A Deniability Analysis of Signal's Initial Handshake PQXDH^{*}

Rune Fiedler Christian Janson

Cryptoplexity, Technische Universität Darmstadt, Germany {firstname.lastname}@cryptoplexity.de

Abstract. Many use messaging apps such as Signal to exercise their right to private communication. To cope with the advent of quantum computing, Signal employs a new initial handshake protocol called PQXDH for post-quantum confidentiality, yet keeps guarantees of authenticity and deniability classical. Compared to its predecessor X3DH, PQXDH includes a KEM encapsulation and a signature on the ephemeral key. In this work we show that PQXDH does not meet the same deniability guarantees as X3DH due to the signature on the ephemeral key. Our analysis relies on plaintext awareness of the KEM, which Signal's implementation of PQXDH does not provide. As for X3DH, both parties (initiator and responder) obtain different deniability guarantees due to the asymmetry of the protocol.

For our analysis of PQXDH, we introduce a new model for deniability of key exchange that allows a more fine-grained analysis. Our deniability model picks up on the ideas of prior work and facilitates new combinations of deniability notions, such as deniability against malicious adversaries in the big brother model, i.e. where the distinguisher knows all secret keys. Our model may be of independent interest.

1 Introduction

Nowadays, messaging applications are widely adopted: Billions of people use messaging apps such as Signal and WhatsApp.¹ In order to provide confidentiality of the users' communication, both apps offer end-to-end encrypted messaging by using the Signal protocol. The Signal protocol consists of an initial handshake, which allows two parties to derive a shared session key, and a ratcheting mechanism [MP16a], which provides forward secrecy and post-compromise security. Until 2023, Signal used the X3DH protocol [MP16b] for the initial handshake. To resist the threat of "harvest now decrypt later attacks" [KS23] posed by future quantum computers, Signal replaced X3DH with the PQXDH protocol; pursuing post-quantum confidentiality (without addressing quantum attackers against authenticity or deniability). Similarly, this paper does not consider deniability against quantum adversaries.

Deniability is a subtle privacy property that allows a party to plausibly deny an action. Consider the following example in a scenario without computers and smartphones present, where Alice and Bob have a conversation. Later, Alice tells Charlie she talked to Bob, and Bob tells Charlie he never talked to Alice. Charlie does not know whose word to trust; hence, Bob can *deny* his conversation with Alice since their conversation did not leave any evidence behind. For sensitive topics such as medical conditions, sexuality, or coordinating a protest against an oppressive regime, this type of privacy guarantee can be not just helpful but even vital.

Conversations in the digital world entail sending data over the internet. Moving our previous example to conversing via messaging, Alice and Bob exchange messages as specified by the messaging protocol,

 $^{^*}$ An extended abstract of this paper appears in the Proceedings on Privacy Enhancing Technologies (PoPETS) 2024. This is the full version.

¹https://en.wikipedia.org/wiki/WhatsApp#User_statistics

resulting in a *protocol transcript*. This transcript may convince Charlie that Alice and Bob *did* converse and *what* they conversed about. However, if Alice and Bob use a deniable protocol, then either party can deny involvement in a protocol run.

A deniable key exchange is necessary (but not sufficient, as pointed out in [CCHD23]) to build deniable messaging. Signal desires deniability² and, therefore, a deniable initial handshake, i.e. previously X3DH and now PQXDH [MP16b, KS23]. Before X3DH, Signal used a derivative of Off-the-Record Messaging (OTR) [OTR20], which also employs a deniable key exchange.

For a key exchange protocol to be deniable, the transcript must not prove Bob's involvement in the conversation. In particular, there must exist a Fake algorithm that *simulates* Bob's involvement. Hence, Alice acting as the *adversary* can simulate a transcript on her own by using the Fake algorithm. Charlie acting as *distinguisher* must not be able to tell the simulated transcript apart from a real transcript. Honest parties would likely use the freshly established session key in some subsequent protocol. Hence, the Fake algorithm must simulate the session key indistinguishably as well.

Several deniability definitions were produced by prior works [DGK06, CF11, HKKP22, DFG⁺13, BFG⁺22, BFG⁺21, YZ13, JS08, Jia14, JCL⁺22]; each definition follows the outline in the previous paragraph. Additionally, [DGK06] introduces *auxiliary information* (e.g. valid protocol transcripts anybody could have observed a priori) that all algorithms (adversary, their equivalent of the Fake algorithm, and the distinguisher) can access. The definitions from the literature are incomparable, vary in the precise guarantee they capture and which formalisms they use. In consequence, capturing several deniability notions may require using several formalisms. Instead, we propose a new model for deniable key exchange, which integrates cherry-picked ideas from prior works: Our model is parameterized in the auxiliary info (*Are some protocol transcripts available?*), the power of the adversary (*Does the adversary follow the protocol?*), the power of the Fake algorithm (*How much help do you need to simulate a transcript?*), the power of the distinguisher (*Is the distinguisher a* big brother *who can force everybody to give up their secret keys? Does it learn the auxiliary info?*). Our new model facilitates new combinations of ideas, e.g. combining malicious adversaries [DGK06] with a big brother distinguisher [BFG⁺22], thereby augmenting the design space for deniability.

Starting a conversation in Signal. Observe that the Signal protocol allows starting a conversation with offline users through the asynchronicity of the initial handshake: Bob prepares a so-called pre-key bundle and uploads it to Signal's key server. To initiate a conversation with Bob, Alice retrieves his pre-key bundle from the key server and sends a message to Bob. (See Figure 3 on page 17 for an algorithmic description of X3DH and PQXDH.)

The X3DH protocol uses the Diffie–Hellman (DH) key exchange, where each user has DH key pairs of different lifetimes (long-term, semi-static, and ephemeral). Bob's pre-key bundle contains his public keys (of all three lifetimes) and a signature certifying his semi-static key under the long-term key. Upon retrieving Bob's pre-key bundle, Alice verifies the signature and generates a fresh ephemeral DH key pair. She computes four (or in case Bob's ephemeral keys run out, only three) DH shared secrets by combining DH keys of different lifetimes, uses the resulting session key to encrypt her first user message with an AEAD scheme, and sends her ephemeral public key and the AEAD ciphertext to Bob. Bob computes the session key in the same manner, decrypts the AEAD ciphertext, and aborts if decryption fails.

The new PQXDH protocol does the same operations as X3DH and adds a KEM encapsulation to achieve post-quantum confidentiality: Bob includes an ephemeral KEM public key in his pre-key bundle and a signature for his KEM key (again under Bob's long-term key). If the ephemeral KEM keys run out, a semi-static KEM key and the corresponding signature are instead included in the pre-key bundle. Alice verifies the signatures on the semi-static DH key and the KEM key and encapsulates against Bob's KEM

²https://signal.org/blog/simplifying-otr-deniability/

key. She derives the session key from the DH shared secrets and the KEM shared secret, and encrypts her first user message with the AEAD scheme under the session key. She sends the resulting (KEM and AEAD) ciphertexts and the public key of her freshly generated ephemeral DH key pair to Bob. Bob decapsulates the KEM ciphertext. Bob computes the session key in the same manner, decrypts the AEAD ciphertext, and aborts if decryption fails.

In order to show deniability, we need to specify a Fake algorithm that simulates protocol messages and the session key indistinguishably. To simulate Bob's message, the Fake algorithm needs to produce signatures under Bob's long-term key without Bob's long-term secret key. For signatures on a semi-static key, the Fake algorithm can reuse any signature from the auxiliary info (in practice, this can be a transcript that the attacker learned from eavesdropping or a pre-key bundle) or from an independent session with Bob. For PQXDH, the Fake algorithm needs a signature on an ephemeral key, i.e. a signature that the distinguisher is not aware of. In Section 4 we call these *private signatures*.

To simulate the session key against malicious adversaries, the Fake algorithm requires all DH shared secrets – and for PQXDH additionally the KEM shared secret. In particular, the Fake algorithm does not know the DH secret keys (some are known to the honest party and some to the adversary); following [VGIK20], we extract the DH shared secrets from the adversary under a knowledge of DH assumption (see Section 2.3 for details; assume if a party produces a DH public key, this party can produce a DH shared secret with this public key, or nobody can). On Alice's side, the Fake algorithm learns the KEM shared secret by encapsulating; on Bob's side the corresponding KEM secret key is not available for decapsulation. Hence, the Fake algorithm extracts the KEM shared secret from the adversary under the plaintext awareness assumption (see Definition 2.7 for details; assume anybody who produces an encapsulation knows the encapsulated secret).

1.1 Related Work

Built on the literature for deniable authentication [DNS98, Kat03, Pas03, DGK05, DG05], Di Raimondo, Gennaro, and Krawczyk [DGK06] have initiated the formal study of deniable key exchange: A distinguisher cannot conclude whether some evidence stems from an actual protocol execution between two parties or whether it was produced by some other means. Many papers have considered deniability for key exchange [Kra96, HKK⁺02, MP02, BMP03, BMP04, LLPM07, YZ10, JS08, DKSW09, YZ13, Jia14, BMS20, JCL⁺22]. Deniability was also named *repudiability*, mostly in the context of Off-the-Record (OTR) messaging [BGB04, UDB⁺15, UG15, UG18]. Next, we revisit some more prominent work on deniability notions in general and their application to Signal's initial handshake.

1.1.1 (Non-)Interactive distinguisher

Deniable authentication [DNS98, DG05] and some forms of deniable key exchange [DGK06] are defined with respect to a distinguisher³ that is presented with evidence after the fact, i.e. an *offline distinguisher*. Dodis, Katz, Smith, and Walfish [DKSW09] have introduced the notion of an interactive, *online distinguisher*. Unger and Goldberg [UG15] argue that online deniability is not achievable for an asynchronous key exchange protocol (i.e. one party uploads pre-keys to a central server) with forward secrecy.⁴ Hence, in this paper we write *deniability* and refer to offline deniability, unless explicitly stated otherwise.

1.1.2 Deniable Key Exchange

Di Raimondo et al. [DGK06] have introduced *concurrent deniability* and *partial deniability* of key exchange protocols, following the simulation-based approach of deniable authentication. They further share the

³Prior work uses the term *judge*, which we avoid for ambiguity with the legal profession.

⁴They call this assumed impossibility "Iron Triangle".

observation of Pass [Pas03] that a simulator for deniability needs to be implementable; in particular, the simulator needs to refrain from rewinding the adversary, reprogramming a random oracle, and similar techniques for simulation in thought experiments. Their deniability definition was since used in [VGIK20] and slightly adapted in [CF11, HKKP22]. Dagdelen, Fischlin, Gagliardoni, Marson, Mittelbach, and Onete [DFG⁺13] have introduced the first game-based definition for a weak version of deniability of key exchange, named *outsider deniability*. Dodis et al. [DKSW09] defined online deniability for key exchange in the Generalized UC framework, which was further refined by Unger and Goldberg [UG15, UG18]. Several works [CF11, JS08, Jia14, JCL⁺22] consider (adaptive) corruptions in the sense that the attacker, the simulator, and possibly the distinguisher learn the corrupted party's secrets. Appendix A discusses the relation between notions of prior work and our new model in more detail.

1.1.3 Signal's initial handshake

Until recently, Signal employed X3DH [MP16b] as initial handshake protocol, whose security can be based on the GapDH assumption [CCD⁺17]. Hashimoto, Katsumata, Kwiatowski, and Prest [HKKP22] have proposed a generic post-quantum construction for Signal's initial handshake based on KEMs for confidentiality and either signatures (named SC-AKE), or ring signatures (SC-DAKE), or ring signatures and NIZKs (SC-DAKE') for authentication and increasing levels of deniability. In concurrent work Brendel, Fiedler, Günther, Janson, and Stebila [BFG⁺22] have proposed a similar generic construction named SPQR based on KEMs and designated verifier signatures (DVS). Dobson and Galbraith [DG22] adopt the X3DH handshake to isogenies, resulting in a protocol named SI-X3DH, which was broken by the SIDH attack [CD23, MMP⁺23, Rob23]. Collins, Huguenin-Dumittan, Nguyen, Rolin, and Vaudenay [CHDN⁺24] have proposed K-Waay, which is based on a primitive named split-KEMs that was introduced by [BFG⁺20]. In 2023, Signal deployed PQXDH [KS23], which uses a more conservative hybrid approach: It combines the classically secure X3DH handshake with a post-quantum secure KEM. This modification aims to add confidentiality against future quantum attackers but is not concerned with quantum attacks against authentication or deniability. Bhargavan, Jacomme, Kiefer, and Schmidt [BJKS23, BJK23] formally verified PQXDH. Fiedler and Günther [FG24] provide a reductionist key indistinguishability analysis of PQXDH, adding coverage for maximum exposure security to the analysis of [BJKS23, BJK23].

1.1.4 Deniability of Signal's initial handshake

Vatandas, Gennaro, Ithurburn, and Krawczyk [VGIK20] have shown that Signal's initial handshake X3DH is deniable in the model of [DGK06] under a novel knowledge of DH assumption. Dobson and Galbraith [DG22] sketch a similar deniability argument for SI-X3DH. Hashimoto et al. [HKKP22] show deniability of SC-DAKE and SC-DAKE' under the model of [DGK06] against semi-honest adversaries and malicious adversaries, respectively, assuming plaintext awareness of the KEM and a related knowledge assumption. Brendel et al. [BFG⁺22, BFG⁺21] show that SPQR is deniable under a new deniability notion (1-out-of-2 deniability against semi-honest adversaries in the big brother model), which they model after the intuition of Signal's specification [MP16b]. Collins et al. [CHDN⁺24] show deniability of K-Waay under a new deniability notion close to the one of Brendel et al. [BFG⁺22, BFG⁺21].

1.1.5 Deniable Secure Messaging

A protocol combining a key exchange with a subsequent exchange of user messages is called *Secure Messaging*. Unger, Dechang, Bonneau, Fahl, Perl, Goldberg, and Smith [UDB⁺15] cover several forms of deniability for Secure Messaging; among them the distinction between *participation deniability* (a user can deny participation in a specific session) and *message deniability* (a user can deny having written a particular (user) message). Cremers and Zhao [CZ24] extended the game-based deniability model of [BFG⁺22]

for key exchange to cover (extended) secure messaging. Reitinger, Malkin, Akgul, Mazurek, and Miers [RMA⁺23], as well as Yadav, Gosain, and Seamons [YGS23] study the social implications of cryptographic deniability for secure messaging.

1.2 Contributions

To capture the precise deniability guarantees of PQXDH, we introduce a new, more expressive deniability model for key exchange protocols (cp. Section 3) that facilitates combinations of previous deniability notions. Our new, game-based deniability model is designed for, but not limited to, asynchronous key exchange. Our model is parameterized to cover deniability notions for a multitude of scenarios: circumstantial knowledge (auxiliary info aux), several adversarial capabilities (malicious, semi-honest, and combinations), how much help the Fake algorithm needs to simulate (partial deniability, or only the peer can simulate, i.e. 1-out-of-2 deniability), and distinguisher capabilities (big brother model, knowledge of auxiliary info). The flexibility of our model allows for easy comparisons of deniability notions and easy adaptability of a proof of deniability to another notion (simply by changing the four aforementioned parameters while keeping the structure of the game, and, hence, the proof).

Furthermore, we provide the first deniability analysis of Signal's new initial handshake PQXDH (cp. Section 4). We show that—due to the signed ephemeral key on Bob's side—the deniability guarantees of PQXDH fall short of those of X3DH. Assuming plaintext awareness of the KEM, which PQXDH adds onto X3DH, and knowledge of a signature hidden from the distinguisher, PQXDH retains the deniability guarantees of X3DH. Table 2 gives the full results.

2 Preliminaries

In this section we review necessary preliminaries.

2.1 Notation

We denote the empty list by [], append *element* to *list* by *list* $\stackrel{+}{\leftarrow}$ *element*, and append the content of *list* to *list* 1 by *list* 1 $\stackrel{+}{\leftarrow}$ *list* 2. In abuse of notation, we allow combining appends, e.g. $(l1, l2) \stackrel{+}{\leftarrow} (l3, l4)$ returns $((l1 \stackrel{+}{\leftarrow} l3), (l2 \stackrel{+}{\leftarrow} l4))$, and allow the same syntax for adding to sets. We denote by $a, b, c \in S$ that S contains a, b, c. Given a map a, we refer to its entry at position i with a[i].

Two distributions are said to be $(t_{\mathcal{D}}, \epsilon_{\mathcal{D}})$ -indistinguishable if no algorithm running in time $t_{\mathcal{D}}$ has an advantage better than $\epsilon_{\mathcal{D}}$ in distinguishing the two distributions. We use $y \leftarrow A(x)$ to denote the random output y of algorithm A for input x, where the probability is over A's internal randomness. We use the arrow \leftarrow for any assignment statements.

2.2 Signatures

We briefly review the syntax for digital signature schemes.

Definition 2.1 (Signature schemes). A digital signature scheme is a triple of algorithms SIG = (KGen, Sig, Vf) with associated message space \mathcal{M}_{SIG} , defined as follows:

- $\mathsf{KGen}() \mathfrak{s} \to (\mathsf{pk}, \mathsf{sk})$: This probabilistic algorithm returns a key pair $(\mathsf{pk}, \mathsf{sk})$
- Sign(sk, m) s→ σ: On input a secret key sk and a message m ∈ M_{SIG}, this probabilistic algorithm returns a signature σ;

• Vf(pk, m, σ) $\rightarrow d$: On input a public verification key pk, a message m, and a candidate signature σ , this deterministic algorithm returns a bit $d \in \{0, 1\}$. If d = 1 we say that the signature is valid, otherwise not.

We say that a digital signature scheme SIG is *correct* if, for every $(\mathsf{pk}, \mathsf{sk}) \leftarrow \mathsf{KGen}()$, every $m \in \mathcal{M}_{\mathsf{SIG}}$, and random $\sigma \leftarrow \mathsf{sSign}(\mathsf{sk}, m)$, it holds that $\Pr[\mathsf{Vf}(\mathsf{pk}, m, \sigma) = 1] = 1$.

We consider the security notion unforgeability under chosen message attacks for digital signature schemes.

Definition 2.2 (Unforgeability for digital signature Schemes). A digital signature scheme SIG is (q_{SIG}, t_{SIG}) -unforgeable if, for every adversary \mathcal{A} running in time t_{SIG} and making at most q_{SIG} queries to the SIGN oracle, it holds that $\Pr[\mathcal{G}^{uf}_{SIG}(\mathcal{A}) = 1] \leq \epsilon_{SIG}$, where $\mathcal{G}^{uf}_{SIG}(\mathcal{A})$ is defined as:

 $\begin{array}{ll} \mathcal{G}_{\mathsf{SIG}}^{uf}(\mathcal{A}): \\ 1 \ Q \leftarrow \emptyset \\ 2 \ (\mathsf{pk}, \mathsf{sk}) \leftarrow \$ \ \mathsf{SIG.KGen}() \\ 3 \ (m^*, \sigma^*) \leftarrow \$ \ \mathcal{A}^{SIGN}(\mathsf{pk}) \\ 4 \ \mathbf{return} \ \mathsf{Vf}(\mathsf{pk}, m^*, \sigma^*) \land m^* \not\in Q \end{array} \qquad \begin{array}{ll} \underbrace{SIGN(m):}{5 \ \sigma \leftarrow \$ \ \mathsf{Sign}(\mathsf{sk}, m)} \\ 6 \ Q \leftarrow Q \cup \{m\} \\ 7 \ \mathbf{return} \ \sigma \end{array}$

2.3 Diffie–Hellman Key Exchange

We briefly review the syntax for the Diffie–Hellman key exchange.

Definition 2.3 (Diffie-Hellman Key Exchange). A Diffie-Hellman Key Exchange (DH) scheme *is a tuple of algorithms* (KGen, DH), *defined as follows:*

- KGen() s→ (pk, sk): This probabilistic algorithm returns a key pair (pk, sk)
- DH(pk_A, sk_B) → DH_{AB}: On input a public key pk_A and a secret key sk_B, this deterministic algorithm returns the shared DH secret DH_{AB}.

We say that a DH key exchange DH is *correct* if, for every $(pk_A, sk_A), (pk_B, sk_B) \leftarrow \mathsf{KGen}()$, it holds that $\Pr[\mathsf{DH}(pk_A, sk_B) = \mathsf{DH}(pk_B, sk_A)] = 1$. For notational convenience we allow the arguments to be in arbitrary order. Hence, $\mathsf{DH}(pk_A, sk_B) = \mathsf{DH}(sk_A, pk_B)$.

The following assumption differs slightly from the EKDH assumption in Definition 2.5: Here, \mathcal{A} is first given two DH public keys and outputs a third DH public key. Then, the extractor tries to extract the two DH shared secrets between the output public key and each of the input public keys.

Definition 2.4 (K2DH Assumption [VGIK20]). Let G be a cyclic group with generator g and AuxPrep a sampler for auxiliary inputs. Let \mathcal{A} be any algorithm running in time $t_{\mathcal{A}}$ which runs on input (U, W, aux) where $U, W \in G$, and $\mathsf{aux} \leftarrow \mathsf{sAuxPrep}$, and outputs $Z \in G$; we denote with $Z = \mathcal{A}(U, W, \mathsf{aux}, r)$ the output of running \mathcal{A} on input U, W, aux with coins r.

We say that the $(t_{\mathcal{A}}, t_{\mathcal{E}}, t_{\mathcal{D}}, \epsilon_{\mathcal{D}})$ -Knowledge of 2DH (K2DH) Assumption holds over group G and for auxiliary info AuxPrep, if for every such \mathcal{A} , there exists a companion algorithm $\mathcal{E}_{\mathcal{A}}^{DH}$ (called the extractor for \mathcal{A}) such that: $\mathcal{E}_{\mathcal{A}}^{DH}$ runs on input U, W, aux, r in time $t_{\mathcal{E}}$ and outputs $\hat{Z}_1, \hat{Z}_2 \in G$ or \perp such that

- If $\mathcal{E}_A^{DH}(U, W, \mathsf{aux}, r) \neq \bot$ then $\hat{Z}_1 = DH(U, Z)$ and $\hat{Z}_2 = DH(W, Z)$
- For every algorithm C running in time t_D, we have that Pr[C(U, W, Z, r, aux) ∈ {DH(U, Z), DH(W, Z)} | E_A^{DH}(U, W, aux, r) = ⊥] ≤ ε_D, where Z = A(U, W, aux, r) and the probability is taken over the coins of C and uniform distribution on (U, W, r).

We express whether a party can compute functions of its own DH secret key with the following assumption: Consider an algorithm \mathcal{A} that outputs a DH public key. Then an extractor $\mathcal{E}_{\mathcal{A}}^{DH}$ with runtime $t_{\mathcal{E}}$ given the same inputs as \mathcal{A} and an extra DH public key can extract the DH shared secret of both public keys from \mathcal{A} . If $\mathcal{E}_{\mathcal{A}}^{DH}$ fails to extract the shared DH secret, then no other algorithm with runtime $t_{\mathcal{D}}$ succeeds with probability better than $\epsilon_{\mathcal{D}}$. To put it in a nutshell, extraction succeeds within $t_{\mathcal{E}}$, or not even within $t_{\mathcal{D}}$. We denote the DH shared secret between two DH public keys U, Z with respect to some generator g as DH(U, Z).

Observe that in [VGIK20], Vatandas et al. show that the knowledge of exponent assumption and their novel knowledge of discrete logarithm assumption imply the knowledge of DH (KDH) assumption. Compared to the EKDH assumption, the algorithm \mathcal{A} of the KDH assumption receives the public key U as input.

Definition 2.5 (EKDH Assumption [VGIK20]). Let G be a cyclic group with generator g and AuxPrep a sampler for auxiliary inputs. Let \mathcal{A} be any algorithm running in time $t_{\mathcal{A}}$ which runs on input aux where aux \leftarrow s AuxPrep, and outputs $Z \in G$; we denote with $Z = \mathcal{A}(\mathsf{aux}, r)$ the output of running \mathcal{A} on input aux with coins r.

We say that the $(t_{\mathcal{A}}, t_{\mathcal{E}}, t_{\mathcal{D}}, \epsilon_{\mathcal{D}})$ -Extended Knowledge of DH (EKDH) Assumption holds over group G and for auxiliary info AuxPrep, if for every such \mathcal{A} , there exists a companion algorithm $\mathcal{E}_{\mathcal{A}}^{DH}$ (called the extractor for \mathcal{A}) such that: $\mathcal{E}_{\mathcal{A}}^{DH}$ runs on input aux, r and an additional input $U \in G$ in time $t_{\mathcal{E}}$ and outputs $\hat{Z} \in G$ or \perp such that

- If $\mathcal{E}^{DH}_{\mathcal{A}}(U, \mathsf{aux}, r) \neq \bot$ then $\hat{Z} = DH(U, Z)$
- For every algorithm C running in time $t_{\mathcal{D}}$, we have that $\Pr[C(U, r, \mathsf{aux}) = DH(U, Z) | \mathcal{E}^{DH}_{\mathcal{A}}(U, \mathsf{aux}, r) = \bot] \leq \epsilon_{\mathcal{D}}$, where $Z = \mathcal{A}(\mathsf{aux}, r)$ and the probability is taken over the coins of C and uniform distribution on r.

2.4 Key Encapsulation Mechanisms (KEMs)

We briefly review the syntax for KEMs.

Definition 2.6 (Key Encapsulation Mechanisms). A key encapsulation mechanism (KEM) is a triple of algorithms KEM = (KGen, Enc, Dec) with associated ciphertext space C and key space K. In more detail:

- $\mathsf{KGen}() \mathrel{\mathfrak{s}}{\to} (\mathsf{pk}, \mathsf{sk})$: A probabilistic algorithm that outputs a key pair $(\mathsf{pk}, \mathsf{sk})$
- Enc(pk) s→ (ct, ss): A probabilistic algorithm taking as input a public key pk and outputs a ciphertext ct ∈ C and the therein encapsulated key ss ∈ K.
- Dec(sk, ct) → ss': A deterministic algorithm taking as input a ciphertext ct ∈ C and secret key sk and outputs ss ∈ K ∪ {⊥}, where ⊥ indicates an error.

We define the key collision probability γ_{coll} for a KEM KEM as the probability of the public keys of two independently sampled key pairs to coincide.

We say that $\mathsf{KEM} = (\mathsf{KGen}, \mathsf{Enc}, \mathsf{Dec})$ is δ -correct if, for every key pair $(\mathsf{pk}, \mathsf{sk}) \leftarrow \mathsf{sKGen}()$ and every encapsulation $(ct, ss) \leftarrow \mathsf{sEnc}(\mathsf{pk})$, it holds that $\Pr[ss' \neq ss \mid ss' \leftarrow \mathsf{Dec}(\mathsf{sk}, ct)] \leq \delta$.

Plaintext awareness for KEMs denotes whether one can create a valid KEM encapsulation without knowing the corresponding shared secret. Based on plaintext awareness for PKE [BR95, BP04], and for KEMs [Den06, JW10], Hashimoto et al. [HKKP22] have adapted the definition to plaintext awareness for KEMs in the multi-key setting. We adapt their definition from asymptotic security to concrete security.

```
\mathcal{G}^{dec}_{\mathsf{KEM}}(\mathcal{C}, \mathcal{D}):
                                                                                                                                                                  \mathcal{G}_{\mathsf{KEM}}^{ext}(\mathcal{C}, \mathcal{E}_{\mathcal{C}}^{PA}, \mathcal{D}):
1 for i \in [n_k]
                                                                                                                                                                   8 for i \in [n_k]
           (\mathsf{pk}_i, \mathsf{sk}_i) \leftarrow \mathsf{KEM}.\mathsf{KGen}()
                                                                                                                                                                  9
                                                                                                                                                                         (\mathsf{pk}_i, \mathsf{sk}_i) \leftarrow \mathsf{KEM.KGen}() \quad DECAPS(i, ct):
                                                                             DECAPS(i, ct):
3 \vec{\mathsf{pk}} \leftarrow \{\mathsf{pk}_i\}_{i \in [n_k]}
                                                                                                                                                                10 \vec{\mathsf{pk}} \leftarrow \{\mathsf{pk}_i\}_{i \in [n_k]}
                                                                               7 return KEM.Dec(sk_i, ct)
                                                                                                                                                                                                                                               14 return \mathcal{E}_{\mathcal{C}}^{PA}(i, ct, r)
                                                                                                                                                                11 r \leftarrow \mathcal{R}_{\mathcal{C}}
4 r \leftarrow \mathcal{R}_{\mathcal{C}}
5 v \leftarrow C^{\text{Decaps}}(\vec{\mathsf{pk}}; r)
                                                                                                                                                                 12 v \leftarrow C^{\text{Decaps}}(\vec{\mathsf{pk}}; r)
6 return \mathcal{D}(v)
                                                                                                                                                                13 return \mathcal{D}(v)
```

Figure 1: Games for plaintext awareness, see Definition 2.7.

Definition 2.7 (Plaintext Awareness for KEMs [HKKP22]). A KEM scheme KEM is $(n_k, t_c, t_{\mathcal{E}}, t_{\mathcal{D}}, \epsilon_{\mathcal{D}})$ plaintext aware (PA1-secure) if for all (non-uniform) ciphertext creators \mathcal{C} running in time t_c , there exists an efficient extractor $\mathcal{E}_{\mathcal{C}}^{PA}$ running in time $t_{\mathcal{E}}$ such that for any distinguisher \mathcal{D} running in time $t_{\mathcal{D}}$, the two experiments $\mathcal{G}_{\mathsf{KEM}}^{dec}(\mathcal{C}, \mathcal{D})$ and $\mathcal{G}_{\mathsf{KEM}}^{ext}(\mathcal{C}, \mathcal{E}_{\mathcal{C}}^{PA}, \mathcal{D})$ in Figure 1 are indistinguishable except with a probability of $\epsilon_{\mathcal{D}}$.

Intuitively, a KEM is plaintext aware if creating a KEM ciphertext implies knowledge of the corresponding shared secret. This is formally modeled with a ciphertext creator who has access to a decryption oracle. This decryption oracle is implemented either with the actual decryption or with an extractor that depends on the ciphertext creator and its randomness. Finally, the ciphertext creator outputs a string, which is fed into a distinguisher. The distinguisher needs to distinguish how the decryption oracle is implemented.

Note that Kyber [SAB⁺22] using the FO transform with explicit rejection instead of implicit rejection is plaintext aware.

2.5 Authenticated Encryption with Associated Data (AEAD)

We follow the exposition of [Rog 02].

Definition 2.8 (AEAD). An Authenticated Encryption scheme with Associated Data is a pair of algorithms AEAD = (Enc, Dec) with associated key space \mathcal{K} , nonce space \mathcal{N} , associated data (header) space \mathcal{H} , message space \mathcal{M}_{AE} , and ciphertext space \mathcal{C} , defined as follows:

- Enc(k, n, AD, μ) → ct: On input a key k, a nonce n, associated data AD, and a message μ, this deterministic algorithm returns a ciphertext ct.
- Dec(k, n, AD, ct) → μ: On input a key k, a nonce n, associated data AD, and a ciphertext ct, this deterministic algorithm returns a message μ or a distinguished error symbol ⊥ ∉ M_{AE}.

We say that an AEAD scheme is correct if, for every $k \in \mathcal{K}, n \in \mathcal{N}, AD \in \mathcal{H}, \mu \in \mathcal{M}_{AE}$ it holds that $\Pr[\mathsf{Dec}(k, n, AD, \mathsf{Enc}(k, n, AD, \mu)) = \mu] = 1.$

Note that throughout the paper we omit the nonce since Signal's implementations of X3DH and PQXDH derive the nonce deterministically from the key.

2.6 Key Exchange Protocols

We review the syntax for key exchange protocols, essentially following [BFG⁺22, BFG⁺21]. We differentiate between protocol messages m and user (plaintext) messages μ (such as "Hi Bob, how are you doing?"). A key exchange protocol may allow a user message to be sent along with a protocol message during the execution, e.g. X3DH and PQXDH.

Definition 2.9 (Key Exchange Protocol). A 2-party key exchange protocol is a triple of algorithms KE = (KGenLT, KGenSS, Run), defined as follows:

- KGenLT() s→ (ltpk_U, ltsk_U): A probabilistic long-term key generation algorithm that outputs a publickey/secret-key pair attributed to user U.
- KGenSS() s→ (sspk^{ssid}_U, sssk^{ssid}_U): A probabilistic semi-static key generation algorithm that outputs a public-key/secret-key pair attributed to user U and semi-static identifier ssid.
- Run(ltsk_U, sssk_U, ltpk, sspk, π, m, μ) s→ (π', m', μ'): A probabilistic session execution algorithm that takes as input a party's long-term secret key ltsk_U, a list of that party's semi-static secret keys sssk_U, lists of long-term and semi-static public keys for all honest parties ltpk and sspk, a session state π, an incoming message m, and a (potentially empty) user message μ to be sent, and outputs an updated session state π', a (possibly empty) outgoing message m', and a (possibly empty) incoming user message μ'.

We assume a distinguished message $m = (\text{create}, \text{info}_{\text{create}})$ upon which Run sets up the session and produces the first message. Note that an ephemeral key generation algorithm is not mandatory and can happen inside of Run.

Parties and sessions. Let \mathcal{P} be the set of n_p parties, each of whom has a long-term public-key/secretkey pair generated by an algorithm KGenLT. Each party may run multiple instances of the protocol simultaneously or sequentially, each of which is called a session. The *i*th session at party P is denoted π_P^i . For each session, the party maintains the following collection of session-specific information:

- $\mathsf{oid} \in \mathcal{P}$: The identity of the session owner.
- $pid \in \mathcal{P} \cup \{\star\}$: The identity of the intended peer, which may initially be unknown (indicated by \star).
- role \in {initiator, responder}: The role of the party in this session.
- $K \in \mathcal{K}_{KE} \cup \{\bot\}$: The session key established in this session, initialized to \bot .
- info_{create} is a tuple of elements that indicates the components of the first message. We use the following elements:
 - ssid $\in [n_{ss}]$ denotes the identifier of the initiator's semi-static key in this session. If π .role = initiator this refers to $sspk_{\pi,oid}^{ssid}$, if π .role = responder this refers to $sspk_{\pi,oid}^{ssid}$.
 - $-e_{\mathsf{DH}} \in {\mathsf{true}, \mathsf{false}}$ denotes whether an ephemeral DH key is included in the first message.
 - $-e_{\mathsf{KEM}} \in {\mathsf{true}, \mathsf{false}}$ denotes whether an ephemeral KEM key is included in the first message.

Asynchronous key exchange. In principle, a key exchange protocol can have an arbitrary number of message flows n_m , which correspond to multiple calls to Run for a single session. In normal execution of an *asynchronous* authenticated key exchange protocol, the following three calls to Run occur: 1) a call to Run at the initiator Bob with $m = (\text{create}, \text{info}_{\text{create}})$, which sets up the initiator session and outputs the pre-key bundle, possibly including ephemeral public keys; 2) a call to Run at the responder Alice with the initiator's pre-key bundle, which generates a session key and outputs a key exchange message; and 3) a call to Run at the initiator with the responder's long-term public key and key exchange message, which generates a session key and has no output message. If a party outputs an empty message before the protocol execution is finished, it aborts the protocol.

3 Our deniability model

We introduce our deniability model that facilitates combining ideas of deniability definitions from prior literature, hence gaining expressiveness. We show how our model can match definitions from prior literature in Appendix A.

We want to capture that artifacts of a key exchange between Alice and Bob do not convince a third party that the key exchange actually took place. Consider Alice, who tries to prove Bob's involvement. To defend against Alice's claim, Bob must be able to argue that the artifacts were produced without him, e.g. by Alice or by just anybody. Hence, we require the existence of a Fake algorithm that produces these artifacts, where the artifacts are indistinguishable from artifacts of an actual execution of the key exchange protocol.

Our model captures this with the following game: First, the challenger prepares the game by initializing variables, sampling keys for all users, preparing auxiliary info aux according to AuxPrep, and sampling a secret bit b. Second, the adversary interacts with some oracles. One of the oracles depends on the secret bit b and embeds the challenge: To distinguish whether (a part of) a protocol transcript was generated honestly or with a dedicated Fake algorithm, which needs to be given in the proof. Third, the distinguisher guesses the challenge bit, based on the previous actions of the adversary. In particular, the distinguisher gets the transcript of the interactions between the adversary and the challenge oracle and the (honestly computed) session keys for all sessions. Note that \mathcal{A} and Fake get access to the auxiliary info aux.

The game is parameterized by the adversary's capabilities O_A , i.e. the oracles it can query, the capabilities of the Fake algorithm O_F , the capabilities of the distinguisher O_D , and the sampling method for auxiliary information AuxPrep. Extending O_A or O_D strengthens the deniability guarantee, while extending O_F or AuxPrep weakens the guarantee. In Sections 3.1 to 3.4 we explain these parameters in detail. In Section 3.5 we discuss how to compare definitions with different parameters. The flexibility and expressiveness of these parameters is the novelty of our model.

Definition 3.1 (Deniability for Key Exchange). We say that a key exchange protocol KE = (KGenLT, KGenSS, Run) is $(n_p, n_{ss}, q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D, \epsilon_D)$ -deniable wrt. (O_A, O_F, O_D) -oracles, where $O_A \subseteq \{REG, INIT, CHALLINIT, CHALLRESP, REGHON, CHALLHONINIT, CHALLHONRESP\}$, $O_F \subseteq \{SK, USER_{x,y}\}$, $O_D \subseteq \{SKS, AUX\}$, and auxiliary inputs sampled with AuxPrep, if for any adversary A running in time t_A , making queries to O_A limited by q_{O_A} , there exists an algorithm Fake making queries to O_F limited by q_{O_F} such that for any distinguisher D making queries to O_D limited by q_{O_D} and running in time t_D , it holds that

$$\Pr[\mathcal{G}_{\mathsf{KE},\mathsf{AuxPrep},n_p,n_{ss}}^{O_A,O_F,O_D\text{-}den}(\mathcal{A},\mathcal{D})=1] \leq \frac{1}{2} + \epsilon_{\mathcal{D}},$$

where $\mathcal{G}_{\mathsf{KE},\mathsf{AuxPrep},n_p,n_{ss}}^{O_A,O_F,O_D-den}(\mathcal{A},\mathcal{D})$ is defined in Figure 2.

Intuitively, the Fake algorithm models that a protocol message was not necessarily produced by the alleged sender. Consider an adversary \mathcal{A} acting as Alice (i.e. the adversary generates keys for Alice), who produces Alice's message(s) according to Run and Bob's message(s) and the session key according to the Fake algorithm. Note that the Fake algorithm also needs to produce the session key to ensure deniability of the subsequent protocol. Since the Fake algorithm is public, Bob can argue that the transcript and session key were produced by Alice using the Fake algorithm, without his involvement. On a technical note, the Fake algorithm knows the code and randomness of the adversary \mathcal{A} since they are in cahoots. This allows the Fake algorithm to extract knowledge from the adversary under some knowledge assumptions.

In the following we describe the choices for using our model and how our model allows easy comparisons of deniability notions.

 $\mathcal{G}_{\mathsf{KE},\mathsf{AuxPrep},n_p,n_{ss}}^{\mathsf{O}_A,\mathsf{O}_F,\mathsf{O}_D\text{-}den}(\mathcal{A},\mathcal{D})\text{:}$ CHALLHONINIT $(U, V, info_{create}, \vec{\mu}, r_C)$: 42 if $U \in C$ //do not allow challenging a party under \mathcal{A} 's control $1 \quad Q[\cdot] \leftarrow []; \quad K[\cdot] \leftarrow \bot; \quad C \leftarrow \emptyset$ 43 $\mathbf{return} \perp$ 2 $(\vec{\mathsf{pk}}, \vec{\mathsf{sk}}) \leftarrow \mathsf{KeyPrep}(n_p, n_{ss})$ 44 π_U .oid $\leftarrow U$; π_U .pid $\leftarrow V$ 3 aux \leftarrow \$ AuxPrep(\vec{pk}, \vec{sk}) 45 π_V .oid $\leftarrow V$; π_V .pid $\leftarrow U$ 4 $b \leftarrow \{0, 1\}$ 46 π_U .role \leftarrow initiator; π_V .role \leftarrow responder 5 $r \leftarrow R_A$ //randomness for the adversary 47 $m_1 \leftarrow (\text{create}, \text{info}_{\text{create}})$ 6 $\mathcal{A}^{\mathcal{O}_A}(\vec{\mathsf{pk}},\mathsf{aux};r)$ 48 $(\mu_1,\ldots,\mu_{n_m}) \leftarrow \vec{\mu}$ 7 $b' \leftarrow \mathcal{D}^{O_D}(\vec{\mathsf{pk}}, r, Q, K)$ 49 for $i \in [1, 3, \dots, n_m - 1]$ //until n_m if n_m is odd 8 return $\llbracket b' = b \rrbracket$ 50 if b = 0KeyPrep (n_p, n_{ss}) : $(\pi_{U}, m_{i+1}) \leftarrow \text{\$Run}(\mathsf{sk}_{U}, \overrightarrow{\mathsf{pk}}, \pi_{U}, m_{i}, \mu_{i})$ 51 9 for $U \in [n_p]$ //number of keys to prepare 52 else $(ltpk_{II}, ltsk_{II}) \leftarrow \mathsf{KGenLT}()$ $(\pi_{II}, m_{i+1}) \leftarrow \text{Fake}^{O_F}(\vec{\mathsf{pk}}, \pi_{II}, m_i, \mu_i, \mathsf{aux}, r_C)$ 53 11 for ssid $\in [n_{ss}]$ 54 $(\pi_V, m_{i+2}) \leftarrow \mathsf{Run}(\mathsf{sk}_V, \vec{\mathsf{pk}}, \pi_V, m_{i+1}, \mu_{i+1}; r_C)$ $(sspk_{II}^{ssid}, sssk_{II}^{ssid}) \leftarrow \mathsf{KGenSS}()$ 12 55 $Q[\pi_U] \leftarrow (m_i)_{i=1}^{n_m};$ $K[\pi_U] \leftarrow \pi_U.K$ 13 $lt\vec{p}k \leftarrow \{ltpk_U\}_{U \in [n_p]}; \quad s\vec{spk} \leftarrow \{sspk_U^{ssid}\}_{U \in [n_p]}^{ssid \in [n_{ss}]}$ 56 **return** $(Q[\pi_{II}], K[\pi_{II}])$ 14 $lt\vec{s}k \leftarrow \{ltsk_U\}_{U \in [n_p]}; \quad s\vec{s}k \leftarrow \{sssk_U^{\text{ssid}}\}_{U \in [n_p]}^{\text{ssid} \in [n_{ss}]}$ ChallHonResp $(U, V, \mathsf{info}_{\mathsf{create}}, \vec{\mu}, r_C)$: 57 if $V \in C$ //do not allow challenging a party under \mathcal{A} 's control 15 return $((lt\vec{p}k, s\vec{spk}), (lt\vec{s}k, s\vec{ssk}))$ 58 return | $\operatorname{Reg}(U, \vec{\mathsf{pk}}_{U})$: 59 π_U .oid $\leftarrow U$; π_U .pid $\leftarrow V$ 16 if $U \in [n_p] \cup C$ //abort if another key is already known for U 60 π_V .oid $\leftarrow V$; π_V .pid $\leftarrow U$ 17 return \perp 61 π_{II} .role \leftarrow initiator; π_{V} .role \leftarrow responder 18 $C \leftarrow C \cup \{U\}$ 62 $m_1 \leftarrow (\text{create}, \text{info}_{\text{create}})$ 19 $\vec{pk} \leftarrow \vec{pk}_{II}$ //includes semi-static keys 63 $(\mu_1,\ldots,\mu_{n_m}) \leftarrow \vec{\mu}$ 64 for $i \in [1,3,\ldots,n_m-1]$ //until n_m if n_m is odd 20 return success $(\pi_{U}, m_{i+1}) \leftarrow \mathsf{Run}(\mathsf{sk}_{U}, \vec{\mathsf{pk}}, \pi_{U}, m_{i}, \mu_{i}; r_{C})$ 65 INIT (U, s, role, V) : if b = 066 21 if $\pi_U^s = \bot \land U \notin C$ $(\pi_V, m_{i+2}) \leftarrow \mathbb{R}un(\mathsf{sk}_V, \vec{\mathsf{pk}}, \pi_V, m_{i+1}, \mu_{i+1})$ 67 22 π^s_{II} .role \leftarrow role 68 else $\pi^s_{\scriptscriptstyle II}.{\sf oid} \leftarrow U$ //set owner identity 23 $(\pi_V, m_{i+2}) \leftarrow \text{*} \mathsf{Fake}^{\mathcal{O}_F}(\vec{\mathsf{pk}}, \pi_V, m_{i+1}, \mu_{i+1}, \mathsf{aux}, r_C)$ 69 24 π^s_{II} .pid $\leftarrow V$ //set partner identity (* if post-specified peers) 70 $Q[\pi_V] \leftarrow (m_i)_{i=1}^{n_m};$ $K[\pi_V] \leftarrow \pi_V.K$ 25 return success 71 return $(Q[\pi_V], K[\pi_V])$ 26 return ⊥ USER_{x,y} (π, u, m, μ) : $\operatorname{RegHon}(U)$: 72 $U \leftarrow \pi.oid$ 27 if $U \in [n_p] \cup C$ //abort if another key is already known for U 73 $H \leftarrow H^u_\pi$ //easier notation for $u^{ ext{th}}$ honest dummy for session π return | 28 74 if $\pi_{U}^{H} = \bot$ 29 $(\vec{\mathsf{pk}}_U, \vec{\mathsf{sk}}_U) \leftarrow \text{KeyPrep}(1, n_{ss})$ 75 $(\mathsf{pk}_H, \mathsf{sk}_H) \leftarrow \mathsf{KeyPrep}(1, 1)$ 30 $\vec{\mathsf{pk}} \xleftarrow{+}{\leftarrow} \vec{\mathsf{pk}}_U; \quad \vec{\mathsf{sk}} \xleftarrow{+}{\leftarrow} \vec{\mathsf{sk}}_U$ $\vec{\mathsf{pk}} \leftarrow \vec{\mathsf{pk}} \cup \{\mathsf{pk}_H\}$ 76 31 $C \leftarrow C \cup \{U\}$ if π .role = initiator 77 32 return $(\vec{pk}_{II}, \vec{sk}_{II})$ //simulates \mathcal{A} honestly generating keys 78 π_{U}^{H} .role $\leftarrow x$ 79 else CHALLINIT (U, s, m, μ) CHALLRESP (U, s, m, μ) 80 π_{U}^{H} .role $\leftarrow y$ 33 if $\pi_{II}^s = \perp \lor \pi_{II}^s$.role = responder initiator //initialized sessions only 81 π_{U}^{H} .oid $\leftarrow U$ 34 return \perp $\pi_{U}^{H}.\mathsf{pid} \leftarrow H$ 82 35 **if** b = 083 $(\pi_U^H, m', \mu') \leftarrow \text{$\mathsf{Run}(\mathsf{sk}_U, \vec{\mathsf{pk}}, \pi_U^H, m, \mu)$}$ $(\pi_U^s,m') \xleftarrow{} \mathsf{Run}(\mathsf{sk}_U,\vec{\mathsf{pk}},\pi_U^s,m,\mu)$ 36 84 return m'37 else $(\pi_U^s, m') \leftarrow \mathsf{Fake}^{\mathcal{O}_F}(\vec{\mathsf{pk}}, \pi_U^s, m, \mu, \mathsf{aux}, r)$ $SK(\pi)$: SKs(): AUX(): 38 86 return sk 85 return $sk_{\pi,pid}$ 39 $Q[\pi_U^s] \xleftarrow{+} (m, m')$ 87 return aux 40 $K[\pi_{II}^s] \leftarrow \pi_{II}^s.\mathsf{K}$ 41 return m'

Figure 2: Our deniability game parameterized in the adversary's capabilities $O_A \subseteq \{\text{REG}, \text{INIT}, \text{CHALLINIT}, \text{CHALLRESP}, \text{REGHON}, \text{CHALLHONINIT}, \text{CHALLHONRESP}\}$, the capabilities of the Fake algorithm $O_F \subseteq \{\text{SK}, \text{USER}_{x,y}\}$, the capabilities of the distinguisher $\mathcal{D} \cup O_D \subseteq \{\text{SKs}, \text{AUX}\}$, and the sampling algorithm for auxiliary info AuxPrep.

Type of Deniability	$ O_A$	O_F	O_D
against malicious adversaries	Reg, Init, Chall	\diamond	\diamond
against semi-honest adversaries	REGHON, CHALLHON	\diamond	\diamond
against passive adversaries	ChallHon	\diamond	\diamond
partial deniability wrt. an x, y oracle	\diamond	$USER_{x,y}$	\diamond
1-out-of-2 deniability	\Diamond	SK	\diamond
in the big brother model	\Diamond	\diamond	SKs
with aux known to the distinguisher	\diamond	\diamond	AUX

Table 1: Translating oracles into names for adversarial capabilities and deniability guarantees. We refer to an arbitrary set with \Diamond . For the challenge oracles we omit the suffix indicating which role they act as.

3.1 AuxPrep: Sampling auxiliary information

The auxiliary info aux models that an adversary may obtain relevant information prior to the experiment. To the best of our knowledge, this idea was introduced by [DGK06]: They use aux to model that the adversary may have obtained transcripts of previous protocol executions by eavesdropping. Following [DGK06], the auxiliary info is also available to the Fake algorithm, i.e. anybody who knows the auxiliary info can simulate a transcript and session key. Optionally, the distinguisher gets access to aux as well (cp. Section 3.4.2).

3.2 O_A : capabilities of \mathcal{A}

The oracles in O_A model the behavior of the adversary, i.e. whether it is malicious or semi-honest,⁵ and define the challenge. In the following, we remark on how we handle the roles of parties and detail the usage of the respective subsets of oracles.

3.2.1 Initiator and receiver deniability

We obtain *initiator deniability* by having CHALLINIT $\in O_A$ or CHALLHONINIT $\in O_A$, and *responder deniability* by having CHALLRESP $\in O_A$ or CHALLHONRESP $\in O_A$; i.e. we give separate challenge oracles per role. This allows separate analyses of deniability per role. (Section 4 shows that the deniability guarantees for both roles differ for X3DH and PQXDH.) Also, this enables more precise specifications, e.g. deniability for the responder may have stricter requirements than for the initiator.

3.2.2 Malicious adversaries

We obtain deniability against malicious adversaries with $O_A = \{\text{REG}, \text{INIT}, \text{CHALLINIT}, \text{CHALLRESP}\}$, allowing at most q_R queries to the REG oracle, q_I queries to the INIT oracle, q_{CI} queries to the CHALLNIT oracle, and q_{CR} queries to the CHALLRESP oracle, for $q_R, q_I, q_{CI}, q_{CR} \in q_{O_A}$. The REG oracle allows the adversary to register maliciously generated keys with the challenger. The INIT oracle allows the adversary to initialize sessions for honest users and to partner these sessions arbitrarily. The CHALLINIT and CHALLRESP oracles allow the adversary to send arbitrary messages to sessions (including, e.g. messages replayed from other sessions) and to interleave sessions at will. The CHALLINIT and CHALLRESP oracles respond either according to the protocol specification via Run or "simulate" via Fake, depending on the challenge bit b. For the case of Signal's initial handshake, note that this also captures an adversarially

 $^{{}^{5}}$ In Appendix B we also treat passive adversaries for the sake of completeness. We do not consider passive adversaries for our analysis and therefore defer it to the appendix.

controlled server, i.e. a server who forwards messages wrongly. Only sessions of honestly generated users can be challenged since the INIT oracle does not create sessions for users with maliciously generated keys.

3.2.3 Semi-honest Adversaries

Next, we obtain the notion of deniability against semi-honest adversaries by setting the oracle $O_A = \{\text{RegHon}, \text{CHALLHONINIT}, \text{CHALLHONRESP}\}$, allowing at most q_{RH} queries to the RegHon oracle, q_{CI} queries to the CHALLHONINIT oracle, and q_{CR} queries to the CHALLHONRESP oracle, for q_{RH} , $q_{CI}, q_{CR} \in q_{O_A}$. To register a key pair for a new user, the adversary can query the RegHon oracle, which returns the user's private key(s) to the adversary. A semi-honest adversary partners sessions correctly and obeys the protocol flow. Hence, we merge the INIT oracle and the challenge oracles. These CHALLHONINIT, CHALLHONRESP oracles partner the sessions correctly and directly produce the complete transcript, thereby enforcing correct message flow.

Only honest users can be challenged (cp. lines 42 and 57). Note that the adversary provides the randomness to the challenge oracles to create messages "in place of the adversary"⁶. The adversary also provides the protocol-specific $info_{create}$ to indicate the components of the first message.

3.2.4 Combining different adversarial capabilities

It is possible to combine the aforementioned oracles at will, e.g. achieving responder deniability against malicious adversaries and initiator deniability against semi-honest adversaries with malicious keys via $O_A = \{\text{Reg, INIT, CHALLRESP, CHALLHONINIT}\}$. Note that O_A should always contain at least one challenge oracle; and for the CHALLINIT, CHALLRESP oracles the INIT oracle is necessary. Oracles for both malicious and semi-honest adversaries seem redundant but may have some use cases.

3.2.5 Adaptive Corruptions

Several works [CF11, JS08, Jia14, JCL⁺22] allow the adversary to adaptively corrupt parties. The attacker, the simulator, and possibly the distinguisher learn the corrupted parties' secrets. Fundamentally, deniability is a statement that a Fake algorithm can simulate a transcript (and session key). The Fake algorithm needs to work even if the adversary does not corrupt any party. Hence, the strategy of the Fake algorithm should not rely on secret keys of corrupted parties and that is why our model does not consider adaptive corruptions. Appendices A.2 and A.8 describe how to represent non-adaptive corruptions in our model.

3.3 O_F : capabilities of Fake

The oracles in O_F are accessed by the Fake algorithm. They offer help in simulating a transcript and session key and should reflect the capabilities of an adversary who wants to frame a user for a particular transcript. If fewer oracles are required, then more people can simulate, i.e. the deniability guarantee gets stronger.

3.3.1 Partial deniability (User_{x,y} oracle (parameterized by x, y))

We obtain partial deniability wrt. x, y oracles by having $USER_{x,y} \in O_F$, allowing at most q_U sessions with the $USER_{x,y}$ oracle per session π , where $x, y \in \{\text{initiator}, \text{responder}, \bot\}$. Each invocation of the Fake algorithm (lines 53, 69, and 38) has access to exactly one session π . The Fake algorithm can query the $USER_{x,y}$ oracle with this session π to get access to another session with π .oid (the owner of π) with the

⁶In consequence, the distinguisher learns this randomness as well.

role x (if π .oid is the initiator in π) or y (if π .oid is the responder in π), where \perp denotes no access. The Fake algorithm can have up to $q_U \in q_{O_F}$ sessions with the $\text{USER}_{x,y}$ oracle per session π . For $u \in [q_U]$, the user π .oid is partnered with a (freshly generated) honest user H^u_{π} (and in particular not with π .pid), where $H^u_{\pi} \in \mathcal{H}$ and $\mathcal{H} \cap [n_p] = \emptyset$.

This idea is sourced from *partial deniability* [DGK06] (informally known as *peer independence* [DGK06] or *receiver obliviousness* [HKKP22], following the idea of *post-specified peers* [CK02]).

This oracle models that an adversary can start a separate session with the victim π .oid to produce the transcript. This oracle can be helpful if messages do not contain session-specific data.

3.3.2 1-out-of-2 or 1-out-of- ∞ deniability (SK oracle)

We obtain 1-out-of-2 deniability by having $SK \in O_F$. If $SK \notin O_F$, then we obtain 1-out-of- ∞ deniability. The SK oracle grants the Fake algorithm access to the peer's secret keys, thereby modeling that the peer (and only the peer) can execute the Fake algorithm. Hence, 1-out-of- ∞ deniability is more desirable. Note that this refers to long-term and semi-static keys, i.e. keys that are registered with the challenger. In particular, the SK oracle cannot yield secret keys of users that the adversary registered via the REG oracle.

If a transcript (and session key) can stem from either an honest execution or the peer creating the transcript on his or her own (with the Fake algorithm) then there are two options for this; hence, the name 1-out-of-2 deniability. Absence of this oracle indicates that anybody could have produced a transcript, hence, we dub it 1-out-of- ∞ deniability. These terms were coined by [HLLC11] for authentication protocols⁷ and adopted for key exchange by [BFG⁺22].

The idea of using one party's secret key for faking a transcript was also used by [VGIK20], though coming from a different angle: They allow the simulator (roughly corresponding to our Fake algorithm) to use the adversary's long-term secret key to fake a transcript. They justify this with the adversary having registered its long-term key to some form of PKI, possibly with an extractable proof of knowledge, from which the simulator can learn the secret key.

3.4 O_D : capabilities of \mathcal{D}

The distinguisher always gets all users' public keys, the transcripts and session keys from the interaction with the CHALL oracle, and the randomness of the adversary. The oracles in O_D define the extra power that the distinguisher may have in the form of access to all (honestly generated) secret keys (i.e. the big brother model) or the auxiliary info. The more oracles the distinguisher has access to, the stronger the deniability guarantee we obtain.

3.4.1 Big brother model

We obtain deniability in the big brother model by having $SKs \in O_D$. By querying the SKs oracle, the distinguisher gets access to all secret keys that the challenger prepared for all users as well as secret keys registered by the adversary via the REGHON oracle. This models that the distinguisher can subpoen all parties into giving up their private keys or can exert social pressure in a different manner. Note that the SKs oracle does not yield secrets of keys registered via the REG oracle—the challenger does not know these secret keys and possibly the adversary does not either. Similarly, the SKs oracle does not yield ephemeral secret keys, as they are not registered with the challenger (and, for X3DH and PQXDH, honest parties erase ephemeral secret keys directly after use).

⁷They also call a ring signature with a ring of n members 1-out-of-n deniable. This allows for a philosophical debate whether this applies to messaging services with n registered users. We argue that anybody can sign up for the service and, hence, there is no practical difference to 1-out-of- ∞ deniability.

The term *big brother* was introduced by [Nao02] in the context of deniable ring authentication. The adversary divulging its secret key to the distinguisher was coined complete deniability by [MP02]. Informal descriptions [UDB⁺15, UG15] and the UC-based definitions of [Ung21] are set in this model without using this name. The game-based definition of [BFG⁺22] first uses the term big brother model for deniability of key exchange protocols.

3.4.2 Auxiliary info known to the distinguisher

We obtain *deniability with* aux *known to the distinguisher* by having $AUX \in O_D$. By querying the AUX oracle, the distinguisher gets access to the auxiliary information (cp. Section 3.1). This models that the distinguisher learns whatever the adversary has eavesdropped on before the beginning of the experiment. If $AUX \notin O_D$ then the deniability guarantee *relies on* the distinguisher never learning aux. This is important for analyzing PQXDH, where the Fake algorithm reuses signatures from aux, which are not reused in honest executions, and thereby allows the distinguisher to tell Run and Fake apart.

3.5 Comparability of deniability notions

Our model allows fairly easy comparisons between any two deniability notions by weighing O_A, O_F , O_D , AuxPrep between the two notions. If the adversary \mathcal{A} or the distinguisher \mathcal{D} gain more capabilities, i.e. if O_A, O_D contain more or stronger oracles, the deniability guarantee becomes stronger (since it holds in more cases). Though, the AUX oracle does not make a difference in case $\mathsf{aux} = \bot$. However, if Fake gains more capabilities, i.e. if O_F contains more oracles, the deniability guarantee becomes weaker: Simulating a transcript requires access to the oracles in O_F . Similarly, the more information is included in the auxiliary information aux , the harder it is to match these circumstances and the weaker the guarantee.

It is easy to see that an algorithm has more capabilities if it gains an extra oracle. The adversary \mathcal{A} can have access to oracles for malicious adversaries (REG, CHALLINIT, CHALLRESP) and for semi-honest adversaries (REGHON, CHALLHONINIT, CHALLHONRESP). It is clear that the former oracles are stronger than the latter, yielding a stronger deniability guarantee.⁸

This leaves some relations open, e.g. how does deniability against semi-honest adversaries in the big brother model compare to deniability against malicious adversaries? While both notions remain incomparable, our model makes the differences explicit: We can argue about the "upper bound" (malicious adversaries in the big brother model) and the "lower bound" (semi-honest adversaries). Note that a deniability guarantee also holds for smaller numbers of users and semi-static keys: If it holds for n_p, n_{ss} , it also holds for $n'_p \leq n_p, n'_{ss} \leq n_{ss}$.⁹

4 Deniability of Signal's initial handshake

We look at Signal's initial handshake protocols—X3DH and its recent replacement PQXDH adding postquantum confidentiality—and which deniability notions they fulfill. Note that both protocols aim for classical authentication and deniability, and not post-quantum deniability. Figure 3 gives an algorithmic overview of both protocols and the respective subsections give a textual description. The protocol specification [MP16b, KS23] uses only one long-term key per user, which is used for both the DH scheme and the signature scheme XEdDSA [Per16]. We follow [BJKS23] in treating them as two separate keys. We follow the protocol specification in using the session key directly as key for the AEAD scheme. Signal's implementation derives several values from the session key, among them the key and nonce used for the

⁸Assuming that the CHALLINIT, CHALLRESP oracles come with an INIT oracle, since they are otherwise useless.

⁹A reduction can simply withhold some keys from the inner adversary.

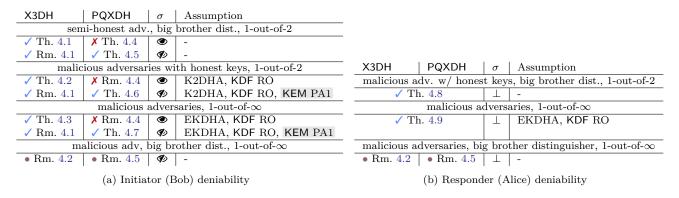


Table 2: Comparing deniability of X3DH and PQXDH: Theorems (and Remarks) in the first two columns address the model in the continuous line above (indicating adversarial capabilities, big brother model, and how many people can fake a transcript, i.e. whether $SK \in O_F$), the third column indicates if Fake requires a signature that the distinguisher is not aware of with \mathcal{O} (USER_{initiator,y} $\in O_F$ or both AuxPrep yields appropriate signatures and AUX $\notin O_D$) or a signature that the distinguisher may learn with \mathfrak{O} (USER_{initiator,y} $\in O_F$ or AuxPrep yields appropriate signatures) or no signature at all with \bot , and the final column states the assumptions needed. We denote that a property holds by \checkmark , does not hold by \bigstar , or that we currently do not have a proof by \bullet . For the knowledge of DH assumptions (K2DH, EKDHA), see Definitions 2.4 and 2.5; the plaintext awareness assumption on KEM (gray background, see Definition 2.7) only applies to PQXDH.

AEAD scheme. Beware that we consider Bob, who creates the pre-key bundle, as initiator, following e.g. [UG15, CCD⁺17, UG18] and contrasting [BFG⁺22].

We summarize our findings for both X3DH and PQXDH in Table 2, stated separately per role. For initiator (Bob) deniability, we have two results per adversary model to differentiate whether the Fake algorithm has access to a signature that the distinguisher does not know, i.e. a *private signature*, or if Fake and the distinguisher know the same signatures, i.e. *public signatures*. For X3DH either case works (see Remark 4.1), while for PQXDH we *must have* a private signature for the ephemeral KEM key (see Theorem 4.4 and Remark 4.4). Without further restrictions, we achieve 1-out-of-2 deniability against semi-honest adversaries. We use a knowledge of DH assumption (see Definitions 2.4 and 2.5) against malicious adversaries and additionally plaintext awareness of the KEM (see Definition 2.7) for PQXDH to simulate the session key. Concerning responder (Alice) deniability, PQXDH matches the guarantees of X3DH without limitations. For malicious adversaries in the big brother model, i.e. the strongest possible adversary model, we give our thoughts in Remarks 4.2 and 4.5. Note that Table 2 implies all other notions: Each theorem still holds if you weaken the adversary, weaken the distinguisher, add an oracle to O_F, or add auxiliary info (without the Fake algorithm using it).

While our analysis relies on the random oracle model, we do not program the RO, following [DGK06] and in line with the results of [Pas03]. All Fake algorithms abort on receiving malformed inputs. We state the expected message format explicitly in the pseudocode.

4.1 X3DH

We recall how Signal's classical initial handshake protocol X3DH [MP16b] works, see Figure 3. First, Bob signs his semi-static public key (or retrieves a previously created signature) under his long-term key. Bob samples an ephemeral DH key pair if the boolean e_{DH} is set. Bob's three public keys (long-term, semi-static with id ssid, and optionally ephemeral) and the signature are Bob's first message. We call this first message Bob's *pre-key bundle*, which he sends to the key server and is not tied to any particular peer. Second, Alice parses the message, verifies the signature, samples an ephemeral DH key herself, and computes four DH shared secrets (or three if Bob's message does not include an ephemeral key) between her keys and Bob's keys (long-term-semi-static, ephemeral-long-term, ephemeral-semi-static, ephemeral),

KGenLT:

- 1 $(ltpk_U^{\text{DH}}, ltsk_U^{\text{DH}}) \leftarrow \text{SDH.KGen}()$ 2 $(ltpk_U^{\text{SIG}}, ltsk_U^{\text{SIG}}) \leftarrow \text{SIG.KGen}()$
- 3 return $\left((ltpk_U^{\mathsf{DH}}, ltpk_U^{\mathsf{SIG}}), (ltsk_U^{\mathsf{DH}}, ltsk_U^{\mathsf{SIG}})\right)$
- Alice

 m_1 $\operatorname{\mathsf{Run}}(ltsk_A, ss\vec{sk}_A, lt\vec{pk}, ss\vec{pk}, \pi_A, m_1, \mu_1)$ $(epk_A^{\mathsf{DH}}, esk_A^{\mathsf{DH}}) \leftarrow \mathsf{SDH}.\mathsf{KGen}()$ $(B, \mathsf{ssid}, \sigma_{\mathsf{DH}}^{\mathsf{ssid}}, epk_B, \sigma_{\mathsf{KEM}}) \leftarrow m_1$ $(sspk_B^{\mathsf{DH}}, sspk_B^{\mathsf{KEM}}) \leftarrow sspk_B^{\mathsf{ssid}}$ $(epk_B^{\mathsf{DH}}, epk_B^{\mathsf{KEM}}) \leftarrow epk_B$ if SIG.Vf $(ltpk_B, sspk_B^{\mathsf{DH}}, \sigma_{\mathsf{DH}}^{\mathsf{ssid}}) = \mathsf{false}$ return $(\pi_A, \epsilon, \epsilon)$ $DH_1 \leftarrow \mathsf{DH}(ltsk_A, sspk_B^{\mathsf{DH}})$ $DH_2 \leftarrow \mathsf{DH}(esk_A^{\mathsf{DH}}, ltpk_B)$ $DH_3 \leftarrow \mathsf{DH}(esk_A^{\mathsf{DH}}, sspk_B^{\mathsf{DH}})$ if $epk_B^{\mathsf{DH}} \neq \perp$ //ephemeral DH key present $DH_4 \leftarrow \mathsf{DH}(esk_A^{\mathsf{DH}}, epk_B^{\mathsf{DH}})$ else $DH_4 \leftarrow \epsilon$ if $epk_B^{\mathsf{KEM}} \neq \bot$ //ephemeral KEM key present $\mathbf{if} ~ \mathsf{SIG}.\mathsf{Vf}(ltpk_B,epk_B^{\mathsf{KEM}},\sigma_{\mathsf{KEM}}) = \mathsf{false}$ return $(\pi_A, \epsilon, \epsilon)$ $(ct_{\mathsf{KEM}}, ss) \leftarrow \mathsf{KEM}.\mathsf{Enc}(epk_B^{\mathsf{KEM}})$ else //no ephemeral KEM key present if SIG.Vf $(ltpk_B, sspk_B^{\text{KEM}}, \sigma_{\text{KEM}}) = false$ return $(\pi_A, \epsilon, \epsilon)$ $(ct_{\mathsf{KEM}}, ss) \leftarrow \mathsf{KEM}.\mathsf{Enc}(sspk_B^{\mathsf{KEM}})$ $\mathsf{ms} \leftarrow DH_1 \| DH_2 \| DH_3 \| DH_4 \| ss$ $\pi_A.\mathsf{K} \leftarrow \mathsf{KDF}(\mathsf{ms})$ $AD \leftarrow ltpk_A || ltpk_B$ $ct_{\mathsf{AE}} \leftarrow \mathsf{Enc}(\pi_A.\mathsf{K}, AD, \mu_1)$ π_A .pid $\leftarrow B$ $m_2 \leftarrow (A, epk_A^{\mathsf{DH}}, ct_{\mathsf{AE}}, ct_{\mathsf{KEM}})$ return (π_A, m_2, ϵ)

KGenSS:

```
4 (sspk_U^{\mathsf{DH}}, sssk_U^{\mathsf{DH}}) \leftarrow \mathsf{SDH}.\mathsf{KGen}()
```

- 5 $(sspk_{U}^{\mathsf{KEM}}, sssk_{U}^{\mathsf{KEM}}) \leftarrow \mathsf{KEM}.\mathsf{KGen}()$
- 6 return $((sspk_{U}^{\mathsf{DH}}, sspk_{U}^{\mathsf{KEM}}), (sssk_{U}^{\mathsf{DH}}, sssk_{U}^{\mathsf{KEM}}))$

Bob

```
\operatorname{Run}(ltsk_B, lt\vec{p}k, ss\vec{s}k_B, s\vec{s}pk, \pi_B, m_0, \mu_0)
  (create, (ssid, e_{DH}, e_{KEM})) \leftarrow m_0
  \pi_B.pid \leftarrow \star
  (sspk_B^{\mathsf{DH}}, sspk_B^{\mathsf{KEM}}) \leftarrow sspk_B^{\mathsf{ssid}}
  \begin{array}{l} \text{if } \sigma_{\text{DH}}^{\text{ssid}} = \bot \text{ //saved from a previous run?} \\ \sigma_{\text{DH}}^{\text{ssid}} \leftarrow \hspace{-0.5mm} \text{s} \hspace{-0.5mm} \text{Sign}(ltsk_B, sspk_B^{\text{DH}}) \end{array}
  if e_{DH} = true
        (epk_B^{\mathsf{DH}}, esk_B^{\mathsf{DH}}) \leftarrow \mathsf{SDH}.\mathsf{KGen}()
   else epk_B^{\mathsf{DH}} \leftarrow \bot
   if e_{\text{KEM}} = \text{true}
         (epk_B^{\text{KEM}}, esk_B^{\text{KEM}}) \leftarrow \text{sKEM.KGen}()
         \sigma_{\mathsf{KEM}} \leftarrow \mathsf{s} \operatorname{Sign}(ltsk_B, epk_B^{\mathsf{KEM}})
   else
         if \sigma_{\text{KEM}}^{\text{ssid}} = \perp //saved from a previous run?
               \sigma_{\mathsf{KEM}}^{\mathsf{ssid}} \gets \hspace{-0.15cm} \$ \operatorname{Sign}(ltsk_B, sspk_B^{\mathsf{KEM}})
         \sigma_{\text{KEM}} \leftarrow \sigma_{\text{KEM}}^{\text{ssid}}
        epk_B^{\mathsf{KEM}} \gets \bot
  epk_B \leftarrow (epk_B^{\mathsf{DH}}, epk_B^{\mathsf{KEM}})
  m_1 \leftarrow (B, \mathsf{ssid}, \sigma_{\mathsf{DH}}^{\mathsf{ssid}}, epk_B, \sigma_{\mathsf{KEM}})
  return (\pi_B, m_1, \epsilon)
m_2
  \operatorname{Run}(ltsk_B, ss\vec{sk}_B, lt\vec{pk}, ss\vec{pk}, \pi_B, m_2, \mu_2)
  (sssk_B^{\mathsf{DH}}, sssk_B^{\mathsf{KEM}}) \leftarrow sssk_B^{\mathsf{ssid}}
  (A, epk_A^{\mathsf{DH}}, ct_{\mathsf{AE}}, ct_{\mathsf{KEM}}) \leftarrow m_2
```

 $DH_1 \leftarrow \mathsf{DH}(ltpk_A, sssk_B^{\mathsf{DH}})$ $DH_2 \leftarrow \mathsf{DH}(epk_A^{\mathsf{DH}}, ltsk_B)$ $DH_3 \leftarrow \mathsf{DH}(epk_A^{\mathsf{DH}}, sssk_B^{\mathsf{DH}})$ if $epk_B^{\mathsf{DH}} \neq \perp$ //ephemeral DH key present $DH_4 \leftarrow \mathsf{DH}(epk_A^{\mathsf{DH}}, esk_B^{\mathsf{DH}})$ else then $DH_4 \leftarrow \epsilon$ if $epk_B^{\mathsf{KEM}} \neq \bot$ //ephemeral KEM key present $ss \leftarrow \mathsf{KEM}.\mathsf{Dec}(esk_B^{\mathsf{KEM}}, ct_{\mathsf{KEM}})$ else //no ephemeral KEM key present $ss \leftarrow \mathsf{KEM}.\mathsf{Dec}(sssk_B^{\mathsf{KEM}}, ct_{\mathsf{KEM}})$ $\mathsf{ms} \leftarrow \mathsf{KDF}(DH_1 \| DH_2 \| DH_3 \| DH_4 \| ss)$ $\pi_B.\mathsf{K} \leftarrow \mathsf{KDF}(\mathsf{ms})$ $AD \leftarrow ltpk_A \| ltpk_B$ $\mu' \leftarrow \mathsf{Dec}(\pi_A.\mathsf{K}, AD, ct_{\mathsf{AE}})$ if $\mu' = \bot$ then $\pi_B.\mathsf{K} \leftarrow \bot$ $\pi_B.\mathsf{pid} \leftarrow A$ return (π_B, ϵ, μ')

Figure 3: Signal's initial handshake protocols X3DH and PQXDH. The KEM with gray background is exclusive to PQXDH.

and finally derives the session key from the DH shared secrets. She encrypts her first user message with the AEAD scheme under the session key¹⁰ and sends this AEAD ciphertext and her ephemeral public key back to Bob. Third, Bob computes the DH shared secrets and derives the session key in the same way. He tries to decrypt the AEAD ciphertext and rejects the session key if decryption fails.

To simulate a transcript and session key, we face two challenges: obtaining all DH shared secrets and obtaining a valid signature on Bob's semi-static key. We briefly discuss our approaches for each of the two challenges, as well as the difficulties arising with big brother distinguishers, before going into the actual proofs.

4.1.1 Obtaining the DH shared secrets

For adversaries that create keys honestly, the Fake algorithm knows the necessary DH secret keys to compute all DH shared secrets itself (via the SK oracle or the semi-honest peer's randomness), as shown in Figure 5 and Figure 7 on the left-hand side. Against malicious adversaries, the Fake algorithm uses the knowledge of DH extractor (of the K2DH or EKDH assumption, see Definitions 2.4 and 2.5) to extract the remaining DH shared secrets from the adversary, as shown in Figures 5 and 6 on the right-hand side, following [VGIK20]. By assumption, if the extractor fails, no other extractor can succeed in a given time with more than a given probability. This allows Fake to set a random session key in case the extraction fails.

4.1.2 Obtaining a signature on Bob's semi-static key

Recall that the semi-static key can be used for several sessions, so neither the semi-static key nor the signature on the semi-static key are tied to a particular session. The Fake algorithm has two options (shown in Figure 4) to obtain a signature on Bob's semi-static key: First, from a pre-key bundle in the auxiliary info (the pre-key bundle may be part of a complete transcript), as previously done by [VGIK20]. Second, Fake may query the USER_{initiator,⊥} oracle (i.e. the oracle acts as initiator) to get a pre-key bundle including the signature on the semi-static key. Note that the Fake algorithm may use the SK oracle to learn the peer's secret key but does not have an option to learn the secret key of the session owner.

Both options rely on Fake obtaining a signature from the pre-key bundle. The first option (using aux) models that the adversary is able to eavesdrop on previous protocol executions and, hence, Fake can use aux as well. The second option (the USER_{initiator, \perp} oracle) models that Fake may start an independent session with the victim, i.e. anybody who can start a session with the victim can produce a transcript and session key. If Fake uses aux and aux is known to the distinguisher, then we have *public signatures*. Otherwise, we have *private signatures*, i.e. Fake knows signatures that the distinguisher does not. For X3DH either setting works. Looking ahead, this will make a difference for PQXDH.

4.1.3 Big Brother distinguishers

Ideally, we want to have deniability results against malicious adversaries in the big brother model. Against malicious adversaries, we rely on the knowledge of DH assumption, which we cannot apply in the big brother model. We discuss the challenges to show responder deniability in the big brother model in Remark 4.2.

4.1.4 Theorems for X3DH

Here we give our results for initiator deniability of X3DH. For responder deniability, we show the same results for both X3DH and PQXDH. We state them in Theorems 4.8 and 4.9.

¹⁰In practice, Signal uses an elaborate key scheduling algorithm. The actual key and nonce used for encryption are deterministically derived from the session key.

```
Fake(\vec{pk}, \pi, m, \mu, aux, r_C):
   1
         (create, (ssid, e_{DH}, e_{KEM}) \leftarrow m
                                                                                                                                                                                             Fake(\vec{pk}, \pi, m, \mu, aux, r_C):
   2
         \pi.pid \leftarrow \star
         (\sigma_{\mathsf{DH}}, \sigma_{\mathsf{KEM}}^{\mathsf{ssid}}) \leftarrow \text{from aux for } sspk_B^{\mathsf{ssid}}
                                                                                                                                                                                            14 (create, (ssid, e_{\mathsf{DH}}, e_{\mathsf{KEM}})) \leftarrow m
  3
  \begin{array}{l} 4 \quad \text{if} \ e_{\mathsf{DH}} = \mathsf{true} \\ 5 \quad (epk_B^{\mathsf{DH}}, esk_B^{\mathsf{DH}}) \leftarrow \mathsf{SDH}.\mathsf{KGen}() \end{array}
                                                                                                                                                                                            15 \pi.pid \leftarrow \star
                                                                                                                                                                                           16 (B, \text{ssid}, \sigma_{\mathsf{DH}}^{\mathsf{ssid}}, epk_B, \sigma_{\mathsf{KEM}}) \leftarrow \operatorname{USER}_{\mathsf{initiator}, \perp}(\pi, 1, m, \epsilon)
   6 else epk_B^{\mathsf{DH}} \leftarrow \bot
                                                                                                                                                                                            17 (epk_B^{\mathsf{DH}}, epk_B^{\mathsf{KEM}}) \leftarrow epk_B
         \mathbf{if} \ e_{\mathsf{KEM}} = \mathsf{true}
   7
                                                                                                                                                                                            \begin{array}{ll} \text{18} & \textbf{if} \ e_{\mathsf{DH}} = \textbf{true} \\ \text{19} & (epk_B^{\mathsf{DH}}, esk_B^{\mathsf{DH}}) \gets \texttt{\$DH}.\mathsf{KGen}() \end{array} 
                 (epk_B^{\mathsf{KEM}}, \sigma_{\mathsf{KEM}}) \leftarrow \text{from aux}
   8
                                                                                                                                                                                           20 else epk_B^{\mathsf{DH}} \leftarrow \bot
   9
           else
                epk_B^{\mathsf{KEM}} \leftarrow \bot
10
                                                                                                                                                                                           \texttt{21} \hspace{0.1in} epk_B \leftarrow (epk_B^{\mathsf{DH}}, \hspace{0.1in} epk_B^{\mathsf{KEM}})
             \sigma_{\mathsf{KEM}} \leftarrow \sigma_{\mathsf{KEM}}^{\mathsf{ssid}}
11
                                                                                                                                                                                            22 return (\pi, (B, \text{ssid}, \sigma_{\text{DH}}^{\text{ssid}}, epk_B, \sigma_{\text{KEM}}), \epsilon)
\texttt{12} \hspace{0.2cm} epk_B \leftarrow (epk_B^{\mathsf{DH}}, \hspace{0.2cm} epk_B^{\mathsf{KEM}})
 13 return (\pi, (B, ssid, \sigma_{DH}^{ssid}, epk_B, \sigma_{KEM}), \epsilon)
```

Figure 4: The Fake algorithms simulating Bob's pre-key bundle (initiator deniability) in Theorems 4.1 to 4.3 and 4.5 to 4.7.

We show that Bob can deny (*initiator deniability*) participation in a protocol run against a semi-honest Alice (*against semi-honest adversaries*, Theorem 4.1), against a malicious Alice that honestly generated her keys (*against malicious adversaries with honest keys*, Theorem 4.2), and against a malicious Alice (*against malicious adversaries*, Theorem 4.3). He does so by showing that Alice (1-*out-of*-2, i.e. anybody using Alice's secret keys, Theorems 4.1 and 4.2) or anybody $(1-out-of-\infty, \text{Theorem 4.3})$ could have produced the transcript and session key herself, assuming she has previously observed a valid transcript with Bob as initiator (AuxPrep yielding a valid pre-key bundle per ssid and user). This holds even against a big brother distinguisher (*in the big brother model*, Theorem 4.1), who shares Alice's observations (aux known to the distinguisher, Theorems 4.1 to 4.3).

Theorem 4.1. The X3DH protocol as shown in Figure 3 is $(n_p, n_{ss}, q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D, \epsilon_D)$ -deniable with respect to ({REGHON, CHALLHONINIT}, {SK}, {SKS, AUX})-oracles and auxiliary info sampled with AuxPrep yielding a valid pre-key bundle per ssid and user, and aux known to the distinguisher, where n_p is the number of parties and n_{ss} the number of semi-static keys per party, and $q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D$ are arbitrary and $\epsilon_D = 0$.

Proof. We give the Fake algorithm in Figure 4 on the left-hand side and Figure 5 on the left-hand side.

The Fake algorithm simulates Bob's pre-key bundle by sampling a fresh ephemeral key (if needed) and learns a signature on the semi-static key from the auxiliary info aux. Furthermore, to process Alice's message, the Fake algorithm learns Alice's long-term secret key from the SK oracle and Alice's ephemeral secret key esk_A^{DH} from the randomness that was previously used to create Alice's message. Using these secret keys, it can compute DH_1 (with esk_A^{DH}), and DH_2 , DH_3 , DH_4 (with $ltsk_A^{DH}$), and in consequence the session key. Using the session key, Fake decrypts the AEAD ciphertext. If decryption fails, it aborts.

Since Fake produces a transcript and session key in the same way as Run, the distinguisher cannot have an advantage in winning his game, even if the distinguisher has access to all long-term and semi-static secret keys.

Remark 4.1 (Obtaining the signature). In Theorems 4.1 to 4.3, the Fake algorithm obtains the signature on the semi-static DH key from the auxiliary info aux (cp. the left-hand side of Figure 4). Instead, the Fake algorithm may also learn the signature via the USER_{initiator,⊥} oracle (cp. the right-hand side of Figure 4), resulting in partial deniability wrt. a USER_{initiator,⊥} oracle. Hence, these theorems also hold for AuxPrep = \perp if O_F additionally includes USER_{initiator,⊥}.

The following two theorems use knowledge of DH assumptions, which [VGIK20] have introduced to show deniability of X3DH in their model. The first theorem adapts the proof of [VGIK20, Theorem 6] to

41 π .pid $\leftarrow A$ 42 $(ltsk_A^{\mathsf{DH}}, ltsk_A^{\mathsf{SIG}}) \leftarrow \text{extract from SK}(\pi)$ $\mathsf{Fake}(\vec{\mathsf{pk}}, \pi, m = (A, epk_A^{\mathsf{DH}}, ct_{\mathsf{AE}}, ct_{\mathsf{KEM}}), \mu, \mathsf{aux}, r_C):$ 43 $esk_B^{\text{DH}} \leftarrow \text{from previous run}$ 23 π .pid $\leftarrow A$ 44 $(sspk_B^{\mathsf{DH}}, sspk_B^{\mathsf{KEM}}) \leftarrow sspk_B^{\mathsf{ssid}}$ 24 $(ltsk_A^{\mathsf{DH}}, ltsk_A^{\mathsf{SIG}}) \leftarrow \text{extract from SK}(\pi)$ 25 $esk_A^{\mathsf{DH}} \leftarrow \text{from } r_C$ $\begin{array}{l} 45 \quad DH_1 \leftarrow \mathsf{DH}(ltsk_A^{\mathsf{DH}}, sspk_B^{\mathsf{ssid}}) \\ 46 \quad (DH_2, DH_3) \leftarrow \mathcal{E}_{\mathcal{A}}^{DH}(ltgk_B^{\mathsf{DH}}, sspk_B^{\mathsf{ssid}}, \mathsf{aux}, r) \\ \end{array}$ 26 $(sspk_B^{\mathsf{DH}}, sspk_B^{\mathsf{KEM}}) \leftarrow sspk_B^{\mathsf{ssid}}$ 47 if $epk_B \neq \perp$ //ephemeral DH key present 27 $DH_1 \leftarrow \mathsf{DH}(ltsk_A^{\mathsf{DH}}, sspk_B^{\mathsf{ssid}})$ $DH_4 \leftarrow \mathsf{DH}(epk_A^{\mathsf{DH}}, esk_B^{\mathsf{DH}})$ 48 $28 \quad DH_1 \leftarrow \mathsf{DH}(esk_A^{\mathsf{DH}}, ltpk_B^{\mathsf{DH}})$ 49 else 29 $DH_3 \leftarrow \mathsf{DH}(esk_A^{\mathsf{DH}}, sspk_B^{\mathsf{ssid}})$ 50 $DH_4 \leftarrow \epsilon$ if $epk_{B}^{\text{KEM}} \neq \perp$ //ephemeral KEM key present 30 if $epk_B \neq \perp$ //full handshake 51 $DH_4 \leftarrow \mathsf{DH}(esk_A^{\mathsf{DH}}, epk_B^{\mathsf{DH}})$ 31 $ss \leftarrow \mathcal{E}_{AOA}^{PA}(epk_B^{\mathsf{KEM}}, ct_{\mathsf{KEM}}, r)$ 52 32 else //reduced handshake 53 else //no ephemeral KEM key present $DH_4 \leftarrow \epsilon$ 33 $ss \leftarrow {}^{\$} \mathcal{E}_{\mathcal{A}^{\mathcal{O}_{A}}}^{PA}(sspk_{B}^{\mathsf{KEM}}, ct_{\mathsf{KEM}}, r)$ 54 34 $ss \leftarrow \text{from } r_C$ //encapsulation against epk_B^{KEM} or $sspk_B^{\text{KEM}}$ 35 $\mathsf{ms} \leftarrow DH_1 \| DH_2 \| DH_3 \| DH_4 \| ss$ 55 **if** $DH_2 = \bot \lor DH_3 = \bot$ 36 $\pi.K \leftarrow KDF(ms)$ 56 37 $AD \leftarrow ltpk_A \| ltpk_B$ 57 else 38 $\mu' \leftarrow \mathsf{Dec}(\pi,\mathsf{K},AD,ct_{\mathsf{AE}})$ $\mathsf{ms} \leftarrow DH_1 \| DH_2 \| DH_3 \| DH_4 \| ss$ 58 39 if $\mu' = \bot$ then $\pi.K \leftarrow \bot$ $\pi.\mathsf{K} \leftarrow \mathsf{KDF}(\mathsf{ms})$ 59 60 $AD \leftarrow ltpk_A \| ltpk_B$ 40 return (π, ϵ, μ') 61 $\mu' \leftarrow \mathsf{Dec}(\pi,\mathsf{K},AD,ct_{\mathsf{AE}})$ 62 if $\mu' = \bot$ then $\pi.\mathsf{K} \leftarrow \bot$ 63 return (π, ϵ, μ')

Figure 5: The Fake algorithm processing Alice's message (initiator deniability) in Theorems 4.1 and 4.5 (on the left-hand side) and Theorems 4.2 and 4.6 (on the right-hand side).

 $\mathsf{Fake}(\vec{\mathsf{pk}}, \pi, m = (A, epk_A^{\mathsf{DH}}, ct_{\mathsf{AE}}, ct_{\mathsf{KEM}}), \mu, \mathsf{aux}, r)$:

our model. Since [VGIK20] considers only keys of two lifetimes (which they call long-term and ephemeral), we extrapolate to long-term, semi-static, and ephemeral keys. Their ephemeral key gets signed, just like the semi-static key in our description of X3DH. And if their ephemeral keys were used only once, then the simulator using a signature from the auxiliary info would raise the distinguisher's suspicion (since the distinguisher knows the auxiliary information as well). Hence, we consider their "ephemeral" keys as semi-static. Furthermore, they do not consider the AEAD ciphertext.

The proof of [VGIK20] relies on key registration: They argue that at the time of key registration a user needs to provide an extractable proof of knowledge of the secret key. Hence, they provide the long-term secret key of the adversary to the simulator. We model this by limiting the adversary to honestly generated keys with the REGHON oracle and by giving the Fake algorithm access to the SK oracle.¹¹

Theorem 4.2. If the $(t_{A^{DH}}, t_{\mathcal{E}^{DH}}, t_{\mathcal{D}^{DH}}, \epsilon_{\mathcal{D}^{DH}})$ -Knowledge of 2DH (K2DH) Assumption holds and KDF is a random oracle, then the X3DH protocol as shown in Figure 3 is $(n_p, n_{ss}, q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D, \epsilon_D)$ -deniable with respect to ({REGHON, INIT, CHALLINIT}, {SK}, {AUX})-oracles and auxiliary inputs sampled with AuxPrep yielding a valid pre-key bundle per ssid and user, and aux known to the distinguisher, where n_p is the number of parties and n_{ss} the number of semi-static keys per party, $q_{O_A}, q_{O_F}, q_{O_D}$ are arbitrary, and $t_A \approx t_{A^{DH}}, t_D \approx t_{D^{DH}}$, and $\epsilon_D \leq \epsilon_{D^{DH}}$.

Proof. We give the Fake algorithm in Figure 4 on the left-hand side and Figure 5.

The Fake algorithm uses Alice's secret key from the SK oracle to compute DH_1 , Bob's ephemeral secret, which it previously sampled, to compute DH_4 , and the extractor from the K2DH assumption to learn DH_2 and DH_3 . If the extraction fails, Fake sets a random key as session key. Using the session key, Fake decrypts the AEAD ciphertext. If decryption fails, it aborts.

¹¹These oracles also return semi-static secret keys, which we deem in spirit consistent with the idea of [VGIK20].

Figure 6: The Fake algorithm processing Alice's message (initiator deniability) in Theorems 4.3 and 4.7.

Formally, the adversary \mathcal{A} never outputs any value. Here, we apply the extractor $\mathcal{E}_{\mathcal{A}}^{DH}$ in line 46 to \mathcal{A} querying the CHALL oracle with the message $(A, epk_{\mathcal{A}}^{DH}, ct_{\mathsf{AE}})$ after receiving a message from Bob with semi-static key ssid.

Since Fake produces the transcript in the same way as Run, it does not help the distinguisher to win. If the extractor succeeds, then Fake has computed the session key in the same way as Run, allowing no advantage for the distinguisher. If the extractor fails, then the distinguisher succeeds in extracting the DH_2 or DH_3 with a maximum probability of $\epsilon_{\mathcal{D}^{DH}}$. Hence, the distinguisher can distinguish the real session key from the simulated session key with a probability of at most $\epsilon_{\mathcal{D}^{DH}}$.

Theorem 4.3. If the $(t_{A^{DH}}, t_{\mathcal{E}^{DH}}, t_{\mathcal{D}^{DH}})$ -Extended Knowledge of DH (EKDH) Assumption holds and KDF is a random oracle, then the X3DH protocol as shown in Figure 3 is $(n_p, n_{ss}, q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D, \epsilon_D)$ -deniable with respect to ({REG, INIT, CHALLINIT}, \emptyset , {AUX})-oracles and auxiliary info sampled with AuxPrep yielding a valid pre-key bundle per ssid and user, and aux known to the distinguisher, where n_p is the number of parties and n_{ss} the number of semi-static keys per party, $q_{O_A}, q_{O_F}, q_{O_D}$ are arbitrary, and $t_A \approx 3 \cdot t_{A^{DH}}, t_D \approx t_{D^{DH}}$, and $\epsilon_D \leq \epsilon_{D^{DH}}$.

Proof. We give the Fake algorithm in Figure 4 on the left-hand side and Figure 6.

The Fake algorithm simulates Bob's pre-key bundle by sampling a fresh ephemeral key (if needed) and learns a signature on the semi-static key from the auxiliary info aux. Furthermore, to process Alice's message the Fake algorithm uses two techniques: It relies on the extractor from the EKDH assumption to learn DH_1 through DH_3 and uses Bob's ephemeral key (from its previous run) to compute DH_4 . If any of the three extractions fail, Fake sets a random key as session key; otherwise it computes the session key normally. It decrypts the AEAD ciphertext using the session key and erases the session key in case decryption fails.

Formally, the adversary \mathcal{A} never outputs any value. Here, we apply the extractor $\mathcal{E}_{\mathcal{A}}^{DH}$ to those cases were \mathcal{A} queries some oracle with a message: In line 67, the adversary's output refers to \mathcal{A} querying the REG oracle (with a long-term and a semi-static public key, but we only care about the long-term key); in lines 68 and 69, the adversary's output refers to \mathcal{A} querying the CHALLINIT oracle on message $(A, epk_{\mathcal{A}}^{\mathsf{DH}}, ct_{\mathsf{AE}})$.

The Fake algorithm produces the transcript in the same way as Run and, hence, the transcript is indistinguishable. If the extractor succeeds in all three cases, then the Fake algorithm has computed the session key in the same way as Run, i.e. indistinguishably. If the extractor fails in any of the three cases, then the distinguisher succeeds in extracting one of the shared DH secrets with a maximum probability of $\epsilon_{\mathcal{D}^{DH}}$ as per the EKDH assumption. Hence, the distinguisher can distinguish the real session key from the simulated session key with a probability of at most $\epsilon_{\mathcal{D}^{DH}}$. Note that the second and third extractor call refer to the same output of the adversary, i.e. they are not independent. Conservatively, we bound with $\epsilon_{\mathcal{D}^{DH}}$ and not $(\epsilon_{\mathcal{D}^{DH}})^3$.

Remark 4.2 (Deniability against malicious adversaries in the big brother model). In the big brother setting, we cannot apply the EKDH assumption (since the distinguisher gets extra inputs compared to the adversary and Fake). Hence, if the extractor fails a big brother distinguisher may be able to compute the correct session key with probability $> \epsilon_{DDH}$. Thereby, the big brother distinguisher can tell the two cases apart. We discuss another (unsuccessful) strategy in Section 5.

4.2 PQXDH

We review PQXDH [KS23] with the help of Figure 3. PQXDH extends X3DH by including a KEM to achieve post-quantum confidentiality. In particular, Bob's pre-key bundle includes an ephemeral KEM key (if the boolean e_{KEM} is set) or a semi-static KEM key (otherwise), and a signature on this KEM public key. Alice verifies the signature on the KEM key (and on the semi-static DH key as before) and encapsulates against Bob's KEM public key. Alice's message consists of her ephemeral DH key, the AEAD ciphertext, and the KEM shared secret, and use the session key for encryption and decryption of the AEAD ciphertext, respectively.

Note that in practice, Bob's pre-key bundle includes a semi-static KEM key only if no ephemeral KEM key is used. We model this by not using the semi-static KEM key if an ephemeral KEM key exists; though we formally unwrap it from the semi-static key.

Since PQXDH builds on top of X3DH, simulating a transcript bears the same difficulties as for X3DH with some new ones added on top: obtaining a signature for the KEM key (which may be ephemeral or semi-static) and the KEM shared secret.

4.2.1 Obtaining a signature on Bob's semi-static and ephemeral keys

We need a signature on Bob's semi-static DH key for X3DH, and also on Bob's ephemeral or semi-static KEM key for PQXDH. So we need to cover signatures on semi-static *and* ephemeral keys.

We have the same two options as for X3DH to obtain a signature: from aux and from the USER_{initiator,y} oracle, see Figure 4. If Fake uses aux and aux is known to the distinguisher, then we have *public signatures*. Otherwise, we have *private signatures*, i.e. Fake knows signatures that the distinguisher does not. The distinguisher can detect reuse of public signatures, see Theorem 4.4. Hence, we require private signatures, see Remark 4.4.

4.2.2 Obtaining the KEM shared secret

For responder deniability, Fake can learn the KEM shared secret by encapsulating against the KEM public key. For initiator deniability against malicious adversaries, we rely on plaintext awareness of the KEM (see Definition 2.7), following [HKKP22].¹²

4.2.3 Theorems for PQXDH

Here we give our results for initiator deniability of PQXDH and or responder deniability of both X3DH and PQXDH. We start with initiator deniability and show that faking a transcript requires a signature that the distinguisher is not aware of, i.e. a private signature. In particular, we show that Bob cannot deny (*initiator deniability*) participation in a protocol run against a semi-honest Alice (*against semi-honest*)

 $^{^{12}}$ For initiator deniability against semi-honest adversaries, Fake can learn the randomness used for encapsulation and thereby obtain the KEM shared secret.

adversaries). This follows from nobody, not even Alice (1-out-of-2, i.e. using Alice's secret keys) being able to produce the transcript and session key herself, even assuming she has previously observed a valid transcript with Bob as initiator (AuxPrep yielding a valid pre-key bundle per ssid and user). This holds even against a big brother distinguisher (in the big brother model), who shares Alice's observations (aux known to the distinguisher).

An indistinguishable transcript contains a valid signature on the ephemeral KEM key and the distinguisher must not have seen that signature before. Since the Fake algorithm cannot get this signature from anywhere, it must be a forgery of the signature scheme.

Theorem 4.4. If SIG is $(q_{SIG}, t_{SIG}, \epsilon_{SIG})$ -unforgeable, KEM has a key collision probability bounded by γ_{coll} , and the public key spaces of KEM and DH are disjoint, then the PQXDH protocol as shown in Figure 3 is not $(n_p, n_{ss}, q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D, \epsilon_D)$ -deniable with respect to ({REGHON, CHALLHONINIT}, {SK}, {SKS, AUX})-oracles and auxiliary inputs sampled with AuxPrep yielding a valid pre-key bundle per ssid and user and q_{CI} pre-key bundles per user with ephemeral KEM keys, where $n_p \geq 2$ is the number of parties and n_{ss} the number of semi-static keys per party, $q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D$ are arbitrary, $q_{SIG} = 2 \cdot n_{ss} + q_{CI}$, $t_{SIG} \approx t_A$, and $\epsilon_D \not\leq 1 - n_{ss} \cdot \gamma_{coll} - \epsilon_{SIG}$.

Proof. For contradiction, we assume PQXDH was $(n_p, n_{ss}, q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D, \epsilon_D)$ -deniable and build a reduction to $(q_{SIG}, t_{SIG}, \epsilon_{SIG})$ -unforgeability of SIG. Note that the reduction relies on the Fake algorithm to win its unforgeability game. As such, the reduction simulates both the deniability game and the deniability adversary around the Fake algorithm. The reduction receives a challenge key pk_{SIG} as input and proceeds to run KeyPrep (n_p, n_{ss}) to generate key pairs for all users. Wlog. we determine two distinct users A and B. The reduction saves the long-term signing key pair of B as $(ltpk_B^{SIG}, ltsk_B^{SIG})$, replaces the public key with the challenge key pk_{SIG} , and saves the resulting list of public keys as pk. The reduction queries its SIGN oracle to obtain signatures under pk_{SIG} for all of B's semi-static DH keys and KEM keys and for q_{CI} ephemeral KEM keys, totaling $2 \cdot n_{ss} + q_{CI}$ queries. Since the reduction controls both the (deniability) challenger and the $(ltsk_A, sssk_A)$. The Fake algorithm has to answer for this query. If Fake queries its SK oracle, the reduction replies with $(ltsk_A, sssk_A)$. The Fake algorithm has to answer for this query. If Fake queries its SK oracle, the reduction replies with $(ltsk_A, sssk_A)$. The Fake algorithm returns a protocol message containing a signed ephemeral KEM key for B. The reduction returns this ephemeral KEM key and signature tuple $(epk_B^{KEM}, \sigma_{KEM})$.

We include the running time of the oracles $O_A = \{\text{REGHON}, \text{CHALLHONINIT}\}, O_F = \{\text{SK}\} \text{ in } t_A, allowing for <math>t_{\mathsf{SIG}} \approx t_A$. The reduction simulates the game $\mathcal{G}_{\mathsf{PQXDH},\mathsf{AuxPrep},n_p,n_{ss}}^{O_A,O_F,O_D-den}(\mathcal{A},\mathcal{D})$ around Fake faithfully, since it can mimic the behavior of the game and answer as the game would.

The distinguisher can tell if Fake uses a signature on the ephemeral key from aux since the distinguisher knows aux itself (via the AUX oracle). We exclude this by assuming PQXDH to be deniable. The probability of epk_B^{KEM} colliding with one of B's n_{ss} KEM keys signed in AuxPrep is $n_{ss} \cdot \gamma_{coll}$ (and 0 to collide with a DH public key). We need to exclude this probability since these signatures were created by the SIGN oracle of the reduction. If the transcript contained an invalid signature, then the transcript could not be $(t_{\mathcal{D}}, \epsilon_{\mathcal{D}})$ -indistinguishable. Hence, the signature must be valid and on a fresh message. Therefore, the reduction wins if the distinguisher loses and we get $\epsilon_{\mathsf{SIG}} \geq 1 - n_{ss} \cdot \gamma_{coll} - \epsilon_{\mathcal{D}}$ and $\epsilon_{\mathcal{D}} \geq 1 - n_{ss} \cdot \gamma_{coll} - \epsilon_{\mathsf{SIG}}$, which is not close to $\frac{1}{2}$.

Remark 4.3 (Domain separators). If signatures on the KEM public keys include a domain separator for ephemeral and semi-static use, then we do not have the loss of $n_{ss} \cdot \gamma_{coll}$ in the above theorem.

¹³via CHALLHONINIT(B, A, (create, (ssid = 1, e_{DH} = true, e_{KEM} = true)))

Remark 4.4 (Fake needs a private signature for initiator deniability). In Theorem 4.4 the distinguisher detects Fake since Fake does not learn any "private" signatures on ephemeral keys that the distinguisher does not know. In Theorems 4.5 to 4.7, the Fake algorithm obtains signatures on the semi-static DH key and on the KEM key from the auxiliary info aux (cp. the left-hand side of Figure 4), while the distinguisher does not, i.e. AUX $\notin O_D$. Alternatively, the Fake algorithm may also learn the signatures via the USER_{initiator,⊥} oracle oracle (cp. the right-hand side of Figure 4), resulting in partial deniability wrt. a USER_{initiator,⊥} oracle. Hence, these theorems also hold for AuxPrep = \perp if O_F additionally includes USER_{initiator,⊥}. In either case the distinguisher does not see the signature, i.e. the signature remains "private".

Bob can deny (*initiator deniability*) participation in a protocol run against a semi-honest Alice (*against semi-honest adversaries*, Theorem 4.5), against a malicious Alice that honestly generated her keys (*against malicious adversaries with honest keys*, Theorem 4.6), and against a malicious Alice (*against malicious adversaries*, Theorem 4.7). He does so by showing that Alice (1-out-of-2, i.e. anybody using Alice's secret keys, Theorems 4.5 and 4.6) or anybody (1-out-of- ∞ , Theorem 4.7) could have produced the transcript and session key herself, assuming she has previously observed a valid transcript with Bob as initiator (AuxPrep yielding a valid pre-key bundle per ssid and user and q_{CI} pre-key bundles per user with ephemeral KEM keys, Theorem 4.5 to 4.7). This holds against big brother distinguishers (*in the big brother model*, Theorem 4.5) who do not learn Alice's observations (aux unknown to the distinguisher, Theorems 4.5).

Theorem 4.5. The PQXDH protocol as shown in Figure 3 is $(n_p, n_{ss}, q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D, \epsilon_D)$ -deniable with respect to ({REGHON, CHALLHONINIT}, {SK}, {SKs})-oracles and auxiliary inputs sampled with AuxPrep yielding a valid pre-key bundle per ssid and user and q_{CI} pre-key bundles per user with ephemeral KEM keys, where n_p is the number of parties and n_{ss} the number of semi-static keys per party, $q_{CI} \in q_{O_A}$, $q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D$ are arbitrary, and $\epsilon_D = 0$.

Proof. We give the Fake algorithm in Figure 4 on the left-hand side and Figure 5.

The Fake algorithm simulates Bob's pre-key bundle by sampling a fresh ephemeral DH key (if needed) and learns a signature on the semi-static DH key and on the KEM key from the auxiliary info aux. Note that the Fake algorithm may need a signature for any semi-static KEM key or a signed ephemeral key for each of the q_{CI} queries. To process Alice's message, the Fake algorithm computes the DH shared secrets as for Theorem 4.1 (using Alice's freshly sampled ephemeral secret and the long-term secret obtained from the SK oracle). Furthermore, Fake extracts the KEM shared secret from the randomness that was previously used for the encapsulation. It computes the session key honestly as well as decrypts the AEAD ciphertext and aborts if decryption fails.

Since Fake produces a transcript and session key in the same way as Run and the distinguisher cannot detect reuse of the signature on the ephemeral KEM key, the distinguisher cannot have an advantage in winning the game, even if the distinguisher has access to all long-term and semi-static secret keys. \Box

For the next two theorems we need plaintext awareness.

Theorem 4.6. If the $(t_{ADH}, t_{\mathcal{E}DH}, t_{\mathcal{D}DH}, \epsilon_{\mathcal{D}DH})$ -Knowledge of 2DH (K2DH) Assumption holds, KEM is $(n_k, t_{\mathcal{C}}, t_{\mathcal{E}PA}, t_{\mathcal{D}PA}, \epsilon_{\mathcal{D}PA})$ -PA1 secure, and KDF is a random oracle, then the PQXDH protocol as shown in Figure 3 is $(n_p, n_{ss}, q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D, \epsilon_D)$ -deniable with respect to ({REGHON, INIT, CHALLINIT}, {SK}, \emptyset)-oracles and auxiliary inputs sampled with AuxPrep yielding a valid pre-key bundle per ssid and user and q_{CI} pre-key bundles per user with ephemeral KEM keys, where n_p is the number of parties and n_{ss} the number of semi-static keys per party, the number of keys for the PA assumption is $n_k = n_p \cdot n_{ss} + q_{CI}$, $q_{CI} \in q_{O_A}, q_{O_F}, q_{O_D}$ are arbitrary, $t_A \approx t_C \approx t_{\mathcal{E}DH}, t_D \approx t_{\mathcal{D}DH} \approx t_{\mathcal{D}PA}$, and $\epsilon_D \leq \epsilon_{\mathcal{D}DH} + \epsilon_{\mathcal{D}PA}$.

Proof. We give the Fake algorithm in Figure 4 on the left-hand side and Figure 5.

The Fake algorithm simulates Bob's pre-key bundle by sampling a fresh ephemeral DH key (if needed) and learns a signature on the semi-static DH key and on the KEM key from the auxiliary info aux. Note that the Fake algorithm may need a signature for any semi-static KEM key or a signed ephemeral key for each of the q_{CI} queries. When processing Alice's message, Fake uses Alice's long-term secret key from the SK oracle and the previously sampled ephemeral key to compute DH_1 and DH_4 . It extracts DH_2 , DH_3 from the adversary under the K2DH assumption. Furthermore, it uses the PA1 extractor in line 52 or 54 (depending on whether an ephemeral or semi-static KEM key is used) to learn the KEM shared secret. In contrast, Run computes DH_2 , DH_3 with Bob's DH secret keys and the KEM shared secret ss with KEM.Dec() and Bob's corresponding KEM secret key. If the K2DH extraction fails, Fake sets a random key as session key. Else, it computes the session key honestly. It decrypts the AEAD ciphertext and aborts if decryption fails.

In a nutshell, the Fake algorithm creates Bob's pre-key bundle virtually identical to Run.

We argue over a series of game hops that the distinguisher cannot distinguish whether Alice's message was processed by Run or Fake.

Game 0 = Run. The initial game is the original deniability game $\mathcal{G}_{\mathsf{PQXDH},\mathsf{AuxPrep},n_p,n_{ss}}^{\mathsf{O}_A,\mathsf{O}_F,\mathsf{O}_D-den}(\mathcal{A},\mathcal{D})$ with challenge b = 0, executing Run.

Game 1 (SK oracle). In this game Fake obtains Bob's long-term and semi-static secret key using the SK oracle, instead of receiving Bob's secret keys as input. This change is purely syntactical and not noticeable to an attacker or distinguisher. Hence, the winning probabilities for Games 0 and 1 are identical.

Game 2 (PA1 extractor). We replace the honest decapsulation of the KEM ciphertext with the PA extractor in lines 52 and 54. (In each invocation of Fake for processing Alice's message exactly one of these two lines is executed, depending on e_{DH} .)

We bound the advantage difference introduced by this step by the $(n_k, t_c, t_{\mathcal{E}PA}, t_{\mathcal{D}PA}, \epsilon_{\mathcal{D}PA})$ -PA1 assumption. We consider the reduction \mathcal{B} – running the adversary with its oracles, excluding calls to the DECAPS oracle – as ciphertext creator \mathcal{C} of the PA1 game. The reduction \mathcal{B} is started with a list of KEM keys $\vec{\mathsf{pk}}_{\mathsf{KEM}}$ and explicit randomness $r_{\mathcal{C}} = (r_{\mathcal{A}}, r_{RH}, r_{C})$.

The reduction samples $n_p + n_p \cdot n_{ss}$ DH keys and sorts them by user into \vec{pk} . Next, it prepares the auxiliary info aux, which consists of one pre-key bundle per user and semi-static key and q_{CI} pre-key bundles per user with ephemeral KEM keys. Note that for both sampling the DH keys and preparing the auxiliary info the reduction uses hard-wired randomness. (Without hard-wiring the randomness the PA1 extractor would require this randomness as input since we view the reduction as ciphertext creator.) Then, it starts the adversary \mathcal{A} on the auxiliary info aux and the list of public keys \vec{pk} with the randomness $r_{\mathcal{A}}$. The reduction provides the adversary with access to its oracles $O_A = \{\text{REGHON}, \text{INIT}, \text{CHALLINIT}\}$. To answer queries to the REGHON oracle the reduction uses the randomness r_{RH} . To answer queries to the CHALLINIT oracle the reduction always uses Fake as described for the previous game with the following three changes: First, instead of sampling a new ephemeral KEM key for Bob's pre-key bundle, it uses the next KEM public key from its input list. Second, the reduction replaces lines 52 and 54 (the decapsulation or extraction of the KEM shared secrets) with a call to its own DECAPS oracle of the PA1 game. Third, it uses the randomness r_C where needed. After the adversary terminates, the reduction outputs (aux, \vec{pk}, r, Q, K). The PA1-distinguisher \mathcal{D} then runs the deniability-distinguisher \mathcal{D} .

Hence, depending on the secret bit of the PA1 game, the reduction simulates either the previous game or the current game. We include the running time of the $O_A = \{\text{RegHon}, \text{INIT}, \text{CHALL}\}$ oracles in t_A , allowing for $t_A \approx t_C$.

The PA1 distinguisher \mathcal{D}' executes the deniability distinguisher \mathcal{D} , resulting in essentially the same run time $t_{\mathcal{D}} \approx t_{\mathcal{D}^{PA}}$, and directly returns the guess of the deniability distinguisher \mathcal{D} . If the deniability distinguisher \mathcal{D} succeeds, then the PA1 distinguisher \mathcal{D}' succeeds as well and can distinguish between the current and the previous game. By the PA1 assumption, the PA1 distinguisher succeeds with an advantage $\leq \epsilon_{\mathcal{D}^{PA}}$, which also limits the advantage of the deniability distinguisher.

Game 3 = Fake (K2DH assumption). We replace the honest computations of DH_2 , DH_3 with the K2DH extractor in line 46, finally resulting in Fake as described in Figure 5 on the right-hand side. We bound the advantage difference introduced by this step by the $(t_{A^{DH}}, t_{\mathcal{E}^{DH}}, t_{\mathcal{D}^{DH}}, \epsilon_{\mathcal{D}^{DH}})$ -K2DH assumption.

Formally, the adversary \mathcal{A} never outputs any value. Here, we apply the extractor $\mathcal{E}_{\mathcal{A}}^{DH}$ in line 46 to \mathcal{A} querying the CHALL oracle with the message $(A, epk_{\mathcal{A}}^{DH}, ct_{AE}, ct_{KEM})$ after receiving a message from Bob with semi-static key ssid.

If the extractor succeeds, then Fake has computed the session key in the same way as Run, allowing no advantage in distinguishing the current game from the previous game. If the extractor fails, then Fake sets a random session key. The runtime of the attacker in the K2DH assumption limits the runtime of the deniability attacker, i.e. $t_{\mathcal{A}} \approx t_{\mathcal{A}^{DH}}$. Furthermore, any K2DH distinguisher running in time $t_{\mathcal{D}^{DH}}$ succeeds in extracting DH_2 or DH_3 with a maximum probability of $\epsilon_{\mathcal{D}^{DH}}$. Hence, the deniability distinguisher with the same limit in runtime can distinguish the previous game from the current game with a probability of at most $\epsilon_{\mathcal{D}^{DH}}$.

Theorem 4.7. If the $(t_{ADH}, t_{\mathcal{E}DH}, t_{\mathcal{D}DH}, \epsilon_{\mathcal{D}DH})$ -Extended Knowledge of DH (EKDH) Assumption holds, KEM is $(n_k, t_{\mathcal{C}}, t_{\mathcal{E}PA}, t_{\mathcal{D}PA}, \epsilon_{\mathcal{D}PA})$ -PA1 secure, and KDF is a random oracle, then the PQXDH protocol as shown in Figure 3 is $(n_p, n_{ss}, q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D, \epsilon_D)$ -deniable with respect to ({REG, INIT, CHALLINIT}, \emptyset, \emptyset)-oracles and auxiliary inputs sampled with AuxPrep yielding a valid pre-key bundle per ssid and user and q_{CI} pre-key bundles per user with ephemeral KEM keys, where n_p is the number of parties and n_{ss} the number of semi-static keys per party, $q_{CI} \in q_{O_A}, q_{O_A}, q_{O_F}, q_{O_D}$ are arbitrary, and $t_A \approx t_{\mathcal{C}} \approx 3 \cdot t_{\mathcal{E}DH},$ $t_{\mathcal{D}} \approx t_{\mathcal{D}DH} \approx t_{\mathcal{D}PA}$, and $\epsilon_{\mathcal{D}} \leq \epsilon_{\mathcal{D}DH} + \epsilon_{\mathcal{D}PA}$.

Proof. We give the Fake algorithm in Figure 4 on the right-hand side and Figure 6.

The Fake algorithm reuses the pre-key bundle obtained from the USER_{initiator,⊥} oracle but replaces the ephemeral DH key with a freshly sampled key (if mandated by e_{DH}). Furthermore, to process Alice's message, the Fake algorithm requires several techniques: The Fake algorithm relies on the extractor from the EKDH assumption to learn DH_1 through DH_3 (lines 67, 68, and 69) and uses Bob's ephemeral key (from its previous run) to compute DH_4 . Additionally, it utilizes the plaintext extractor for the KEM to learn the KEM shared secret (line 76 or 78). If any of the three extractions under the EKDH assumption fail, Fake sets a random key as session key. Else, it computes the session key honestly. It decrypts the AEAD ciphertext and aborts if decryption fails.

The Fake algorithm creates Bob's pre-key bundle virtually identical to Run. We argue via a series of game hops that the distinguisher cannot tell whether Alice's message was processed by Run or Fake.

Game 0 = Run. The initial game is the original deniability game $\mathcal{G}_{\mathsf{PQXDH},\mathsf{AuxPrep},n_p,n_{ss}}^{\mathsf{O}_A,\mathsf{O}_F,\mathsf{O}_D}(\mathcal{A},\mathcal{D})$ with challenge b = 0, executing Run.

Game 1 (PA1 extractor). We replace the honest decapsulation of the KEM ciphertext with the PA extractor in lines 76 and 78. (In each invocation of Fake for processing Alice's message, exactly one of these two lines is executed, depending on e_{DH} .)

We bound the advantage difference introduced by this step by the $(n_k, t_c, t_{\mathcal{E}^{PA}}, t_{\mathcal{D}^{PA}}, \epsilon_{\mathcal{D}^{PA}})$ -PA1 assumption. We consider the reduction \mathcal{B} —running the adversary with its oracles, excluding calls to the DECAPS oracle—as ciphertext creator C of the PA1 game. The reduction \mathcal{B} is started with a list of KEM keys $\vec{\mathsf{pk}}_{\mathsf{KEM}}$ and explicit randomness $r_{\mathcal{C}} = (r_{\mathcal{A}}, r_R, r_C)$.

The reduction samples $n_p \cdot n_{ss}$ DH keys and sorts them by user into $\mathbf{p}\mathbf{k}$. Next, it prepares the auxiliary info aux, which consists of one pre-key bundle per user and semi-static key. Note that for both sampling the DH keys and preparing the auxiliary info, the reduction uses hard-wired randomness. (Without hardwiring the randomness, the PA1 extractor would require this randomness as input since we view the reduction as ciphertext creator.) Then, it starts the adversary \mathcal{A} on the auxiliary info aux and the list of public keys $\mathbf{p}\mathbf{k}$ with the randomness $r_{\mathcal{A}}$. The reduction provides the adversary with access to its oracles $O_A = \{\text{REG, INIT, CHALL}\}$. To answer queries to the REG, oracle the reduction uses the randomness r_R . To answer queries to the CHALL oracle, the reduction always uses Fake as described for the current game with the following three changes: First, instead of sampling a new ephemeral KEM key for Bob's pre-key bundle, it uses the next KEM public key from its input list. Second, the reduction replaces lines 76 and 78 (the decapsulation or extraction of the KEM shared secrets) with a call to its own DECAPS oracle of the PA1 game. Third, it uses the randomness r_C where needed. After the adversary terminates, the reduction outputs (aux, $\mathbf{p}\mathbf{k}, r, Q, K$). The PA1-distinguisher \mathcal{D}' then runs the deniability-distinguisher \mathcal{D}

Hence, depending on the secret bit of the PA1 game, the reduction simulates either the previous game or the current game. We include the running time of the $O_A = \{\text{REG}, \text{INIT}, \text{CHALL}\}$ oracles in t_A , allowing for $t_A \approx t_C$.

The PA1 distinguisher \mathcal{D}' executes the deniability distinguisher \mathcal{D} , resulting in essentially the same run time $t_{\mathcal{D}} \approx t_{\mathcal{D}^{PA}}$, and directly returns the guess of the deniability distinguisher \mathcal{D} . If the deniability distinguisher \mathcal{D} succeeds, then the PA1 distinguisher \mathcal{D}' succeeds as well and can distinguish between the current and the previous game. By the PA1 assumption, the PA1 distinguisher succeeds with an advantage $\leq \epsilon_{\mathcal{D}^{PA}}$, which also limits the advantage of the deniability distinguisher.

Game 2 = Fake (EKDH assumption). We replace the honest computations of DH_1, DH_2, DH_3 with the EKDH extractor in lines 67, 68, and 69, finally resulting in Fake as described above. Following Theorem 4.3, we bound the advantage difference introduced by this step by the $(t_{ADH}, t_{EDH}, t_{DDH}, \epsilon_{DDH})$ -EKDH assumption.

If the extractor succeeds (in all three cases), then Fake has computed the session key in the same way as Run, and, hence, indistinguishably. If the extractor fails (in any of the three cases), then the distinguisher succeeds in extracting one of the shared DH secrets with a maximum probability of ϵ_{DDH} while running in time t_{DDH} . Hence, the distinguisher can distinguish the real session key from the simulated session key, i.e. the previous game from the current game, with a probability of at most ϵ_{DDH} . As for Theorem 4.3, the second and third extractor call refer to the same output of the adversary, i.e. they are not independent. Conservatively, we bound with ϵ_{DDH} and not $(\epsilon_{DDH})^3$.

For responder deniability we obtain identical guarantees for X3DH and PQXDH. Alice can deny (responder deniability) participation in a protocol run against a malicious Bob that honestly generated his keys (against malicious adversaries with honest keys, Theorem 4.8), and against a malicious Bob (against malicious adversaries, Theorem 4.9). She does so by showing that Bob (1-out-of-2, i.e. using Bob's secret keys, Theorem 4.8), or anybody (1-out-of- ∞ , Theorem 4.9) could have produced the transcript and session key himself. This holds even against a big brother distinguisher (in the big brother model, Theorem 4.8) and does not rely on eavesdropped info (AuxPrep = \bot , Theorems 4.8 and 4.9).

Theorem 4.8. Both the X3DH and PQXDH protocols as shown in Figure 3 are $(n_p, n_{ss}, q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D, \epsilon_D)$ -deniable with respect to ({REGHON, INIT, CHALLRESP}, {SK}, {SKs})-oracles and auxiliary inputs sampled with AuxPrep = \perp , where n_p is the number of parties and n_{ss} the number of semi-static keys per party, $q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D$ are arbitrary, and $\epsilon_D = 0$.

 $\mathsf{Fake}(\vec{\mathsf{pk}}, \pi, m = (B, \mathsf{ssid}, \sigma_{\mathsf{DH}}^{\mathsf{ssid}}, epk_B, \sigma_{\mathsf{KEM}}), \mu, \mathsf{aux}, r)$: $\mathsf{Fake}(\vec{\mathsf{pk}}, \pi, m = (B, \mathsf{ssid}, \sigma_{\mathsf{DH}}^{\mathsf{ssid}}, epk_B, \sigma_{\mathsf{KEM}}), \mu, \mathsf{aux}, r)$: 116 π .pid $\leftarrow B$ 88 π .pid $\leftarrow B$ 117 $(epk_A^{\mathsf{DH}}, esk_A^{\mathsf{DH}}) \leftarrow \mathsf{SDH}.\mathsf{KGen}()$ 89 $sssk_B^{ssid} \leftarrow extract from SK(\pi)$ 118 $(sspk_B^{\mathsf{DH}}, sspk_B^{\mathsf{KEM}}) \leftarrow sspk_B^{\mathsf{ssid}}$ 90 $(sssk_B^{\mathsf{DH}}, sssk_B^{\mathsf{KEM}}) \leftarrow sssk_B^{\mathsf{ssid}}$ 91 $(epk_A^{\mathsf{DH}}, esk_A^{\mathsf{DH}}) \leftarrow$ DH.KGen() 119 $(epk_B^{\mathsf{DH}}, epk_B^{\mathsf{KEM}}) \leftarrow epk_B$ if SIG.Vf $(ltpk_B^{SIG}, sspk_B^{DH}, \sigma_{DH}) = false$ 92 $(sspk_B^{\mathsf{DH}}, sspk_B^{\mathsf{KEM}}) \leftarrow sspk_B^{\mathsf{ssid}}$ 120 return $(\pi, \epsilon, \epsilon)$ 121 $(epk_B^{\mathsf{DH}}, epk_B^{\mathsf{KEM}}) \gets epk_B$ 122 $DH_1 \leftarrow \mathcal{E}_{\mathcal{A}}^{DH}(ltpk_A^{\mathsf{DH}},\mathsf{aux},r)$ 94 if SIG.Vf $(ltpk_B^{SIG}, sspk_B^{DH}, \sigma_{DH}) =$ false $DH_2 \leftarrow \mathsf{DH}(esk_A^{\mathsf{DH}}, ltpk_B^{\mathsf{DH}})$ 123 95 return $(\pi, \epsilon, \epsilon)$ 96 $DH_1 \leftarrow \mathsf{DH}(ltpk_A^{\mathsf{DH}}, sssk_B^{\mathsf{DH}})$ 124 $DH_3 \leftarrow \mathsf{DH}(esk_A^{\mathsf{DH}}, sspk_B^{\mathsf{DH}})$ if $epk_B^{\mathsf{DH}} \neq \bot$ //ephemeral DH key present 97 $DH_2 \leftarrow \mathsf{DH}(esk_A^{\mathsf{DH}}, ltpk_B^{\mathsf{DH}})$ 125 $DH_4 \leftarrow \mathsf{DH}(esk_A^{\mathsf{DH}}, epk_B^{\mathsf{DH}})$ 126 $DH_3 \leftarrow \mathsf{DH}(esk_A^{\mathsf{DH}}, sspk_B^{\mathsf{DH}})$ 98 127 else if $epk_B^{\mathsf{DH}} \neq \perp$ //ephemeral DH key present 99 128 $DH_4 \leftarrow \epsilon$ $DH_4 \leftarrow \mathsf{DH}(esk_A^{\mathsf{DH}}, epk_B^{\mathsf{DH}})$ 100 if $epk_B^{\mathsf{KEM}} \neq \bot$ //ephemeral KEM key present 129 101 elseif SIG.Vf $(ltpk_B^{SIG}, epk_B^{KEM}, \sigma_{KEM}) = false$ 130 $DH_4 \leftarrow \epsilon$ 102 return $(\pi_A, \epsilon, \epsilon)$ if $epk_B^{\mathsf{KEM}} \neq \bot$ //ephemeral KEM key present 131 103 $(ct_{\mathsf{KEM}}, ss) \leftarrow \mathsf{KEM}.\mathsf{Enc}(epk_B^{\mathsf{KEM}})$ if SIG.Vf $(ltpk_{B}^{SIG}, epk_{B}^{KEM}, \sigma_{KEM}) = false$ 132 104 else //no ephemeral KEM key present 133 105 return $(\pi_A, \epsilon, \epsilon)$ if SIG.Vf $(ltpk_B^{SIG}, sspk_B^{KEM}, \sigma_{KEM}) = false$ 134 $(ct_{\mathsf{KEM}}, ss) \leftarrow \mathsf{KEM}.\mathsf{Enc}(epk_B^{\mathsf{KEM}})$ 106 135 return $(\pi_A, \epsilon, \epsilon)$ else //no ephemeral KEM key present 107 $(ct_{\mathsf{KEM}}, ss) \gets \mathsf{KEM}.\mathsf{Enc}(sspk_{B}^{\mathsf{KEM}})$ if SIG.Vf $(ltpk_B^{SIG}, sspk_B^{KEM}, \sigma_{KEM}) = false$ 136 108 137 if $DH_1 = \bot$ 109 return $(\pi_A, \epsilon, \epsilon)$ $\pi.\mathsf{K} \leftarrow \{0,1\}^{256}$ 138 $(ct_{\mathsf{KEM}}, ss) \leftarrow \mathsf{KEM}.\mathsf{Enc}(sspk_B^{\mathsf{KEM}})$ 110 139 else 111 $\mathsf{ms} \leftarrow DH_1 \| DH_2 \| DH_3 \| DH_4 \| ss$ $\mathsf{ms} \leftarrow DH_1 \| DH_2 \| DH_3 \| DH_4 \| ss$ 140 112 $\pi.K \leftarrow KDF(ms)$ 141 $\pi.K \leftarrow KDF(ms)$ 113 $AD \leftarrow ltpk_A \| ltpk_B$ 142 $AD \gets ltpk_A \| ltpk_B$ 114 $ct_{\mathsf{AE}} \leftarrow \mathsf{Enc}(\pi_A.\mathsf{K}, AD, \mu)$ $ct_{\mathsf{AE}} \leftarrow \mathsf{Enc}(\pi_A.\mathsf{K}, AD, \mu)$ 143 115 return $(\pi, (A, epk_A^{DH}, ct_{AE}, ct_{KEM}), \epsilon)$ return $(\pi, (A, epk_A^{\mathsf{DH}}, ct_{\mathsf{AE}}, ct_{\mathsf{KEM}}), \epsilon)$ 144

Figure 7: The Fake algorithms for producing Alice's message (responder deniability) in Theorem 4.8 (on the left-hand side) and Theorem 4.9 (on the right-hand side).

Proof. We give the Fake algorithm in Figure 7 on the left-hand side. The Fake algorithm considers the KEM parts (obtaining the KEM public key, verifying the signature on the KEM public key, encapsulating against the KEM public key, and sending a KEM ciphertext) only for PQXDH.

The Fake algorithm uses the freshly sampled ephemeral key esk_A^{DH} to compute DH_1 and Bob's longterm secret keys from the SK oracle $ltsk_B^{DH}$ to compute DH_2 , DH_3 , DH_4 . It learns the KEM shared secret from encapsulating against Bob's public key. If either signature (on the semi-static DH key or on the KEM key) does not verify, it aborts. Using those secrets, Fake computes the session key and encrypts the user message with the AEAD scheme under the session key.

Since Fake produces a transcript and session key in the same way as Run, the distinguisher cannot have an advantage in winning his game, even if the distinguisher has access to all long-term and semi-static secret keys.

Theorem 4.9 (Responder deniability against malicious adversaries of PQXDH). If the $(t_{A^{DH}}, t_{\mathcal{E}^{DH}}, t_{\mathcal{D}^{DH}}, \epsilon_{\mathcal{D}^{DH}})$ -Extended Knowledge of DH (EKDH) Assumption holds and KDF is a random oracle, then both the X3DH and the PQXDH protocols as shown in Figure 3 are $(n_p, n_{ss}, q_{O_A}, q_{O_F}, q_{O_D}, t_A, t_D, \epsilon_D)$ -deniable with respect to ({REG, INIT, CHALLRESP}, \emptyset, \emptyset)-oracles and auxiliary inputs sampled with AuxPrep = \bot , where n_p is the number of parties and n_{ss} the number of semi-static keys per party, $q_{O_A}, q_{O_F}, q_{O_D}$ are arbitrary, $t_A \approx t_{A^{DH}}, t_D \approx t_{D^{DH}}$, and $\epsilon_D \leq \epsilon_{D^{DH}}$.

Proof. We give the Fake algorithm in Figure 7 on the right-hand side. The Fake algorithm considers the

KEM parts (obtaining the KEM public key, verifying the signature on the KEM public key, encapsulating against the KEM public key, and sending a KEM ciphertext) only for PQXDH.

The Fake algorithm uses its freshly sampled ephemeral secret key esk_A^{DH} to compute DH_2 through DH_4 . It extracts DH_1 from the adversary under the EKDH assumption. As before, we interpret the adversary's query to the REG oracle as its output and apply the extractor \mathcal{E}_A^{DH} accordingly. It learns the KEM shared secret from encapsulating against Bob's public key. If either signature (on the semi-static DH key or on the KEM key) does not verify, it aborts. If the extraction fails, Fake sets a random key as session key. Otherwise, it computes the session key honestly. It encrypts the user message with the AEAD scheme under the session key.

The Fake algorithm produces the transcript in the same way as Run and, hence, indistinguishably. If the extractor succeeds, then Fake has computed the session key in the same way as Run, and, hence, indistinguishably. If the extractor fails, then the distinguisher succeeds in extracting the shared DH secret with a maximum probability of $\epsilon_{\mathcal{D}^{DH}}$. Hence, the distinguisher can distinguish the real session key from the simulated session key with a probability of at most $\epsilon_{\mathcal{D}^{DH}}$.

Remark 4.5 (Deniability against malicious adversaries in the big brother model). In Remark 4.2 we have pointed out challenges for showing deniability against malicious adversaries in the big brother model for X3DH. They directly apply to PQXDH as well.

5 Conclusion and Future Work

Table 2 summarizes that in comparison to X3DH, PQXDH keeps the deniability guarantees for the responder but not for the initiator: simulating a transcript requires a signature that the distinguisher never sees. It remains an open question if deniability against malicious adversaries holds in the big brother model for both X3DH and PQXDH.

Observe that the AEAD ciphertext influences the assumptions needed to show deniability: Without the AEAD ciphertext, we can show that the DH assumption (and for PQXDH also plaintext awareness of the KEM) is necessary to show deniability against malicious adversaries.¹⁴ When including the AEAD ciphertext, we cannot show the necessity of these assumptions anymore. Also, the KEM used in Signal's implementation, Kyber, does not fulfill plaintext awareness, as pointed out in [KS23].

One can consider a key awareness assumption (similar to the case of SIGMA [DGK06]) that allows extracting the AEAD key from a malicious Alice. However, on the one hand, the Fake algorithm cannot tell if an extracted AEAD key is the real session key. On the other hand, the distinguisher can tell them apart (since it knows what the adversary does). Observing the adversary's queries to the random oracle does not help either since Fake cannot compute the master secret ms. (Nor could the challenger verify if the adversary queries the random oracle on the actual master secret.) Hence, it appears this proof strategy does not work.

Our analysis treats Bob's long-term DH key and signing key separately, which is not the case in practice. A formal proof for the protocol taking this combined primitive into account is still open.

We see it as an interesting problem to examine if our deniability model is directly applicable to messaging protocols or if any modifications are needed. Cremers and Zhao [CZ24] have extended the deniability notion of [BFG⁺22] for key exchange to messaging protocols. Furthermore, one can extend our deniability model to group chats, see e.g. [HW21] and the efforts around the MLS standardization process. An extended deniability model for secure messaging can then be used to analyze the deniability guarantees of the whole Signal protocol, i.e. including the Double Ratchet.

¹⁴A deniability adversary can use any ciphertext creator to obtain a KEM ciphertext, add a freshly sampled DH public key, and send this as Alice's message to the challenge oracle. We can build a PA extractor by observing the RO queries of the Fake algorithm to learn the KEM shared secret, which must be indistinguishable due to deniability.

Acknowledgements

We thank the anonymous reviewers for their constructive feedback, and Marc Fischlin, Ehren Kret, Trevor Perrin, Olga Sanina, and Rolfe Schmidt for helpful discussions and feedback.

R.F. was supported by the German Federal Ministry of Education and Research and the Hessian Ministry of Higher Education, Research, Science and the Arts within their joint support of the National Research Center for Applied Cybersecurity ATHENE. C.J. was (partially) funded by the Deutsche Forschungsgemeinschaft (DFG) – SFB 1119 – 236615297.

References

- [BFG⁺20] Jacqueline Brendel, Marc Fischlin, Felix Günther, Christian Janson, and Douglas Stebila. Towards post-quantum security for Signal's X3DH handshake. In Orr Dunkelman, Michael J. Jacobson Jr., and Colin O'Flynn, editors, SAC 2020: 27th Annual International Workshop on Selected Areas in Cryptography, volume 12804 of Lecture Notes in Computer Science, pages 404–430, Halifax, NS, Canada (Virtual Event), October 21-23, 2020. Springer, Heidelberg, Germany. (Cited on page 4.)
- [BFG⁺21] Jacqueline Brendel, Rune Fiedler, Felix Günther, Christian Janson, and Douglas Stebila. Post-quantum asynchronous deniable key exchange and the Signal handshake. Cryptology ePrint Archive, Report 2021/769, 2021. https://eprint.iacr.org/2021/769. (Cited on pages 2, 4, 8, 36, 38, and 39.)
- [BFG⁺22] Jacqueline Brendel, Rune Fiedler, Felix Günther, Christian Janson, and Douglas Stebila. Post-quantum asynchronous deniable key exchange and the Signal handshake. In Goichiro Hanaoka, Junji Shikata, and Yohei Watanabe, editors, *PKC 2022: 25th International Conference on Theory and Practice of Public Key Cryptography, Part II*, volume 13178 of *Lecture Notes in Computer Science*, pages 3–34, Virtual Event, March 8–11, 2022. Springer, Heidelberg, Germany. (Cited on pages 2, 4, 8, 14, 15, 16, 29, 36, 38, and 39.)
- [BGB04] Nikita Borisov, Ian Goldberg, and Eric A. Brewer. Off-the-record communication, or, why not to use PGP. In Vijay Atluri, Paul F. Syverson, and Sabrina De Capitani di Vimercati, editors, Proceedings of the 2004 ACM Workshop on Privacy in the Electronic Society, WPES 2004, Washington, DC, USA, October 28, 2004, pages 77–84. ACM, 2004. (Cited on page 3.)
- [BJK23] Karthikean Barghavan, Charlie Jacomme, and Franziskus Kiefer. Formal analysis of the PQXDH protocol, 2023. https://github.com/Inria-Prosecco/pqxdh-analysis. (Cited on page 4.)
- [BJKS23] Karthikean Barghavan, Charlie Jacomme, Franziskus Kiefer, and Rolfe Schmidt. An analysis of Signal's PQXDH. https://cryspen.com/post/pqxdh/, October 2023. (Cited on pages 4 and 15.)
- [BMP03] Colin Boyd, Wenbo Mao, and Kenneth G. Paterson. Deniable authenticated key establishment for internet protocols. In *Security Protocols Workshop*, pages 255–271, 2003. (Cited on page 3.)
- [BMP04] Colin Boyd, Wenbo Mao, and Kenneth G. Paterson. Key agreement using statically keyed authenticators. In Markus Jakobsson, Moti Yung, and Jianying Zhou, editors, ACNS 04: 2nd International Conference on Applied Cryptography and Network Security, volume 3089 of Lecture Notes in Computer Science, pages 248–262, Yellow Mountain, China, June 8–11, 2004. Springer, Heidelberg, Germany. (Cited on page 3.)

- [BMS20] Colin Boyd, Anish Mathuria, and Douglas Stebila. Protocols for Authentication and Key Establishment, Second Edition. Information Security and Cryptography. Springer, 2020. (Cited on page 3.)
- [BP04] Mihir Bellare and Adriana Palacio. Towards plaintext-aware public-key encryption without random oracles. In Pil Joong Lee, editor, Advances in Cryptology – ASIACRYPT 2004, volume 3329 of Lecture Notes in Computer Science, pages 48–62, Jeju Island, Korea, December 5–9, 2004. Springer, Heidelberg, Germany. (Cited on page 7.)
- [BR95] Mihir Bellare and Phillip Rogaway. Optimal asymmetric encryption. In Alfredo De Santis, editor, Advances in Cryptology – EUROCRYPT'94, volume 950 of Lecture Notes in Computer Science, pages 92–111, Perugia, Italy, May 9–12, 1995. Springer, Heidelberg, Germany. (Cited on page 7.)
- [CCD⁺17] Katriel Cohn-Gordon, Cas J. F. Cremers, Benjamin Dowling, Luke Garratt, and Douglas Stebila. A formal security analysis of the Signal messaging protocol. In *IEEE European* Symposium on Security and Privacy, EuroS&P 2017, pages 451–466, 2017. (Cited on pages 4 and 16.)
- [CCHD23] Daniel Collins, Simone Colombo, and Loïs Huguenin-Dumittan. Real world deniability in messaging. Cryptology ePrint Archive, Paper 2023/403, 2023. (Cited on page 2.)
- [CD23] Wouter Castryck and Thomas Decru. An efficient key recovery attack on SIDH. In Carmit Hazay and Martijn Stam, editors, Advances in Cryptology – EUROCRYPT 2023, Part V, volume 14008 of Lecture Notes in Computer Science, pages 423–447, Lyon, France, April 23– 27, 2023. Springer, Heidelberg, Germany. (Cited on page 4.)
- [CF11] Cas Cremers and Michele Feltz. One-round strongly secure key exchange with perfect forward secrecy and deniability. Cryptology ePrint Archive, Report 2011/300, 2011. https://eprint. iacr.org/2011/300. (Cited on pages 2, 4, 13, 36, and 37.)
- [CHDN⁺24] Daniel Collins, Loïs Huguenin-Dumittan, Ngoc Khanh Nguyen, Nicolas Rolin, and Serge Vaudenay. K-Waay: Fast and deniable Post-Quantum X3DH without ring signatures. In 33rd USENIX Security Symposium (USENIX Security 24), Philadelphia, PA, August 2024. USENIX Association. (Cited on pages 4, 36, and 39.)
- [CK02] Ran Canetti and Hugo Krawczyk. Security analysis of IKE's signature-based key-exchange protocol. In Moti Yung, editor, Advances in Cryptology – CRYPTO 2002, volume 2442 of Lecture Notes in Computer Science, pages 143–161, Santa Barbara, CA, USA, August 18– 22, 2002. Springer, Heidelberg, Germany. https://eprint.iacr.org/2002/120/. (Cited on page 14.)
- [CZ24] Cas Cremers and Mang Zhao. Secure messaging with strong compromise resilience, temporal privacy, and immediate decryption. In *IEEE Symposium on Security and Privacy*, 2024. (Cited on pages 4 and 29.)
- [Den06] Alexander W. Dent. The Cramer-Shoup encryption scheme is plaintext aware in the standard model. In Serge Vaudenay, editor, Advances in Cryptology – EUROCRYPT 2006, volume 4004 of Lecture Notes in Computer Science, pages 289–307, St. Petersburg, Russia, May 28 – June 1, 2006. Springer, Heidelberg, Germany. (Cited on page 7.)

- [DFG⁺13] Özgür Dagdelen, Marc Fischlin, Tommaso Gagliardoni, Giorgia Azzurra Marson, Arno Mittelbach, and Cristina Onete. A cryptographic analysis of OPACITY - (extended abstract). In Jason Crampton, Sushil Jajodia, and Keith Mayes, editors, ESORICS 2013: 18th European Symposium on Research in Computer Security, volume 8134 of Lecture Notes in Computer Science, pages 345–362, Egham, UK, September 9–13, 2013. Springer, Heidelberg, Germany. (Cited on pages 2, 4, 36, and 38.)
- [DG05] Mario Di Raimondo and Rosario Gennaro. New approaches for deniable authentication. In Vijayalakshmi Atluri, Catherine Meadows, and Ari Juels, editors, ACM CCS 2005: 12th Conference on Computer and Communications Security, pages 112–121, Alexandria, Virginia, USA, November 7–11, 2005. ACM Press. (Cited on page 3.)
- [DG22] Samuel Dobson and Steven D. Galbraith. Post-quantum signal key agreement from SIDH. In Jung Hee Cheon and Thomas Johansson, editors, *Post-Quantum Cryptography - 13th International Workshop*, *PQCrypto 2022*, *Virtual Event, September 28-30, 2022*, *Proceedings*, volume 13512 of *Lecture Notes in Computer Science*, pages 422–450. Springer, 2022. (Cited on page 4.)
- [DGK05] Mario Di Raimondo, Rosario Gennaro, and Hugo Krawczyk. Secure off-the-record messaging. In Vijay Atluri, Sabrina De Capitani di Vimercati, and Roger Dingledine, editors, Proceedings of the 2005 ACM Workshop on Privacy in the Electronic Society, WPES 2005, Alexandria, VA, USA, November 7, 2005, pages 81–89. ACM, 2005. (Cited on page 3.)
- [DGK06] Mario Di Raimondo, Rosario Gennaro, and Hugo Krawczyk. Deniable authentication and key exchange. In Ari Juels, Rebecca N. Wright, and Sabrina De Capitani di Vimercati, editors, ACM CCS 2006: 13th Conference on Computer and Communications Security, pages 400– 409, Alexandria, Virginia, USA, October 30 – November 3, 2006. ACM Press. (Cited on pages 2, 3, 4, 12, 14, 16, 29, 35, 36, 37, and 38.)
- [DKSW09] Yevgeniy Dodis, Jonathan Katz, Adam Smith, and Shabsi Walfish. Composability and online deniability of authentication. In Omer Reingold, editor, TCC 2009: 6th Theory of Cryptography Conference, volume 5444 of Lecture Notes in Computer Science, pages 146– 162. Springer, Heidelberg, Germany, March 15–17, 2009. (Cited on pages 3 and 4.)
- [DNS98] Cynthia Dwork, Moni Naor, and Amit Sahai. Concurrent zero-knowledge. In 30th Annual ACM Symposium on Theory of Computing, pages 409–418, Dallas, TX, USA, May 23–26, 1998. ACM Press. (Cited on page 3.)
- [FG24] Rune Fiedler and Felix Günther. Security analysis of Signal's PQXDH handshake. Cryptology ePrint Archive, Paper 2024/702, 2024. https://eprint.iacr.org/2024/702. (Cited on page 4.)
- [HKK⁺02] Dan Harkins, Charlie Kaufman, Tero Kivinen, Stephen Kent, and Radia Perlman. Design rationale for ikev2. Internet-Draft draft-ietf-ipsec-ikev2-rationale-00.txt, IETF Secretariat, Aug 2002. (Cited on page 3.)
- [HKKP22] Keitaro Hashimoto, Shuichi Katsumata, Kris Kwiatkowski, and Thomas Prest. An efficient and generic construction for Signal's handshake (X3DH): Post-quantum, state leakage secure, and deniable. *Journal of Cryptology*, 35(3):17, July 2022. (Cited on pages 2, 4, 7, 8, 14, 22, 36, and 38.)

- [HLLC11] Lein Harn, Chia-Yin Lee, Changlu Lin, and Chin-Chen Chang. Fully deniable message authentication protocols preserving confidentiality. *Comput. J.*, 54(10):1688–1699, 2011. (Cited on page 14.)
- [HW21] Andreas Hülsing and Fiona Johanna Weber. Epochal signatures for deniable group chats. In 2021 IEEE Symposium on Security and Privacy, pages 1677–1695, San Francisco, CA, USA, May 24–27, 2021. IEEE Computer Society Press. (Cited on page 29.)
- [JCL⁺22] Shaoquan Jiang, Yeow Meng Chee, San Ling, Huaxiong Wang, and Chaoping Xing. A new framework for deniable secure key exchange. *Inf. Comput.*, 285(Part):104866, 2022. (Cited on pages 2, 3, 4, 13, 36, and 39.)
- [Jia14] Shaoquan Jiang. Timed encryption with application to deniable key exchange. *Theor. Comput. Sci.*, 560:172–189, 2014. (Cited on pages 2, 3, 4, 13, 36, and 39.)
- [JS08] Shaoquan Jiang and Reihaneh Safavi-Naini. An efficient deniable key exchange protocol (extended abstract). In Gene Tsudik, editor, FC 2008: 12th International Conference on Financial Cryptography and Data Security, volume 5143 of Lecture Notes in Computer Science, pages 47–52, Cozumel, Mexico, January 28–31, 2008. Springer, Heidelberg, Germany. (Cited on pages 2, 3, 4, 13, 36, and 39.)
- [JW10] Shaoquan Jiang and Huaxiong Wang. Plaintext-awareness of hybrid encryption. In Josef Pieprzyk, editor, Topics in Cryptology – CT-RSA 2010, volume 5985 of Lecture Notes in Computer Science, pages 57–72, San Francisco, CA, USA, March 1–5, 2010. Springer, Heidelberg, Germany. (Cited on page 7.)
- [Kat03] Jonathan Katz. Efficient and non-malleable proofs of plaintext knowledge and applications. In Eli Biham, editor, Advances in Cryptology – EUROCRYPT 2003, volume 2656 of Lecture Notes in Computer Science, pages 211–228, Warsaw, Poland, May 4–8, 2003. Springer, Heidelberg, Germany. (Cited on page 3.)
- [Kra96] Hugo Krawczyk. SKEME: a versatile secure key exchange mechanism for internet. In James T. Ellis, B. Clifford Neuman, and David M. Balenson, editors, *ISOC Network and Distributed System Security Symposium NDSS'96*, pages 114–127, San Diego, CA, USA, February 22–23, 1996. IEEE Computer Society. (Cited on page 3.)
- [KS23] Ehren Kret and Rolfe Schmidt. The PQXDH key agreement protocol, September 2023. https://signal.org/docs/specifications/pqxdh/. (Cited on pages 1, 2, 4, 15, 22, and 29.)
- [LLPM07] Meng-Hui Lim, Sanggon Lee, Youngho Park, and Sangjae Moon. Secure deniable authenticated key establishment for internet protocols. Cryptology ePrint Archive, Report 2007/163, 2007. https://eprint.iacr.org/2007/163. (Cited on page 3.)
- [MMP⁺23] Luciano Maino, Chloe Martindale, Lorenz Panny, Giacomo Pope, and Benjamin Wesolowski. A direct key recovery attack on SIDH. In Carmit Hazay and Martijn Stam, editors, Advances in Cryptology – EUROCRYPT 2023, Part V, volume 14008 of Lecture Notes in Computer Science, pages 448–471, Lyon, France, April 23–27, 2023. Springer, Heidelberg, Germany. (Cited on page 4.)
- [MP02] Wenbo Mao and Kenneth G. Paterson. On the plausible deniability feature of internet protocols. unpublished, https://web.archive.org/web/20220818192033/http://www.isg. rhul.ac.uk/~kp/IKE.ps, 2002. (Cited on pages 3 and 15.)

- [MP16a] Moxie Marlinspike and Trevor Perrin. The double ratchet algorithm, November 2016. https: //www.signal.org/docs/specifications/doubleratchet/. (Cited on page 1.)
- [MP16b] Moxie Marlinspike and Trevor Perrin. The X3DH key agreement protocol, November 2016. https://signal.org/docs/specifications/x3dh/. (Cited on pages 1, 2, 4, 15, and 16.)
- [Nao02] Moni Naor. Deniable ring authentication. In Moti Yung, editor, Advances in Cryptology CRYPTO 2002, volume 2442 of Lecture Notes in Computer Science, pages 481–498, Santa Barbara, CA, USA, August 18–22, 2002. Springer, Heidelberg, Germany. (Cited on page 15.)
- [OTR20] OTR team. OTR version 4, 2020. https://github.com/otrv4/otrv4/blob/master/otrv4. md. (Cited on page 2.)
- [Pas03] Rafael Pass. On deniability in the common reference string and random oracle model. In Dan Boneh, editor, Advances in Cryptology CRYPTO 2003, volume 2729 of Lecture Notes in Computer Science, pages 316–337, Santa Barbara, CA, USA, August 17–21, 2003. Springer, Heidelberg, Germany. (Cited on pages 3, 4, and 16.)
- [Per16] Trevor Perrin. The XEdDSA and VXEdDSA signature schemes, October 2016. https: //signal.org/docs/specifications/xeddsa/. (Cited on page 15.)
- [RMA⁺23] Nathan Reitinger, Nathan Malkin, Omer Akgul, Michelle L. Mazurek, and Ian Miers. Is cryptographic deniability Sufficient f non-expert perceptions of deniability in secure messaging. In 2023 IEEE Symposium on Security and Privacy, pages 274–292, San Francisco, CA, USA, May 21–25, 2023. IEEE Computer Society Press. (Cited on page 5.)
- [Rob23] Damien Robert. Breaking SIDH in polynomial time. In Carmit Hazay and Martijn Stam, editors, Advances in Cryptology – EUROCRYPT 2023, Part V, volume 14008 of Lecture Notes in Computer Science, pages 472–503, Lyon, France, April 23–27, 2023. Springer, Heidelberg, Germany. (Cited on page 4.)
- [Rog02] Phillip Rogaway. Authenticated-encryption with associated-data. In Vijayalakshmi Atluri, editor, ACM CCS 2002: 9th Conference on Computer and Communications Security, pages 98–107, Washington, DC, USA, November 18–22, 2002. ACM Press. (Cited on page 8.)
- [SAB⁺22] Peter Schwabe, Roberto Avanzi, Joppe Bos, Léo Ducas, Eike Kiltz, Tancrède Lepoint, Vadim Lyubashevsky, John M. Schanck, Gregor Seiler, Damien Stehlé, and Jintai Ding. CRYSTALS-KYBER. Technical report, National Institute of Standards and Technology, 2022. available at https://csrc.nist.gov/Projects/post-quantum-cryptography/ selected-algorithms-2022. (Cited on page 8.)
- [UDB⁺15] 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, pages 232–249, San Jose, CA, USA, May 17–21, 2015. IEEE Computer Society Press. (Cited on pages 3, 4, and 15.)
- [UG15] Nik Unger and Ian Goldberg. Deniable key exchanges for secure messaging. In Indrajit Ray, Ninghui Li, and Christopher Kruegel, editors, ACM CCS 2015: 22nd Conference on Computer and Communications Security, pages 1211–1223, Denver, CO, USA, October 12– 16, 2015. ACM Press. (Cited on pages 3, 4, 15, and 16.)

- [UG18] Nik Unger and Ian Goldberg. Improved strongly deniable authenticated key exchanges for secure messaging. *Proceedings on Privacy Enhancing Technologies*, 2018(1):21–66, January 2018. (Cited on pages 3, 4, and 16.)
- [Ung21] Nik Unger. End-to-End Encrypted Group Messaging with Insider Security. PhD thesis, University of Waterloo, Ontario, Canada, 2021. https://hdl.handle.net/10012/17196. (Cited on page 15.)
- [VGIK20] Nihal Vatandas, Rosario Gennaro, Bertrand Ithurburn, and Hugo Krawczyk. On the cryptographic deniability of the Signal protocol. In Mauro Conti, Jianying Zhou, Emiliano Casalicchio, and Angelo Spognardi, editors, ACNS 20: 18th International Conference on Applied Cryptography and Network Security, Part II, volume 12147 of Lecture Notes in Computer Science, pages 188–209, Rome, Italy, October 19–22, 2020. Springer, Heidelberg, Germany. (Cited on pages 3, 4, 6, 7, 14, 18, 19, 20, and 36.)
- [YGS23] Tarun Kumar Yadav, Devashish Gosain, and Kent E. Seamons. Cryptographic deniability: A multi-perspective study of user perceptions and expectations. In Joseph A. Calandrino and Carmela Troncoso, editors, 32nd USENIX Security Symposium, USENIX Security 2023, Anaheim, CA, USA, August 9-11, 2023, pages 3637–3654. USENIX Association, 2023. (Cited on page 5.)
- [YZ10] Andrew Chi-Chih Yao and Yunlei Zhao. Deniable internet key exchange. In Jianying Zhou and Moti Yung, editors, ACNS 10: 8th International Conference on Applied Cryptography and Network Security, volume 6123 of Lecture Notes in Computer Science, pages 329–348, Beijing, China, June 22–25, 2010. Springer, Heidelberg, Germany. (Cited on page 3.)
- [YZ13] Andrew Chi-Chih Yao and Yunlei Zhao. OAKE: a new family of implicitly authenticated Diffie-Hellman protocols. In Ahmad-Reza Sadeghi, Virgil D. Gligor, and Moti Yung, editors, ACM CCS 2013: 20th Conference on Computer and Communications Security, pages 1113–1128, Berlin, Germany, November 4–8, 2013. ACM Press. (Cited on pages 2, 3, 36, 38, and 39.)

A Comparing prior definitions to our model

This section discusses in more detail how our model captures the idea of prior definitions from the literature (sketched by Table 3). Keep in mind that whenever we write CHALL we refer to the challenge oracles for both roles (i.e. CHALLINIT, CHALLRESP), and similarly for CHALLHON and CHALLHONINIT, CHALLHONRESP.

The following deniability models are not concerned with user messages. Though, we expect that they can be adapted similar to our model to handle key exchange protocols with user messages.

A.1 Concurrent and partial deniability [DGK06]

Deniability in the context of key exchange was introduced by Di Raimondo, Gennaro and Krawczyk [DGK06], namely concurrent deniability and partial deniability. Let us first focus on *concurrent deniability* [DGK06, Definition 2]. Both definitions are simulation-based definitions, i.e. a distinguisher has to decide if it is presented with values sampled from the real distribution or the simulated distribution. In either case, key pairs for all honest users and auxiliary info (which may depend on the keys) are sampled first. On a technical note, the probability of the distinguisher succeeding is taken over the randomness in generating the keys (and consequentially the sampling of the auxiliary info) and the random coins of the adversary. Note that the distinguisher gets access to the auxiliary info as well.

Type of Deniability	O_A	O_F	$ O_D$	Comment
concurrent deniability [DGK06]	Reg, Init, Chall	Ø	AUX	
partial deniability [DGK06]	Reg, Init, Chall	$USER_{x,y}$	AUX	$q_U = 1$
peer deniability [CF11]	REG, REGHON, INIT, CHALL	$USER_{x,y}$, SK	Ø	$AuxPrep = \bot$, session keys not deniable
peer-and-time deniability [CF11]	REG, REGHON, INIT, CHALL	SK	Ø	AuxPrep has transcripts, session keys not deniable
deniability (against malicious adv.) [HKKP22]	Reg, Init, Chall	Ø	Ø	$AuxPrep = \bot$
deniability (against semi-honest adv.) [HKKP22]	RegHon, ChallHon	Ø	Ø	$AuxPrep = \bot$
outsider deniability [DFG ⁺ 13]	ChallPassive	Ø	Ø	AuxPrep has valid transcripts
deniability $[BFG^+22, BFG^+21]$	ChallPassive	SK	SKs	$AuxPrep = \bot$
HP-deniability [YZ13]	ChallPassive	Ø	Ø	$AuxPrep = \bot$
$DENY^{false} \text{ and} \\ DENY^{true} \text{ [CHDN^+24]}$	CHALLPASSIVE	SK	SKs	$AuxPrep = \bot$
deniability [JS08]	Reg, Init, Chall	Ø	AUX	AuxPrep has secret keys
deniability [Jia14]	Reg, Init, Chall	Ø	AUX	AuxPrep has secret keys, valid transcripts
deniability [JCL ⁺ 22]	Reg, Init, Chall	Ø	AUX	AuxPrep has secret keys, valid transcripts

Table 3: Relation between our model and prior work: How to achieve deniability notions similar to those in the literature. Challenge oracles not suffixed by a role refer to the oracles for both roles. Deniability $[BFG^+22, BFG^+21]$, outsider deniability $[DFG^+13]$, $DENY^{false}$ and $DENY^{true}$ $[CHDN^+24]$ are game-based definitions; the others are simulation-based.

In the real case, the adversary may then interact with oracles representing the honest users in an arbitrary manner; the output consists of the previously sampled public keys and auxiliary info, the randomness that the adversary used, the transcripts between the adversary and all oracles, and the session keys of all completed sessions. The astute observer notices the similarity to our game-based model with b = 0for $O_A = \{\text{REG}, \text{INIT}, \text{CHALL}\}$ (with syntactic changes). In the simulated case, a simulator (who depends on the adversary) is given the same input as the adversary (i.e. the public keys, the auxiliary info, and the randomness that the adversary used) and needs to produce indistinguishable output. Note that the distinguisher receives the auxiliary info as well.

The proof technique for this definition suggested by [DGK06] and used thenceforth [CF11, VGIK20, HKKP22, JCL⁺22] is that the simulator executes the adversary and answers in place of the oracle. Subsequent work [JS08, Jia14] follows this approach and directly tasks the simulator to answer in place of the oracles. The advantage of this technique is that the final transcript then contains identically distributed messages by the adversary. Our game-based definition covers this case with b = 1 by using the Fake algorithm for simulation. To be more precise, the Fake algorithm of our game-based definition depends only on the currently challenged session and does not have a shared state over all invocations. The authors believe that this is also intended by [DGK06], who focus on concurrent sessions, where the simulation needs to happen regardless of other sessions. However, they need the simulator to have an overall state to manage the queries of the adversary.

Hence, we capture the idea of their concurrent deniability with $O_A = \{\text{REG}, \text{INIT}, \text{CHALL}\}, O_F = \emptyset$, $O_D = \{\text{AUX}\}$, and arbitrary AuxPrep.

Their partial deniability [DGK06, Definition 3] allows the simulator to access an oracle acting as the challenged user partnered with another user. So we capture the idea of their partial deniability with $O_A = \{\text{REG}, \text{INIT}, \text{CHALL}\}, O_F = \{\text{USER}_{x,y}\}, O_D = \{\text{AUX}\}, \text{ arbitrary AuxPrep, and } q_U = 1, \text{ i.e. per session for which the Fake algorithm is queried it can start one session with the target user (who is then partnered with a new, honest user).}$

A.2 Peer deniability and peer-and-time deniability [CF11]

Peer deniability and peer-and-time deniability [CF11] are simulation-based notions: The former is based on partial deniability [DGK06] and the latter is a strengthening of peer deniability. Neither definition uses auxiliary info.

A.2.1 Peer deniability

Starting from partial deniability [DGK06], peer deniability [CF11, Definition 9] is conceptually close with a few important differences:

- Peer deniability does not guarantee deniability of session keys. Hence, in line 7 of Figure 2 the distinguisher does not get access to the session keys computed by the CHALL oracle.
- The simulator can access polynomially many sessions to create messages, i.e. q_U is not fixed to 1.
- They give the attacker access to a corruption oracle, which allows the attacker *and the simulator* to learn the secret key(s) of the corrupted parties. We model the corruption oracle with the REGHON and SK oracles: The former allows the adversary to learn the secret key of a honestly generated key pair, and the latter for the Fake algorithm.¹⁵ However, we retain the REG oracle from partial deniability [DGK06]. In consequence, the Fake algorithm cannot rely on the SK oracle, just like the simulator cannot rely on learning the peer's secret key from the set of corrupted parties (maybe the adversary does not corrupt the party in question).

Hence, we capture the idea of their definition by setting $O_A = \{\text{INIT}, \text{REG}, \text{REGHON}, \text{CHALL}\}, O_F = \{\text{USER}_{x,y}, \text{SK}\}, O_D = \emptyset$, and $\text{AuxPrep} = \bot$. Additionally, we remove K as argument to \mathcal{D} in line 7 of Figure 2.

A.2.2 Peer-and-time deniability

Peer-and-time deniability [CF11, Definition 10] strengthens peer deniability (above) by restricting the simulator (i.e. the Fake oracle): The simulator has to finish interacting with the $USER_{x,y}$ oracle before interacting with the adversary. This ensures that the victim doesn't produce incriminating messages dependent on the adversary's input.

For our game-based approach we split up the simulator into many calls to Fake, each producing a single message. Hence, we cannot directly adapt the notion of Fake first interacting with a user oracle as preparation, and replying to the adversary only afterwards. We note, however, that peer-and-time deniability does not use auxiliary info, allowing us to repurpose the auxiliary info in our game-based definition: Originally, the simulator uses some strategy to query its oracle in an arbitrary manner. We encode this strategy into the AuxPrep algorithm and give aux only to Fake (and not to \mathcal{A}).

Both approaches (the simulator having access to an oracle before interacting with the adversary and our tailored sampling of auxiliary info) provide auxiliary info to Fake that the distinguisher is not aware of. Therefore, the Fake algorithm can use that info without the distinguisher noticing.

Hence, we capture the idea of their definition by setting $O_A = \{\text{INIT}, \text{REG}, \text{REGHON}, \text{CHALL}\}, O_F = \{\text{SK}\}, O_D = \emptyset$, and AuxPrep follows the simulator's strategy for querying USER_{x,y}. Additionally, we remove K as argument to \mathcal{D} in line 7 of Figure 2 and aux as argument to \mathcal{A} in line 6 in Figure 2.

¹⁵Strictly speaking, they give all corrupted secret keys to the simulator, while we give only the peer's secret key to the Fake algorithm. However, the authors do not see how a third party's secret key would be of help.

A.3 Deniability against malicious and semi-honest adversaries [HKKP22]

Deniability against malicious and semi-honest adversaries [HKKP22, Definition 7.1] are two simulationbased notions based on [DGK06]. Neither case allows any auxiliary information.

To capture the idea of their definition against malicious adversaries, we follow the same approach as for [DGK06] and set $O_A = \{\text{REG}, \text{INIT}, \text{CHALL}\}, O_F = \emptyset, O_D = \emptyset$, and $\text{AuxPrep} = \bot$. To capture the idea of their definition against semi-honest adversaries, we limit the adversary to accessing the oracles which enforce adherence to the protocol flow. Hence, we set $O_A = \{\text{REGHON}, \text{CHALLHON}\}, O_F = \emptyset, O_D = \emptyset$, and $\text{AuxPrep} = \bot$.

A.4 Outsider deniability [DFG⁺13]

Outsider deniability [DFG⁺13, Definition 2.5] is a game based-notion. Their notion gives the adversary access to an execute oracle (which yields a transcript) and a test oracle, which - depending on the secret bit - yields either a real transcript and session key or simulated values. Unlike our Fake algorithm, their simulator produces the complete transcript at once. Their notion is concerned with a passive adversary and a transcript should be simulatable with public data only. They do not allow auxiliary info, nor does the distinguisher get the adversary's randomness.

Hence, to capture the idea of their definition we set $O_A = \{CHALLPASSIVE\}, O_F = \emptyset, O_D = \emptyset$, and AuxPrep contains valid transcripts (substituting for the execute oracle). Additionally, we remove r as argument to \mathcal{D} in line 7 of Figure 2.

A.5 Deniability [BFG⁺22, BFG⁺21]

Deniability [BFG⁺21, Definition 11] is a game-based notion of 1-out-of-2 deniability against semi-honest adversaries in the big brother model. They tailor the deniability notion to their use case of a replacement for Signal's initial handshake: They argue that Bob's pre-key bundle does not need to be deniable since it is not bound to any peer or any session. Hence, it suffices for Alice's message to be simulatable by Bob to achieve 1-out-of-2 deniability.

Note that they consider Bob, who prepares the pre-key bundle, the responder, while we consider Bob the initiator.

Furthermore, they merge the adversary and the distinguisher into one algorithm \mathcal{A} . Syntactically, this has the consequence that the adversary-distinguisher already learns all secret keys while having access to the challenge oracle. Since the adversary-distinguisher is semi-honest, we are not aware how the knowledge of secret keys could help the adversary-distinguisher.

While they give a Fake algorithm, their Fake algorithm syntactically differs from ours: It directly produces a complete transcript and session key. Semantically this is equivalent for the use within our CHALLPASSIVE oracle. They do not use auxiliary info.

All in all, to capture the idea of their definition we set $O_A = \{CHALLPASSIVE\}, O_F = \{SK\}, O_D = \{SKS\}, and AuxPrep = \bot$, where the SK oracle only responds to queries for the secret keys of the initiator Bob.

A.6 HP-deniability [YZ13]

HP-deniability (short for honest-player-deniability) [YZ13, Definition 6.1] is a simulation-based definition that targets passive adversaries and exposes "pre-computed and stored session-states" to the distinguisher. In [YZ13] the authors stress that this ephemeral state is not necessarily equivalent to the random coins used for the session. We observe that some protocols, e.g. X3DH and PQXDH, instruct parties to delete

their ephemeral secrets after use. Hence, it is not generally obvious which information should be included in this ephemeral session state, but appears to be protocol-specific.

All in all, to capture the idea of their definition we set $O_A = \{CHALLPASSIVE\}, O_F = \emptyset, O_D = \emptyset$, and $AuxPrep = \bot$. Additionally, we need an extra game variable to save the ephemeral session state in: It is initialized in line 1 of Figure 2, populated in line 100 at the end of the CHALLPASSIVE oracle in Figure 8, and is given to the distinguisher in line 7 of Figure 2.

A.7 $DENY^{\text{false}}$ and $DENY^{\text{true}}$ [CHDN⁺24]

The $DENY^{\text{false}}$ notion adapts the idea of deniability [BFG⁺22, BFG⁺21] to the formalization from [CHDN⁺24]. Similarly, it is another game-based notion of 1-out-of-2 deniability against semi-honest adversaries in the big brother model. The $DENY^{\text{true}}$ notion extends $DENY^{\text{false}}$ by additionally giving the adversary-distinguisher access to the responder's state. They remark that the $DENY^{\text{false}}$ notion intuitively models the responder trying to frame the initiator, and the $DENY^{\text{true}}$ notion adds that the responder, who frames the initiator, cooperates with the judge by handing over his or her ephemeral state after the protocol run. As for HP-deniability [YZ13], we remark that some protocols, e.g. X3DH and PQXDH, require users to delete ephemeral data after completing the protocol run. Hence, it appears that the ephemeral session state needs to be defined and included on a per-protocol basis. They do not use auxiliary info.

All in all, to capture the idea of their definition we set $O_A = \{CHALLPASSIVE\}, O_F = \{SK\}, O_D = \{SKS\}, and AuxPrep = \bot$, where the SK oracle only responds to queries for the secret keys of the initiator Bob.

A.8 Deniability [JS08]

Deniability [JS08] is a simulation-based definition that allows adversarial corruptions. Note that adversary, simulator, and distinguisher all learn the corrupted secret keys. Hence, we model this non-adaptively by encoding the adversary's corruption strategy into AuxPrep, which then yields the corresponding secret keys to all three algorithms \mathcal{A} , Fake, and \mathcal{D} . While their definition leaves it implicit, we assume they allow the adversary to create keys maliciously.

To capture the idea of their definition we set $O_A = \{\text{Reg}, \text{INIT}, \text{CHALLRESP}\}, O_F = \emptyset$, $O_D = \{\text{AUX}\}$, and AuxPrep yields secret keys according to a given corruption strategy.

A.9 Deniability [Jia14]

Deniability [Jia14] is similar to the notion of [JS08] but additionally allows access to valid transcripts for all three algorithms \mathcal{A} , Fake, and \mathcal{D} .

To capture the idea of their definition we set $O_A = \{\text{REG}, \text{INIT}, \text{CHALLRESP}\}, O_F = \emptyset$, $O_D = \{\text{AUX}\}$, and AuxPrep yields secret keys according to a given corruption strategy and valid protocol transcripts.

A.10 Deniability [JCL⁺22]

Deniability [JCL⁺22], follows the ideas of [Jia14]. Additionally, they give the adversary access to a Reveal oracle, which yields the session key. We do not need this extra oracle since our model provides the session key to the adversary and distinguisher directly.

We can capture the idea of their definition as for [Jia14].

```
CHALLPASSIVE(U, V, \mathsf{info}_{\mathsf{create}}, \vec{\mu}):
 88 \pi_U.oid \leftarrow U; \pi_U.pid \leftarrow V
 89 \pi_V.oid \leftarrow V; \pi_V.pid \leftarrow U
 90 \pi_U.role \leftarrow initiator; \pi_V.role \leftarrow responder
 91 m_1 \leftarrow (\text{create}, \text{info}_{\text{create}})
 92 (\mu_1,\ldots,\mu_{n_m}) \leftarrow \vec{\mu}
 93 for i \in [1,3,\ldots,n_m-1] //until n_m if n_m is odd
            if b = 0
 94
                (\pi_{U}, m_{i+1}) \leftarrow \mathbb{R}un(\mathsf{sk}_{U}, \overrightarrow{\mathsf{pk}}, \pi_{U}, m_{i}, \mu_{i})
 95
                (\pi_V, m_{i+2}) \leftarrow \mathbb{R}un(\mathsf{sk}_V, \vec{\mathsf{pk}}, \pi_V, m_{i+1}, \mu_{i+1})
 96
 97
            else
                (\pi_{II}, m_{i+1}) \leftarrow \text{sFake}^{\mathcal{O}_F}(\vec{\mathsf{pk}}, \pi_{II}, m_i, \mu_i, \mathsf{aux})
 98
                (\pi_{II}, m_{i+2}) \leftarrow \text{Fake}^{O_F}(\vec{\mathsf{pk}}, \pi_{II}, m_{i+1}, \mu_{i+1}, \mathsf{aux}) //all Fake invocations share a state
 99
100 Q[\pi_U] \leftarrow (m_i)_{i=1}^{n_m};
                                                  K[\pi_U] \leftarrow \pi_U.K
101 return (Q[\pi_U], K[\pi_U])
```

Figure 8: The challenge oracle for passive adversaries.

B Extending the model to passive adversaries

For the sake of completeness, we extend our model to passive adversaries by adding a CHALLPASSIVE oracle, which the adversary may have access to via O_A , given in Figure 8.

In consequence, we also need to adapt Definition 3.1 to allow CHALLPASSIVE $\in O_A$. We obtain deniability against passive adversaries by setting {CHALLPASSIVE} = O_A , allowing at most q_C queries to the CHALLPASSIVE oracle. The CHALLPASSIVE limits the adversary to passively observing transcripts, without the opportunity to actively participate in the protocol execution or learning a party's secret keys.

It is easy to see that deniability against semi-honest adversaries implies deniability against passive adversaries.