conference*, pages 31–49, 2012.
- [35] Tetsu Iwata and Yannick Seurin. Reconsidering the Security Bound of AES-GCM-SIV. *IACR Transactions on Symmetric Cryptology*, pages 240–267, 2017.
- [36] Dennis Jackson, Cas Cremers, Katriel Cohn-Gordon, and Ralf Sasse. Seems Legit: Automated Analysis of Subtle Attacks on Protocols that Use Signatures. In *ACM Conference on Computer and Communications Security (CCS)*, pages 2165–2180, 2019.
- [37] Jakob Jonsson. On the security of CTR+ CBC-MAC. In *International Workshop on Selected Areas in Cryptography*, pages 76–93, 2002.
- [38] Werner Koch, Paul Wouters, Daniel Huigens, and Justus Winter. OpenPGP Message Format. Internet-Draft draft-ietf-openpgp-crypto-refresh-06, Internet Engineering Task Force, June 2022. Work in Progress.
- [39] Steve Kremer and Robert Künnemann. Automated analysis of security protocols with global state. *Journal of Computer Security*, 24(5):583–616, 2016.
- [40] Steve Kremer and Mark D. Ryan. Analysing the Vulnerability of Protocols to Produce Known-pair and Chosen-text Attacks. In Riccardo Focardi and Gianluigi Zavattaro, editors, *Proceedings of the 2nd International Workshop on Security Issues in Coordination Models, Languages and Systems (SecCo’04)*, pages 84–107, May 2005.
- [41] Ted Krovetz and Phillip Rogaway. The software performance of authenticated-encryption modes. In *International Workshop on Fast Software Encryption*, pages 306–327, 2011.
- [42] Robert Künnemann and Graham Steel. YubiSecure? Formal security analysis results for the Yubikey and YubiHSM. In *International Workshop on Security and Trust Management*, pages 257–272. Springer, 2012.
- [43] Julia Len, Paul Grubbs, and Thomas Ristenpart. Partitioning Oracle Attacks. In *30th USENIX Security Symposium*, 2021.
- [44] David A McGrew and John Viega. The security and performance of the Galois/Counter Mode (GCM) of operation. In *International Conference on Cryptology in India*, pages 343–355, 2004.
- [45] Simon Meier, Benedikt Schmidt, Cas Cremers, and David Basin. The TAMARIN prover for the symbolic analysis of security protocols. In *International conference on computer aided verification*, pages 696–701, 2013.
- [46] Kazuhiko Minematsu, Stefan Lucks, and Tetsu Iwata. Improved authenticity bound of EAX, and refinements. In *International Conference on Provable Security*, pages 184–201, 2013.
- [47] Serge Mister and Robert Zuccherato. An attack on CFB mode encryption as used by OpenPGP. In *International Workshop on Selected Areas in Cryptography*, pages 82–94, 2005.
- [48] Karl Norrman, Vaishnavi Sundararajan, and Alessandro Bruni. Formal Analysis of EDHOC Key Establishment for Constrained IoT Devices. *CoRR*, abs/2007.11427, 2020.
- [49] E. Omara, J. Uberti, A. GOUAILLARD, and S. Murillo. Secure Frame (SFrame) v01. <https://datatracker.ietf.org/doc/html/draft-omara-sframe-01>, 2020. accessed: 2022-08-08.
- [50] Gordon Procter. A Security Analysis of the Composition of ChaCha20 and Poly1305. In *Cryptology ePrint Archive, Paper 2014/613*, 2014. <https://eprint.iacr.org/2014/613>.
- [51] Phillip Rogaway. Authenticated-Encryption with Associated-Data. In *ACM Conference on Computer and Communications Security (CCS)*, CCS ’02, page 98–107, 2002.
- [52] Phillip Rogaway and John Steinberger. Security/Efficiency Tradeoffs for Permutation-Based Hashing. In *EUROCRYPT 2008*, pages 220–236, 2008.
- [53] Saltpack v2. <https://saltpack.org/encryption-format-v2>, 2017. accessed: 2022-08-08.
- [54] Scuttlebot Private Box v0.3.1. <https://scuttlebot.io/more/protocols/private-box.html>, 2019. accessed: 2022-08-08.
- [55] Alon Shakevsky, Eyal Ronen, and Avishai Wool. Trust Dies in Darkness: Shedding Light on Samsung’s Trust-Zone Keymaster Design. In *Cryptology ePrint Archive, Paper 2022/208*, 2022. <https://eprint.iacr.org/2022/208>.
- [56] Martin Thomson. Message Encryption for Web Push. RFC 8291, November 2017.
- [57] Mathy Vanhoef and Frank Piessens. Key reinstallation attacks: Forcing nonce reuse in WPA2. In *ACM Conference on Computer and Communications Security (CCS)*, pages 1313–1328, 2017.
- [58] Serge Vaudenay. Security flaws induced by CBC padding—applications to SSL, IPSEC, WTLS... In *International Conference on the Theory and Applications of Cryptographic Techniques*, pages 534–545, 2002.

- [59] WhatsApp Security Whitepaper. <https://www.whatsapp.com/security/WhatsApp-Security-Whitepaper.pdf>, 2021. accessed: 2022-08-08.
- [60] Jianliang Wu, Ruoyu Wu, Dongyan Xu, Dave Jing Tian, and Antonio Bianchi. Formal Model-Driven Discovery of Bluetooth Protocol Design Vulnerabilities. In *IEEE Symposium on Security and Privacy (S&P)*, pages 2285–2303, 2022.
- [61] YubiHSM. <https://www.yubico.com/resource/hsm-security-for-manufacturing/>, 2021. accessed: 2022-08-08.
- [62] Zerologon – hacking Windows servers with a bunch of zeros . <https://nakedsecurity.sophos.com/2020/09/17/zerologon-hacking-windows-servers-with-a-bunch-of-zeros/>, 2020. accessed: 2023-02-01.

With this, we could see if the protocol relies on collision resistant AEADs.

- As the described protocol explicitly uses AES-GCM tags we would also add the **Tag** models. As collisions on tags are as hard or even easier than finding collisions on the AEAD scheme itself, we would recommend to use at least the same kind of collision types for tags as well, for instance **FullKeyTag**.

A Choosing the correct AEAD model

Whereas using the fully automated methodology from the previous section covers all AEAD models, it can be out-of-scope for complex and detailed protocol models. As complex protocol models often need manual work to aid automation, it might be more feasible to a priori choose the correct AEAD model for the instantiations actually used in the protocol. We demonstrate a way to choose the right combinations of AEAD models on the example of a toy protocol using AES-GCM. Assume that the protocol explicitly adds the functionality that compares the tag instead of using authenticated decryption of the ciphertext:

- As a first step check whether your protocol specification forbids sending the nonce used for AES-GCM. If no, add **Leak** to you AEAD model combination.
- Check **Table 1** and see if the the AEAD is resistant to nonce-reuse attacks. For AES-GCM we see that an XOR of plaintexts can be leaked and there is the possibility to forge ciphertexts. Here, add **k-NR** to the AEAD models. As this is an over-approximation of the before-mentioned weakness, you can also decide to instead of leaking the encryption key, to leak the XOR of plaintext (if your tool of choice allows modeling of XOR) or to output a forged ciphertext under the given key.
- When checking **Table 1** again, AES-GCM is not collision resistant. Then we check **Table 2** and see that AES-GCM is also vulnerable to collisions of type **KeyColl** (**KeyColl**), and **nColl**. As **KeyColl** is strictly stronger than **KeyColl**, we only need to add **KeyColl** and **nColl** to the set of combinations. However, if we would like to future proof the protocol (and we know that AES-GCM is not collision-resistant) we could also decide to add the strongest collision models, e.g. **FullKeyColl**, instead.