

# KeyForge: Non-Attributable Email from Forward-Forgeable Signatures

Michael Specter  
MIT

Sunoo Park  
MIT & Harvard

Matthew Green  
Johns Hopkins University

## Abstract

Email breaches are commonplace, and they expose a wealth of personal, business, and political data whose release may have devastating consequences. Such damage is compounded by email’s strong attributability: today, any attacker who gains access to your email can easily prove to others that the stolen messages are authentic, a property arising from a necessary anti-spam/anti-spoofing protocol called DKIM. This greatly increases attackers’ capacity to do harm by selling the stolen information to third parties, blackmail, or publicly releasing intimate or sensitive messages — all with built-in cryptographic proof of authenticity.

This paper introduces *non-attributable email*, which guarantees that a wide class of adversaries are unable to convince discerning third parties of the authenticity of stolen emails. We formally define non-attributability, and present two system proposals — KeyForge and TimeForge — that provably achieve non-attributability while maintaining the important spam/spoofing protections currently provided by DKIM. Finally, we implement both and evaluate their speed and bandwidth performance overhead. We demonstrate the practicality of KeyForge, which achieves reasonable verification overhead while signing faster and requiring 42% *less* bandwidth per message than DKIM’s RSA-2048.

## 1 Introduction

Email is the world’s largest messaging scheme, used ubiquitously for personal, industry, and government communication. As such, it is a valuable target for attack: a user’s account represents a trove of sensitive information, unauthorized access to which enables spam, fraud, blackmail, and other abuse.

To help protect users from spam and fraud, the IETF developed a widely-adopted standard called DomainKeys Identified Mail (DKIM) [16]. DKIM’s goal is to assure the receiving server that each incoming message was really sent from the domain it appears to be from, enabling inter-domain accountability in case of spam and easy detection of spoofed messages. DKIM’s protocol is simple: the originating server cryptographically signs each outgoing email’s contents and metadata, allowing the receiving server to verify the message after looking up the sending server’s public key via DNS.

While DKIM was an important innovation that continues to be critical to the email ecosystem, its design came with an unintended side-effect: namely, email thieves can credibly convince any third party that stolen messages are authentic and unmodified via DKIM signatures from a reputable service provider. This increases incentives to break into email accounts, as a successful attacker can credibly (and anonymously) sell, publish, or use the stolen data for blackmail.

Email attributability has had real-world impact. For example, Wikileaks publicly asserts [65] that it relies on DKIM signatures to confirm the veracity of their publications: Wikileaks leveraged DKIM to authenticate messages stolen from the Democratic National Committee (DNC) and Hillary Clinton’s campaign chairman during the 2016 U.S. presidential election season [64]. Because of DKIM, any third party could easily confirm the legitimacy of these stolen messages using public keys tied to Google and Microsoft’s email services, despite the information’s questionable origin. Indeed, the practice of using DKIM to verify unauthorized email leaks has now become a standard journalistic practice [55, 42], with the Associated Press releasing a software tool for this purpose [5].

DKIM’s attributability problem has been recognized but unsolved for some time. Jon Callas, one of the original authors of the DKIM RFC, has publicly stated that attributability is an unintended design flaw of the protocol [18, 17], and has since suggested a number of ways to mitigate its impact, but notes that proposals at the time of his writing were insufficient or impractical [19]. Other researchers also flagged the issue as early as 2004, e.g., Adida *et al.* [1], Unger *et al.* [60], and Bellovin [8]; however, designing a practical, non-attributable DKIM replacement has remained an open question.

It is alarming that an unintended result of an ubiquitous messaging protocol has produced a *scalable, by-default* system for credible propagation of illicitly obtained private messages. The specific DNC incident might well have happened with or without DKIM: for a high-value target, interested parties would likely seek to verify the stolen emails in various ways, including non-technical methods (e.g., journalistic corroboration, cross-checking timestamps, geolocation, etc). But just the possibility of manual verification — a possibility that has existed since handwritten letters — is a stark contrast from the easy, inbuilt attribution that has *unintentionally* become ingrained in today’s email ecosystem.

Public figures are not the only victims of email breaches; new reports of email theft seem to surface every few weeks. Astoundingly, all of Yahoo!’s 3 billion email accounts were compromised in a 2013 breach [57]. Although Yahoo!’s users have been spared public dissemination of their messages, others (e.g., Sony and Stratfor), have been less fortunate [63, 62]. Attackers appear to have diverse motives, ranging from financial gain — e.g., selling patient healthcare data gleaned from emails [34] — to industrial espionage and monitoring political dissidents and foreign officials [29].

In light of the potential harm to users, it would be irresponsible to let DKIM’s unintended side-effect of attributability remain unscrutinized: if attributability is to remain a feature of DKIM, it should be as a result of a deliberate decision that takes into account the range of technically feasible alternatives. With the above as motivation, we ask:

*Is it possible to mitigate the potential harms of attributability in DKIM while maintaining the system’s efficient spam and spoofing-resistance?*

An initial intuition may be that attributability of stolen email is an unavoidable side-effect given the indirect and decentralized nature of email: it is intuitively unclear how a recipient with no communication to the sending server can be certain of a message’s origin without also gaining the ability to convince a third party of the same. Under certain conditions, this intuition amounts to an impossibility, as Perhaps surprisingly, our work shows that modern cryptography can reconcile the apparently conflicting goals of spam protection and non-attributability. We construct efficient protocols that achieve the important security guarantees that DKIM provides, while simultaneously *guaranteeing non-attributability* of stolen email. Further, we show that configurations of

our protocols are *practical* for deployment on the Internet today, achieving reasonable efficiency and bandwidth overhead.

## 1.1 Key Ideas

There are two main ideas underlying our proposals: *delayed universal forgeability* and *immediate recipient forgeability*.

**Delayed universal forgeability.** This approach ensures that signatures with respect to past emails “expire” after a time delay  $\Delta$  and thereafter become forgeable by the general public (i.e., arbitrary outsiders or non-parties). This property ensures that no attribution will be credible after the time delay has elapsed. We call this property *delayed universal forgeability*. As long as  $\Delta$  is set larger than the maximum viable time for email latency, the signature will still be convincing to the recipient at the time of receipt, thus maintaining the spam and spoofing-resistance of DKIM.

Signatures that possess delayed universal forgeability retain all the unforgeability properties of a standard signature scheme, until the set time  $\Delta$  has passed. Thus in cases where an attacker gains access to email and shows it to a third party within  $\Delta$  time after the email was sent, a third party will be convinced of the email’s authenticity. Effectively, delayed universal forgeability protects against adversaries that compromise an email account by breaking in and taking a snapshot (“after-the-fact attacks”), but not adversaries that fully control an email account and monitor its email in real time (“real-time attacks”). After-the-fact attacks cover a broad range of realistic attacks, for example, including many data breaches. Next, we discuss how we address real-time attacks.

**Immediate recipient forgeability.** Suppose that the fact of access to a particular client account implies the ability to forge messages from arbitrary other servers *to that recipient only*: that is, the ability to obtain valid DKIM signatures on email content and metadata of one’s choice. We call this *immediate recipient forgeability*. Importantly, the recipient constraint ensures the inability to impersonate any other server for the purposes of email addressed to *other* recipients, thus maintaining DKIM’s spam and spoofing-resistance. This undermines the credibility of attackers claiming ongoing access to a particular email account and attempting to convince third parties of the authenticity of emails supposedly sent to (and from) that account — even for real-time attacks, which may publish allegedly-incoming emails immediately as they are received.

Recipient forgeability is weaker than universal forgeability in the following sense: published emails credibly reveal that the attacker has gained access to some users’ key material, although not that the email content is authentic. Thus, recipient forgeability is not enough by itself; the two definitions are complementary and incomparable.

**Combining both ideas.** Our protocols attempt to achieve the “best of both worlds,” by providing universal forgeability when possible, and falling back on immediate recipient forgeability when necessary. Section 3 defines our threat model, discusses its limitations, and formalizes immediate recipient forgeability and delayed universal forgeability.

## 1.2 Overview of Solutions

This paper constructs and evaluates two base protocols KeyForge and TimeForge, and two enhanced variants KeyForge<sup>+</sup> and TimeForge<sup>+</sup> (which consist of the respective base protocol with a modified signing algorithm and one additional sub-protocol). The two base schemes can be seen as two

different approaches to building a new type of signature scheme that we introduce: *forward-forgeable signatures (FFS)*.

**Forward-forgeable signatures.** An FFS is a digital signature scheme equipped with a method to selectively disclose signature-invalidating “expiry information” for past signatures without similarly damaging the public key for future signatures. *Succinctness* of FFS is a measure of efficiency of disclosure. We present two constructions of FFS, which are the key building blocks of KeyForge and TimeForge respectively. FFS may be of independent interest as a signature primitive for other applications.

**KeyForge.** Our first proposal, KeyForge (§5.1), achieves *delayed universal forgeability* by publishing signing keys after a delay  $\Delta$ . KeyForge relies on an FFS based on hierarchical identity-based signatures (HIBS), which achieves logarithmic succinctness. As a result, KeyForge can efficiently distribute forging keys with minimal bandwidth.

**TimeForge.** Our second protocol, TimeForge (§5.2), assumes a *publicly verifiable timekeeper (PVTk)* model in which a trusted timekeeper periodically issues publicly verifiable timestamps. In a nutshell, the idea of TimeForge is to substitute each signature on a message  $m$  at time  $t$  with a succinct zero-knowledge proof of the statement  $S(m) \vee T(t + \Delta)$ , where:  $S(m)$  denotes knowledge of a valid signature by the sender on  $m$  and  $T(t + \Delta)$  denotes knowledge of a valid timestamp for a time later than  $t + \Delta$ . Including  $T(t + \Delta)$  ensures *delayed universal forgeability*, while  $R(m)$  ensures *immediate recipient forgeability*. TimeForge can be described as a forward-forgeable signature scheme in the PVTk model.

**KeyForge<sup>+</sup>/TimeForge<sup>+</sup>.** The enhanced protocols (§5.4) consist of the respective base protocols with the following modifications: (1) an additional protocol, called *forge-on-request*, that allows parties to request forged emails addressed *only to the requester herself* under limited circumstances; and (2) for multiple-recipient emails, a new signature is produced for each recipient domain (unlike the base protocols and DKIM, which produce one signature per outgoing email).

Among our protocols, KeyForge is the most efficient and would necessitate the least change to existing infrastructure. KeyForge<sup>+</sup> and TimeForge<sup>+</sup> are alternative approaches showing the feasibility of addressing stronger threat models though at significant overhead (in fact, certain overhead is unavoidable in the stronger threat model; see §3). TimeForge could become more practical with advances in the fast-moving area of non-interactive proofs.

### Summary of our Contributions.

1. We define non-attributability in store-and-forward email systems, and propose two *system designs* — KeyForge (§5.1), and TimeForge (§5.2) — that achieve this goal.
2. We *implement* KeyForge and TimeForge and evaluate their signing, verification, and bandwidth costs, and show that KeyForge has acceptable bandwidth and processing overhead for practical deployment (§6).
3. We provide formal *definitions* for email non-attributability and prove that our constructions realize them.
4. Of independent interest, we give *provably secure constructions* of a new cryptographic primitive, *succinct forward-forgeable signatures (FFS)* in both the standard and PVTk models (§4.3, §5.2).

## 2 Background on Email

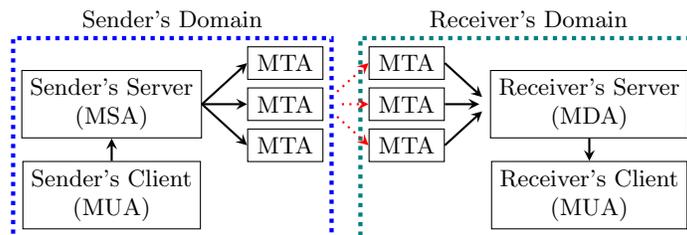


Figure 1: Simplified email routing infrastructure

This section introduces basic terminology of mail routing (as defined in RFC 5598 [25]) and describes how email infrastructure necessitates certain system requirements.

Figure 1 illustrates the architecture of email routing: an asynchronous routing protocol built on top of TCP/IP. Users first establish a relationship with a trusted email service provider, called a Mail Submission Agent (MSA) on the sender side and a Mail Delivery Agent (MDA) on the receiver side. The MDA is responsible for various tasks including verifying the authenticity of incoming messages.

The user's email client is called a Mail User Agent (MUA). Email originates from an MUA, and arrives at the user's trusted MSA. Depending on the system's configuration, the MSA may send the message to intermediary Mail Transfer Agents (MTAs) it trusts. Eventually, an MTA performs a DNS lookup for the receiving domain to discover which MTAs are authorized to process emails for that domain. The email is then sent via SMTP to one of these MTAs. After a number of hops depending on the sending and receiving organizations' infrastructure, the email reaches the receiver's MDA, which is responsible for verifying the message for the receiver's MUA.

### 2.1 Email Authentication

The IETF has developed a number of standards that allow domains to sign and verify incoming and outgoing messages. Next, we overview the three that have seen appreciable adoption: DKIM, SPF, and DMARC. Appendix B discusses a fourth, experimental protocol, ARC, and its potential impact.

**DKIM.** DomainKeys Identified Mail (DKIM) is an IETF standard that requires an MSA to sign outgoing email, and an MDA to verify that email by looking up the MSA's public key in the DNS. This procedure is described informally below:

1. *Setup:* The MSA generates a key pair and uploads the public key to the DNS in a TXT record.
2. *Sign:* The MSA adds the location of its public key to the email's metadata (or *header*), as well as additional metadata needed for signature verification, then signs the email and headers with its private key.<sup>1</sup>
3. *Verify:* On receipt, the MDA does a DNS lookup for the MSA's public key, and uses it to verify the signature.

<sup>1</sup>This usually includes a hash of the whole message, but the specification does allow for portions of the message to go unsigned. This is not default behavior for most DKIM applications, and has seen limited use in practice.

**SPF.** The Sender Policy Framework (SPF) ensures that intermediary MTAs are permitted to send and receive messages as a part of the domain. This solves a somewhat orthogonal problem to DKIM: SPF provides little guarantee that the message has not been modified by an intermediary, but instead provides spoofing protection by limiting what IP addresses are valid accepting MTAs.

**DMARC.** An SPF or DKIM failure as a result of a misconfiguration is indistinguishable from a failure due to an attempted message spoofing, and neither DKIM nor SPF provide mechanisms for alerting the sending domain that there has been a problem. DMARC solves this by adding a DNS TXT record specifying to the receiver what it should do in the case of such failures (such as quarantine, reject, or accept the message despite the failure), as well as providing an email address to send aggregated statistics on such failures.

## 2.2 DKIM Replacement Constraints

This section overviews a number of demands on email, some of which are not common to many other messaging systems. The combination of these demands means that achieving deniability and security presents distinct challenges for email, different from other messaging contexts.

**Indirectness by store and forward.** Email routing is a *store and forward* protocol in which messages are delivered indirectly via multiple hops, and routes, as well as the actual destination addresses, are often not known in advance. To quote the SMTP RFC [38], “[i]t is sometimes difficult for an SMTP server to determine whether or not it is making final delivery since forwarding or other operations may occur after the message is accepted for delivery.” Obvious examples of indirectness include *mail forwarding* (in which users configure their MDA to forward email received from an account on one domain to another), and *remailers* (such as mailing lists, that act as MUAs initially); however, there are other, less obvious, places in the ecosystem where this occurs.<sup>2</sup>

For example, many organizations leverage *third-party MTAs* that they do not own as an initial hop between the Internet and the organization’s self-hosted MDA/MSA.<sup>3</sup> These MTAs often provide security benefits to the MDA, such as protection from spam, malicious attachments, or DDoS attacks. While these intermediaries are allowed to quarantine messages or provide flow control to the MDA, under DKIM, they cannot undetectably modify or spoof emails.

In summary, these constraints inhere in store-and-forward systems: (1) final-destination information (e.g., addresses, keys) may be unknown to the sender, and (2) an MDA may not be certain whether it is the final destination of a message.

**Throughput and scalability.** Email is an any-mesh ecosystem: any domain owner must be able to set up the appropriate DNS records and interoperate with any other domain’s servers. Further, larger domains may sign and verify hundreds to millions of emails per day, and throughput requirements often increase over time. Therefore, beyond good constants on signing and verification time, the service must scale: adding more resources should provide linear or better performance.

---

<sup>2</sup>Similarly, Mail Retrieval Agents (MRAs) like Getmail [23] behave like MUAs to an MDA, but may forward emails on to an alternate, final MDA. Popular email services like Gmail provide services that download messages from other domains via IMAP.

<sup>3</sup>Third-party MTAs are commonplace. We did an informal survey by scraping DNS MX records for the Alexa top 150k. Surprisingly few, 31,615, have an MX record, and 10,260 use an obvious third-party hosting service (e.g., Google’s MTAs), leaving 21,615 that potentially self-host. Of the last category, 31.4% (6,793) are using a confirmed multi-hop third-party MTA. Raw results are in our repo in results.csv: <https://github.com/keyforgery/KeyForge>. This is likely a conservative estimate, as few servers appear to have matching domain names.

Scalability in interconnection with other servers is crucial too.<sup>4</sup>

Such scalability requirements indicate that certain types of overhead that would be trivial in other messaging contexts, (e.g., communication prior to sending a message or per-message round trips between servers), are unlikely to be viable for email. For example, it would be difficult to require the MDA to connect back to the original MSA for every email.

**Long-lived public keys.** One natural approach to short-lived signatures is to leverage correspondingly short-lived keys and publish each secret key at the end of its lifetime, or use short key sizes designed to be able to be brute-forced within the same period (see [19]). This approach has been mentioned in passing outside of the context of email [12]. Unfortunately, too-frequent key rotation entails practical problems that render this tactic unworkable for DKIM. Rotating keys stored in DNS is an often manual process that introduces risk of misconfiguration that can cause stability issues, and storing large amounts of key material that must be published, maintained, and shared among several servers is organizationally difficult and increases risk of key theft. DNS results are also often cached, so replacing an individual record is slow and can yield inconsistent results. Finally, it is hard to bound the time for short keys to be broken by all threat actors.

**Incremental deployment.** Given the myriad existing email servers and the need for interoperability, we consider the majority of the email ecosystem to be entrenched. It would be difficult to require substantial changes to mail routing, and it is unrealistic that every actor would promptly switch to a new scheme. Instead, it is far more realistic that DKIM could be replaced by incrementally updating the signing algorithms.

### 2.2.1 Resulting System Requirements

The particular constraints of email, described earlier in §2.2, rule out many natural approaches to non-attributability, including solutions that might be more feasible in other messaging environments. Since we treat email’s indirect, store-and-forward nature as an entrenched property of the infrastructure, realistic proposals for email protocol modifications must not rely on sender use of final-destination information, such as addresses or keys (“Requirement 1” or “R1”). Moreover, due to the store-and-forward and scalability requirements, email protocols should avoid interactive sender-receiver (MDA–MSA) communication whenever possible; in particular, we consider roundtrip sender-receiver communication per email to be inviable (“R2”). Additionally, email protocols must have long-lived public keys (“R3”).

Notably, *none* of the following approaches adhere to both the above requirements: interactive zero-knowledge proofs (violate R2); ring signatures (proposed for email non-attributability in [2, 12]) (violate R1); designated-verifier signatures (violate R1); short-lived keys with publication of secret keys after use (violate R3); and — importantly — systems based on deniable authenticated key exchange (DAKE) (which violates R2), such as OTR or Signal [12, 59, 60]. Indeed, both the OTR paper [12, §6] and a recent DAKE paper [59, §6.6] dedicate a full subsection to discussing the heightened challenges of non-attributability for email as compared to other messaging environments, and note that their proposals are not adequate for email due to its asynchronous, non-interactive, store-and-forward nature. (Table 3, Appendix A breaks down differences between email and other messaging schemes, and where our solutions fit in, in more detail.)

Finally, we note that the simple approach of relying on MDAs to delete DKIM header information after receipt is flawed *not only* because it fails to address our threat models (§3), which require

---

<sup>4</sup>The IETF standard for DMARC [39] states that pre-sending agreements is a poor scalability choice for this reason. See also [58].

security against *malicious or compromised recipients*, but also because it violates Requirement 1: relying on MDAs for deletion is untenable given that MDAs may not know if they are the final endpoint (and if not, the signatures must be kept for later verification).

**Summary.** A viable non-attributable replacement for DKIM must have: (1) compatibility with indirect, store-and-forward communication (in particular, no reliance on sender knowledge of final destination addresses or keys); (2) no requirement of sender-receiver interaction per email; (3) long-lived public keys; (4) no required behavior for MDAs that depends on whether they are the final destination; (5) little impact on other parts of the email ecosystem; and (6) good systems properties allowing for incremental, scalable deployment.

In addition, a *non-attributable* DKIM replacement must have universal forgeability, and should have recipient forgeability whenever feasible.

### 3 Model and Security Definitions

This section presents threat models and formal definitions of email non-attributability.

**Notation.** “PPT” means “probabilistic polynomial time.”  $|S|$  denotes the size of a set  $S$ .  $[n]$  denotes the set  $\{1, \dots, n\}$  of positive integers up to  $n$ , and  $\mathbb{P}(\cdot)$  denotes powerset.  $\approx_c$  denotes computational indistinguishability.  $\tau||e$  denotes the result of appending an additional element  $e$  to a tuple  $\tau$ .

#### 3.1 Model

**Time** We model time in discrete time-steps and assume fairly consistent (say, within 3 mins) local clocks. This is realistic given NTP [15].

**Synchrony**  $\hat{\Delta}$  is an upper bound on the time required for email delivery. Our parameter settings depend on  $\hat{\Delta}$ , and our evaluation sets  $\hat{\Delta}$  at 15 minutes (see § 5.1.1).

**DNS** Our model assumes all parties and algorithms have access to DNS and can update their own DNS records.

**Publication** We assume each party has a method of publishing persistent, updatable information retrievable by all other parties and algorithms. This could be via DNS or another medium, such as posting on a website.

#### 3.2 Threat Models

We are concerned with attacks that disclose private communications obtained at the MDA (whether because the MDA is compromised or because it is malicious).

We consider two threat models, defined below. KeyForge and TimeForge achieve security against Threat Model 1, which targets scenarios where attackers may gain access to an email server but are unlikely to maintain access for extended periods. The enhanced protocols KeyForge<sup>+</sup> and TimeForge<sup>+</sup> achieve security against Threat Model 2, the stronger of the two threat models, which is necessary in settings where attackers’ access may likely remain undetected for extended periods (e.g., advanced persistent threats).

**Threat Model 1.** (After-the-fact attacks) *In this model, the recipient is presumed honest at the time of email receipt, but is later compromised by an attacker that takes a snapshot of all stored email content.*

**Threat Model 2.** (Real-time attacks) *In this model, the recipient may be malicious at the time of email receipt, with ongoing and immediate intent to disclose received email content to third parties.*

**Ruling out trivial solutions.** A trivial and uninteresting way to achieve non-attributability, in either threat model, is not to sign emails at all. Of course, this is undesirable as it would undermine the spam- and spoofing-resistance for which DKIM was designed. Providing these guarantees is an implicit requirement throughout this paper. Moreover, since our threat models consider malicious receiving servers, any non-attributability that relies on receiving-server behavior — such as DKIM header deletion upon receipt — is unsatisfactory.

**Preventing real-time attacks requires interaction.** Any store-and-forward email protocol that both (1) allows recipients to verify the sending domain’s identity and (2) is secure against *real-time* attacks (Threat Model 2) must be interactive, as more formally detailed in Claim 1, Appendix C. Informally, in the store-and-forward model, a non-interactive protocol transcript (consisting of a single message from the sender), cannot depend on final-destination recipient information, so any operations (such as verification or forgery) that the verifier can run must also be executable by others. This also relates to the intuitive idea that someone who receives a single message  $m$  convincing them of the message’s origin must also be able to use  $m$  to convince others of the same.

In contrast, security against *after-the-fact* attacks (Threat Model 1) is possible non-interactively, as KeyForge and TimeForge exemplify. KeyForge<sup>+</sup> and TimeForge<sup>+</sup> augment KeyForge and TimeForge with an interactive (two-message) protocol, which adds significant overhead and complexity to the non-interactive base protocols. Claim 1 shows that this overhead is, in a sense, unavoidable. The overhead of our constructions is furthermore minimal in certain respects: just two rounds of interaction, and the protocols do not *require* interaction on email receipt, but rather, introduce the *possibility* of interaction by an additional protocol (details in §5.4).

**What’s outside our threat models?** While Threat Model 2 considers powerful real-time adversaries, it too has limits. Definitionally, and unsurprisingly, no deniability is possible against a global passive adversary that can be sure of observing all traffic as it flows over the network. As already mentioned, our threat models are not designed to provide non-attributability against adversaries directly observing email traffic, but rather against those to whom the adversaries might try to pass the stolen emails on.

Our threat models focus on attacks at the receiving server (MDA), because we believe this covers a wide, though not exhaustive, range of attack scenarios of interest. This notably excludes *malicious intermediaries (MTAs)*. Even though our threat models do not focus on MTA-based attacks, our protocols KeyForge<sup>+</sup> and TimeForge<sup>+</sup> do provide a partial non-attributability guarantee against malicious intermediaries (as discussed in §3.3). Nonetheless, malicious intermediaries pose a legitimate concern not fully addressed by this work; achieving stronger non-attributability guarantees against MTAs could be interesting future research.<sup>5</sup>

---

<sup>5</sup> It is also unclear how effective local MTA-based attacks would be to compromise entire email accounts; such attacks’ effectiveness would likely depend on email routing configurations at the servers involved. By entire-account compromise we mean learning all stored emails and/or all real-time emails for a single account over an extended period, as opposed to learning only occasional emails from scattered accounts. Entire-account compromise would be useful to target particular accounts, or to obtain a relatively complete picture of compromised accounts (e.g., for identity theft). In contrast, MDA-based attacks provide a direct way to compromise entire accounts.

Finally, we note that our definitions do not necessarily provide non-attributability against adversaries that can preconfigure the receiving server with custom secure hardware (see also §3.3). We consider such attacks outside our threat model: i.e., we assume servers are compromised after physical setup.

We conclude this section with additional context and explanation for our modeling choices.

**Client-server trust.** Email clients rely heavily on their email servers. A malicious email server could easily and undetectably misbehave in many essential functions: e.g., drop incoming emails, modify outgoing emails (since typically, emails are not signed client-side), or falsify content and metadata of incoming emails (since typically, clients do not perform DKIM verification themselves). Since client-server trust is very high in practice, this paper treats the client and server as a single entity, and relatedly, our threat models do not consider malicious behavior by MSAs that aims to undermine non-attributability of their own clients’ emails. (One might also argue such malicious behavior would quickly lose an MSA its clients.)

**Evidence-based credibility.** In a system where credibility is based on reputation rather than evidence — that is, where certain parties’ statements are taken on faith, or believed simply because of who they are even without supporting evidence — a “reputable” party with the ability to eavesdrop on the communication channel would be able to undermine non-attributability by keeping traffic logs. Our model assumes mutually distrustful parties: i.e., that no party is taken simply on its word as just described. In other words, credibility in our model is evidence-based and not reputation-based.

**Systemic attributability vs. attributability by choice.** The goal of non-attributability is to empower users to choose whether or not their messages are attributable, to disincentivize email theft and misuse in contrast to attributability-by-default (see §1). We are not concerned with preventing attributability when correspondents desire it: e.g., for business transactions or contracts, correspondents may intentionally sign messages to ensure they are binding. Attribution by journalistic investigation is also outside our threat model: confirmation of selected documents by careful investigation is possible even with handwritten letters, but the current *systemic* attributability facilitates scalable, malicious attribution far beyond the handful of high-profile messages that might be published after arduous manual verification.

### 3.3 Defining Non-Attributability

We define email non-attributability as a game involving an email protocol  $\mathbf{Email}$ , adversary  $\mathcal{A}$ , simulator  $\mathcal{S}$ , and distinguisher  $\mathcal{D}$ . For any email server  $S$  with internal state  $s$ ,  $\mathbf{Email}_s(S, R, m, \mu, t)$  denotes the information that recipient  $R$  receives when  $S$  legitimately sends at time  $t$  an email message  $m$  with metadata  $\mu$  to  $R$ . For simplicity, our notation leaves implicit that parties other than the sender participate in transmission and may affect the information received by  $R$  (e.g., the MTAs); however, the final received information  $\mathbf{Email}_s(S, R, m, \mu)$  should be thought to contain any modifications made en route between  $S$  and  $R$ . While this definition refers to “internal state  $s$ ” for generality, the state  $s$  can essentially be thought of as secret key material.

Intuitively, we require indistinguishability between a legitimate email  $\mathbf{Email}_s(S, R, m, \mu, t)$  and a “fake” email that was created without access to the sending server at all: that is, without knowing  $s$ . To model this, we consider a simulator  $\mathcal{S}$  that “aims” to create such an email without knowing  $s$ , and our security definition requires  $\mathbf{Email}_s(S, R, m, \mu, t)$  be indistinguishable from  $\mathcal{S}$ ’s output.

Next, we give two formal definitions of non-attributability. *Recipient non-attributability* (Definition 1) considers a simulator that has access to a particular recipient’s email server, and is required to output email from any sender to that recipient.  $\Delta$ -*universal non-attributability* (Definition 2) is a stronger definition whose simulator is required to output email from any sender to any recipient while having access to neither the sender’s nor the receiver’s email server.

In other words, if an adversary publishes an email allegedly authored by an honest party, Definition 1 guarantees that the victim can credibly argue that, granting that the attacker indeed broke into her correspondent’s account, the attacker’s allegations inherently lack credibility because by the very fact of such access, he could have forged arbitrary emails between them. Definition 2 gives the stronger guarantee that *anyone* can forge past emails after some delay  $\Delta$ : so after  $\Delta$ , any allegations are even less credible and the email accounts may not have been compromised at all. The two definitions are complementary and incomparable.

**Definition 1** (Recipient non-attributability). *Email is non-attributable for recipients if there is a PPT simulator  $\mathcal{S}$  such that for any sender  $S$  and recipient  $R$ , for any email message  $m$  and metadata  $\mu$ ,*

$$\text{Email}_s(S, R, m, \mu, t) \approx_c \mathcal{S}^R(S, m, \mu) ,$$

where the superscript  $R$  denotes that  $\mathcal{S}$  has the capability of sending outgoing mail as  $R$  through  $R$ ’s email server.

**Definition 2** ( $\Delta$ -universal non-attributability). *For  $\Delta \in \mathbb{N}$ , an email protocol  $\text{Email}$  is  $\Delta$ -strongly non-attributable if there is a PPT simulator  $\mathcal{S}$  such that for any sender  $S$  (with internal state  $s$ ) and recipient  $R$ , for any email message  $m$ , metadata  $\mu$ , and timestamp  $t$ , the following holds at any time  $\geq t + \Delta$ :*

$$\text{Email}_s(S, R, m, \mu, t) \approx_c \mathcal{S}(S, R, m, \mu, t) .$$

Definitions 1 and 2 serve to ensure that no attacker can credibly claim to a third party<sup>6</sup> that he is providing her with authentic emails: the third party is in the role of distinguisher.

Note that Definition 2 is inviable if  $\Delta < \hat{\Delta}$ . Otherwise, the spam- and spoofing-resistance provided by DKIM would be undermined, since any outsider could use the simulator in real time to send spam email indistinguishable to the recipient from email actually sent by an honest party.

**Relation to the threat models.**  $\Delta$ -universal non-attributability achieves non-attributability against after-the-fact attacks (Threat Model 1) for all emails sent and received at least  $\Delta$  before the server is compromised.

Combining recipient non-attributability and  $\hat{\Delta}$ -universal non-attributability (Definitions 1 and 2) yields non-attributability against real-time attacks (Threat Model 2). A real-time attacker with ongoing access to an email server can easily make the fact of his access evident by immediately publishing all emails he sees (within time  $\hat{\Delta}$  of receipt), but will be unable to convince third parties of any given email’s authenticity since the fact of his access to the server allows him to forge emails in real time, under Definition 1. For allegedly compromised emails from more than  $\hat{\Delta}$  ago, an attacker’s credibility is even lower, since for such past timestamps *anyone with internet access* can generate seemingly validly signed emails, even without breaking into any email server at all, under Definition 2.

---

<sup>6</sup>E.g., the general public (if the allegedly stolen emails are released publicly) or a specific interested party (such as a potential buyer or disseminator of the information).

**Necessity of recipient forgeries.** It may seem a counterintuitive or risky design choice to enable real-time email forgery in any part of the system. If forgery is restricted only to recipients forging emails to themselves, as in our definition, there is no spam/spoofing vulnerability — but given the choice, one might avoid introducing any forging capability at all, in the interest of a simpler and easier-to-analyze system. However, some sort of real-time forging capability by recipients is definitionally necessary to achieve non-attributability against real-time attacks: if the recipient cannot forge in real time, then any third party to whom a recipient server passes emails in real time must be convinced of the emails’ authenticity.

**Other inherent model constraints.** A practical consequence of recipient non-attributability is that a recipient  $R$ ’s email server can, *unknown to  $R$* , create fraudulent messages that appear to be legitimate emails from any sender to  $R$ , and deliver them to  $R$ . As discussed §3.2, the current email system necessitates heavy client-server trust. In this context, recipient non-attributability does not meaningfully increase the trust a client places in her email server. For example, email servers in the current system could (and often do) omit DKIM headers when delivering emails to clients: this effectively implies the ability to deliver fake messages.

Also, we note that both definitions allow for strong, persistent attackers to convince others of the very fact that they have ongoing access to a particular email account. The definitions guarantee that even so, such attackers cannot make credible claims about email contents, since they gain the ability to falsify emails by the very fact of their access. That attackers with ongoing access can prove their access is unavoidable since universal forgeability is incompatible with spam resistance for too small  $\Delta$ , as discussed above.

**Adversarial secure hardware at recipient.** The requirement of spam- and spoofing-resistance means that any simulator  $\mathcal{S}$  satisfying Definition 1 must use the recipient  $R$ ’s secret state  $r$ : in order to prevent spam, real-time forgery must be limited to messages whose recipient is the forger herself. This suggests that recipient non-attributability would lose meaning in an extreme situation where every use of  $r$  can be monitored and attested to, since then an attacker could prove that  $\mathcal{S}$  was never invoked on  $r$ . This might be plausible assuming secure hardware, e.g., by generating and monitoring all uses of  $r$  within a secure enclave (as suggested in [33]) — but even then, such an attack would likely only be feasible by the unlikely attacker who has designed her recipient email server with this unlikely configuration from its very setup. We note this possibility for completeness, but *such attacks are outside our threat models*, as mentioned earlier in §3.

**Malicious intermediaries and traffic logging.** Although our threat models focus on malicious recipient servers (as discussed earlier in §3), Definition 1 actually provides a meaningful, though limited, guarantee against malicious intermediaries (MTAs) as well. If a malicious MTA were to log all traffic and publish it in real time (perhaps even timestamped in a trustworthy way for future reference), in a system with immediate recipient forgeability, observers of the publications would still be unconvinced of: (1) whether any email the MTA claims is genuine (unforged) is really genuine, since the MTA could have *omitted* evidence of forgery, and (2) whether the MTA omitted any genuine emails from its publications.

**Why (sometimes) settle for weaker non-attributability?** KeyForge and TimeForge achieve only non-attributability against after-the-fact attacks, and their enhanced versions KeyForge<sup>+</sup> and TimeForge<sup>+</sup> are non-attributable against both after-the-fact and real-time attacks. Yet we consider KeyForge to be our main protocol and the most realistic proposal for deployment. In practice, the enhanced protocols’ (unavoidable) interactivity and other overhead would often be compelling

reasons to prefer the simpler base protocols except in contexts where addressing real-time attacks (or malicious intermediaries) is of particular heightened concern.

**Relation to deniability definitions in other contexts.** The cryptographic literature features many works on deniability of signatures and authentication, including (but not at all limited to) [37, 27, 54, 44, 26]. Our constructions could be seen as a practical instantiation of a deniable signature scheme subject to tight systems-based requirements.

Cryptographic deniability definitions tend to come in two flavors: denying communication *content*, or denying having participated in communication at all. Deniable encryption [22] does the former, whereas deniable authentication generally does the latter. Our *recipient-forgeability* is of the former flavor, whereas our *universal forgeability* is of the latter.

## 4 Forward-Forgeable Signatures

### 4.1 Definition

Definition 3 formalizes *forward-forgeable signatures (FFS)*. They are a new primitive that this paper introduces, and are an essential building block for our proposed protocols. Informally, FFS are signature schemes equipped with a method to selectively “expire” past signatures by releasing *expiry information* that makes them forgeable. In an FFS, each signature is made with respect to a *tag*  $\tau$ , which is an arbitrary string. Expiry information can be released with respect to any tag or set of tags. In our context, the tag can be thought to be a timestamp: i.e., each email is signed with respect to the current time  $\tau$ , and at some later time  $\tau + \Delta$ , the signer may publish expiry information for  $\tau$ . FFS have *correctness* and *unforgeability* requirements similar to standard signatures, as well as a new requirement, *forgeability on expiry*, that has no analogue in standard signatures.

The *correctness* requirement of FFS is the same as that of standard signatures. The *unforgeability* requirement is modified to include an *expiry oracle*: that is, unforgeability of signatures w.r.t. non-expired tags must hold even in the presence of arbitrary, adversarially chosen expirations. The *forgeability on expiry* requirement is a feature of FFS that has no analogue in standard signatures.

**Definition 3 (FFS).** A forward-forgeable signature scheme (FFS)  $\Sigma$  is implicitly parametrized by message space  $\mathcal{M}$  and tag space  $\mathcal{T}$ , and consists of five algorithms

$$\Sigma = (\text{KeyGen}, \text{Sign}, \text{Verify}, \text{Expire}, \text{Forge})$$

satisfying the following syntax and requirements.

SYNTAX:

- $\text{KeyGen}(1^\kappa)$  takes as input a security parameter<sup>7</sup>  $1^\kappa$  and outputs a key pair  $(vk, sk)$ .
- $\text{Sign}(sk, \tau, m)$  takes as input a signing key  $sk$ , a tag  $\tau \in \mathcal{T}$ , and a message  $m \in \mathcal{M}$ , and outputs a signature  $\sigma$ .
- $\text{Verify}(vk, \tau, m, \sigma)$  takes as input a verification key  $vk$ , a tag  $\tau \in \mathcal{T}$ , a message  $m \in \mathcal{M}$ , and a signature  $\sigma$ , and outputs a single bit indicating whether or not  $\sigma$  is a valid signature with respect to  $vk$ ,  $m$ , and  $\tau$ .

---

<sup>7</sup>Technically, all five algorithms take  $1^\kappa$  as an input, and  $\mathcal{M}$  and  $\mathcal{T}$  may be parametrized by  $\kappa$ . For brevity, we leave this implicit except in  $\text{KeyGen}$ .

- $\text{Expire}(sk, T)$  takes as input a signing key  $sk$  and a tag set  $T \subseteq \mathcal{T}$ , and outputs expiry info  $\eta$ .
- $\text{Forge}(\eta, \tau, m)$  takes as input expiry info  $\eta$ , a tag  $\tau \in \mathcal{T}$ , and a message  $m \in \mathcal{M}$ , and outputs signature  $\sigma$ .

REQUIRED PROPERTIES:

1. Correctness: For all  $m \in \mathcal{M}, \tau \in \mathcal{T}$ , there is a negligible function  $\varepsilon$  such that for all  $\kappa$ ,

$$\Pr \left[ \begin{array}{l} (vk, sk) \leftarrow \text{KeyGen}(1^\kappa) \\ \sigma \leftarrow \text{Sign}(sk, \tau, m) \\ b \leftarrow \text{Verify}(vk, \tau, m, \sigma) \end{array} : b = 1 \right] \geq 1 - \varepsilon(\kappa) .$$

2. Unforgeability: For any PPT  $\mathcal{A}$ , there is a negligible function  $\varepsilon$  such that for all  $\kappa \in \mathbb{N}$ ,

$$\Pr \left[ \begin{array}{l} (vk, sk) \leftarrow \text{KeyGen}(1^\kappa) \\ (\tau, m, \sigma) \leftarrow \mathcal{A}^{\mathcal{S}_{sk}, \mathcal{E}_{sk}}(vk) \\ b \leftarrow \text{Verify}(vk, \tau, m, \sigma) \\ b' = \tau \notin Q'_E \wedge (\tau, m) \notin Q_S \end{array} : b = b' = 1 \right] \leq \varepsilon(\kappa) ,$$

where  $\mathcal{S}_{sk}$  and  $\mathcal{E}_{sk}$  respectively denote oracles  $\text{Sign}(sk, \cdot, \cdot)$  and  $\text{Expire}(sk, \cdot)$ ,  $Q_S$  and  $Q_E$  denote the sets of queries made by  $\mathcal{A}$  to the respective oracles, and  $Q'_E = \bigcup_{T \in \mathcal{Q}_E} T$ .

3. Forgeability on expiry: For all  $m \in \mathcal{M}, T \subseteq \mathcal{T}$ , for any  $\tau \in T$ , for any “distinguisher” algorithm  $\mathcal{D}$ , there is a negligible function  $\varepsilon$  such that for all  $\kappa$ ,

$$\Pr \left[ \begin{array}{l} (vk, sk) \leftarrow \text{KeyGen}(1^\kappa) \\ \sigma_0 \leftarrow \text{Sign}(sk, \tau, m) \\ \eta \leftarrow \text{Expire}(sk, T) \\ \sigma_1 \leftarrow \text{Forge}(\eta, \tau, m) \\ b \leftarrow \{0, 1\} \\ b' \leftarrow \mathcal{D}(\sigma_b, \eta) \end{array} : b = b' \right] \leq 1/2 + \varepsilon(\kappa) .$$

That is,  $\mathcal{D}$  must not be able to distinguish whether a signature was produced using  $\text{Sign}$  or  $\text{Forge}$ , even in the presence of the expiry information  $\eta$ .

The *forgeability on expiry* property requires computational indistinguishability between a signature produced using  $\text{Sign}$  and a signature produced using  $\text{Forge}$  on valid expiry information. In particular, when combined with the *correctness* property, this implies that any such signature produced with  $\text{Forge}$  must appear to be a valid signature, i.e., cause  $\text{Verify}$  to output 1.

**FFS in the publicly verifiable timekeeper model** Recall the *publicly verifiable timekeeper* (PVTk) model,<sup>8</sup> where a reliable timekeeper periodically issues publicly verifiable timestamps. In this model, expiration may occur “automatically”: signers need not publish additional expiry information for signatures to become forgeable after a delay. Thus, the algorithm  $\text{Expire}$  is unnecessary, and  $\text{Forge}$  need not take input  $\eta$ . In §5.2, we construct FFS in the PVTk model.

**Difference with forward-secure signatures** Both FFS and FSS yield a system of short-lived secret keys all corresponding to one long-lived public key. However, the definitions differ in two main ways, described below and depicted in Figure 2.

<sup>8</sup>Mentioned in § 1.2; discussed further in § 5.2.

1. Forward-*secure* signatures require that past keys cannot be computed from future keys, whereas forward-*forgeable* signatures require that future keys cannot be computed from past keys.
2. Forward-*secure* signatures are designed to prevent compromise of past signatures by compromising a later secret key. All FFS secret keys are short-lived and each secret key must be derivable based solely on the previous short-lived secret key. Forward-*forgeable* signatures, in contrast, may have persistent “master secret key” material used to generate each short-lived key.

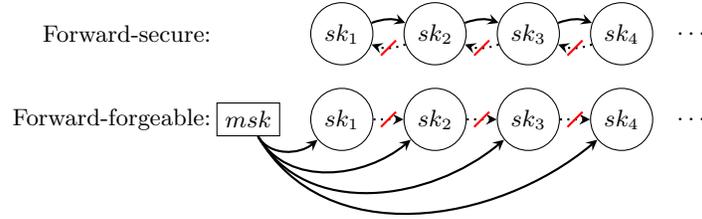


Figure 2: Forward-secure vs. forward-forgeable signatures

## 4.2 Succinctness

Next, we define *succinctness* of FFS, a measure of the efficiency of disclosure in terms of the size of expiry info per tag expired. Concretely, in our application, succinctness measures how expiry info scales as more non-attributable emails are exchanged over time. KeyForge uses a construction of FFS based on hierarchical identity-based signatures (§4.3), which achieves logarithmic succinctness.

**Definition 4.** Let  $z : \mathbb{N} \rightarrow \mathbb{N}$ . Let  $S \subset \mathbb{P}(\mathcal{T})$  be a set of sets of tags. A forward-forgeable signature scheme  $\Sigma$  is  $(S, z)$ -succinct if for any  $T \in S$ , there is a negligible function  $\varepsilon$  such that for all  $\kappa$ ,

$$\Pr_{(vk, sk) \leftarrow \text{KeyGen}(1^\kappa)} \left[ |\text{Expire}(sk, T)| \leq z(|T|) \right] \geq 1 - \varepsilon(\kappa) .$$

*Remark 1.* Definition 4 is a worst-case definition: it guarantees the size of expiry information with overwhelming probability. In certain applications, an average-case definition may be appropriate instead, i.e., defining succinctness by bounding the size on *expectation*. We use a worst-case definition since it is stronger than an average-case definition, and our construction achieves it.

## 4.3 FFS Construction from (Hierarchical) IBS

We first outline a simple FFS construction `BasicFFS` based on identity-based signatures (IBS) [56], as a stepping stone to our main construction from hierarchical IBS (HIBS). The next paragraph assumes familiarity with standard IBS terminology; readers unfamiliar with IBS may skip to the main construction which is explained formally with explicit syntax definitions.

Let tags in the FFS correspond to identities in the IBS. `BasicFFS.KeyGen` outputs IBS master keys. The `BasicFFS` signing and verification algorithms for tag  $\tau$  respectively invoke the IBS signing and verification algorithms for identity  $\tau$ . `BasicFFS.Expire` outputs the secret key for each input tag  $\tau \in T$ , and `BasicFFS.Forge` uses the appropriate secret key from the expiry information to invoke

the IBS signing algorithm. This simple solution has linear succinctness. By leveraging hierarchical IBS (HIBS), our main construction achieves logarithmic succinctness, as described next.

**Definition 5.** A hierarchical identity-based signature scheme HIBS is parametrized by message space  $\mathcal{M}$  and identity space  $\mathcal{I} = \{\mathcal{I}_\ell\}_{\ell \in \mathbb{N}}$ , and consists of four algorithms  $\text{HIBS} = (\text{Setup}, \text{KeyGen}, \text{Sign}, \text{Verify})$  with the following syntax.

- **Setup**( $1^\kappa$ ) takes as input a security parameter<sup>9</sup> and outputs a master key pair  $(\text{mvk}, \text{msk})$ .
- **KeyGen**( $sk_{\vec{id}}, id$ ) takes as input a secret key  $sk_{\vec{id}}$  for a tuple of identities  $\vec{id} = (id_1, \dots, id_\ell) \in \mathcal{I}_1 \times \dots \times \mathcal{I}_\ell$  and an additional identity  $id \in \mathcal{I}_{\ell+1}$  and outputs a signing key  $sk_{\vec{id}'}$  where  $\vec{id}' = (id_1, \dots, id_\ell, id)$ . The tuple may be empty (i.e.,  $\ell = 0$ ): in this case,  $sk_{\vec{id}'} = \text{msk}$ .
- **Sign**( $sk_{\vec{id}'}, m$ ) takes as input a signing key  $sk_{\vec{id}'}$  and a message  $m \in \mathcal{M}$ , and outputs a signature  $\sigma$ .
- **Verify**( $\text{mvk}, \vec{id}, m, \sigma$ ) takes as input master verification key  $\text{mvk}$ , tuple of identities  $\vec{id}$ , message  $m \in \mathcal{M}$ , and signature  $\sigma$ , and outputs a single bit indicating whether or not  $\sigma$  is a valid signature with respect to  $\text{mvk}$ ,  $\vec{id}$ , and  $m$ .

A depth- $L$  HIBS is a HIBS where the maximum length of identity tuples is  $L$ , i.e., the identity space is  $\mathcal{I} = \{\mathcal{I}_\ell\}_{\ell \in [L]}$ .

**Definition 6.** For an identity space  $\mathcal{I} = \{\mathcal{I}_\ell\}_{\ell \in \mathbb{N}}$ , we say  $\vec{id}$  is a level- $\ell$  identity if  $\vec{id} \in \mathcal{I}_1 \times \dots \times \mathcal{I}_\ell$ . For any  $\ell' > \ell$ , let  $\vec{id}$  be a level- $\ell$  identity and  $\vec{id}'$  be a level- $\ell'$  identity. We say that  $\vec{id}'$  is a sub-identity of  $\vec{id}$  if  $\vec{id}$  is a prefix of  $\vec{id}'$ . If moreover  $\ell' = \ell + 1$ , we say  $\vec{id}'$  is a immediate sub-identity of  $\vec{id}$ .

**Deriving subkeys** Given a master secret key of a HIBS, it is possible to derive secret keys corresponding to level- $\ell$  identities for any  $\ell$ , by running **KeyGen**  $\ell$  times. By a similar procedure, given any secret key corresponding to a level- $\ell$  identity  $\vec{id}$ , it is possible to derive any “subkeys” thereof, i.e., secret keys for sub-identities of  $\vec{id}$ . For our construction, it is useful to name this (simple) procedure: we define  $\text{HIBS.KeyGen}^*$  in Algorithm 1. We write the randomness  $\rho_1, \dots, \rho_\ell$  of  $\text{HIBS.KeyGen}^*$  explicitly.

---

**Algorithm 1**  $\text{HIBS.KeyGen}^*$

---

**Input:**  $sk, \ell, \vec{id} = (id_1, \dots, id_\ell)$   $\triangleright$  Require:  $\ell \leq \ell'$   
**Randomness:**  $\rho_1, \dots, \rho_{\ell'}$   
  **for**  $j = \ell + 1, \dots, \ell'$  **do**  
     $sk \leftarrow \text{HIBS.KeyGen}(sk, id_j; \rho_j)$   
**return**  $sk$

---

**HIBS security requirements** Definition 5 gives only syntax and not security requirements. For a formal security definition, see, e.g., [30]. Informally, an HIBS must satisfy the following:

- *Correctness:* For any identity tuple  $\vec{id}$ , an honestly produced signature w.r.t  $\vec{id}$  must verify as valid w.r.t.  $\vec{id}$ .

---

<sup>9</sup>Technically, all four algorithms take  $1^\kappa$  as an input, and  $\mathcal{M}$  and  $\mathcal{I}$  may be parametrized by  $\kappa$ . For brevity, we leave this implicit except in **Setup**.

- *Unforgeability*: For any PPT adversary  $\mathcal{A}$  with access to a  $\text{KeyGen}(msk, \cdot, \cdot)$  oracle, the probability that  $\mathcal{A}$  outputs a valid signature w.r.t. an identity  $\vec{id} \notin Q$  must be negligible, where  $Q$  is the set of all sub-identities of identities  $\mathcal{A}$  has queried to the oracle.

---

**Algorithm 2** Compress
 

---

**Input:**  $\mathcal{I} = \{\mathcal{I}_\ell\}_{\ell \in [L]}, T \subseteq \mathcal{I}_1 \times \dots \times \mathcal{I}_\ell$   
 // Remove redundant sub-identities  
**for all**  $\tau \in T$  **do**  
   **if**  $\exists \tau' \in T$  s.t.  $\tau'$  is a prefix of  $\tau$  **then**  
      $T = T \setminus \{\tau\}$   
 // Replace identities with prefix identities where possible  
**for**  $\ell = L - 1, \dots, 1$  **do**  
   **for all**  $\vec{\tau} = (\tau_1, \dots, \tau_\ell) \in \mathcal{I}_1 \times \dots \times \mathcal{I}_\ell$  **do**  
     //  $X$  represents all level- $(\ell + 1)$  sub-identities of  $\vec{\tau}$   
      $X = \{(\tau_1, \dots, \tau_\ell, \tau')\}_{\tau' \in \mathcal{I}_{\ell+1}}$   
     **if**  $X \subseteq T$  **then**  
        $T = T \setminus X$   
        $T = T \cup \{(\tau_1, \dots, \tau_\ell)\}$   
**return**  $sk$

---

**Succinctly representing expiry information** Given any set  $T$  of tuples of identities, the simplest way to make signatures with respect to  $T$  forgeable would be to release the secret key corresponding to each  $\vec{id} \in T$ , much as in BasicFFS:

$$\eta = \left\{ sk_{\vec{id}} = \text{HIBS.KeyGen}^*(msk, 0, \vec{id}) \right\}_{\vec{id} \in T} . \quad (1)$$

However, leveraging the hierarchical nature of HIBS,  $\eta$  can often be represented more succinctly than (1). Based on the fact that Algorithm 1 allows the derivation of any subkey, we make two optimizations. First, before computing (1), we delete from  $T$  any  $\vec{id} \in T$  that is a sub-identity of some  $\vec{id}' \in T$ . Secondly, if there is any  $\vec{id}' = (id_1, \dots, id_\ell) \in \mathcal{I}_1 \times \dots \times \mathcal{I}_\ell$  such that every immediate subkey of  $\vec{id}'$  is in  $T$ , i.e.,

$$\forall id_{\ell+1} \in \mathcal{I}_{\ell+1}, (id_1, \dots, id_\ell, id_{\ell+1}) \in T ,$$

then all sub-identities of  $\vec{id}'$  can be removed from  $T$  and replaced by  $\vec{id}'$  before computing (1). Such replacement is permissible only when *every* possible subkey of  $\vec{id}'$  is derivable from  $T$ : otherwise, adding  $\vec{id}'$  to  $T$  would implicate additional subkeys beyond those originally in  $T$ .

These two optimizations yield an algorithm **Compress**, which takes as input a set of identity tuples  $T$ , and outputs a (weakly) smaller set of identity tuples  $T'$  such that knowledge of the secret keys corresponding to  $T'$  enables computing valid signatures with respect to exactly the identity tuples in  $T$ . HIBS security guarantees that even given  $T'$ , signatures for identity tuples not in  $T$  remain unforgeable. Next, we describe how **Compress** works using a tree-based representation of identity tuples.

**Tree representation** It is convenient to think of identity tuples represented graphically in a tree. A node at depth  $\ell$  represents a tuple of  $\ell$  identities (considering the root node to be at depth 0).

The set of all depth- $\ell$  nodes corresponds to the set of all  $\ell$ -tuples of identities. The branching factor at level  $\ell$  is  $|\mathcal{I}_{\ell+1}|$ . Given a secret key for a particular node (i.e., identity tuple), the secret keys of all its descendant nodes are easily computable using  $\text{HIBS.KeyGen}^*$ . (The secret key for the root node is the master secret key.) In this language, **Compress** simply takes a set  $T$  of nodes and returns the smallest set  $T'$  of nodes such that (1) all nodes in  $T$  are descendants of some node in  $T'$  and (2) no node not in  $T$  is a descendant of any node in  $T'$ . Figure 3 gives an illustration of the **Compress** algorithm on a small example tree.

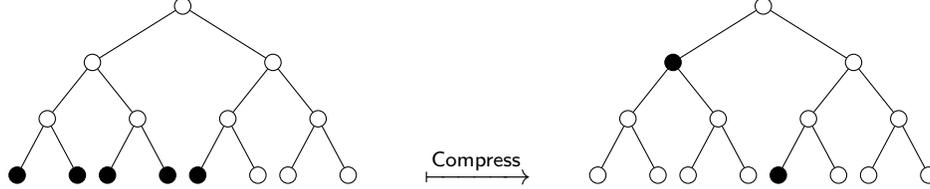


Figure 3: Example application of **Compress**

Our construction of FFS based on HIBS follows. It makes use of Algorithms 1 and 2, which were just defined.

**Construction 1.** Let HIBS be a depth- $L$  HIBS<sup>10</sup> with message space  $\mathcal{M}$  and identity space  $\mathcal{I} = \{\mathcal{I}_\ell\}_{\ell \in [L]}$ . Let  $\mathcal{O}$  be a random oracle,<sup>11</sup> and for any tuple  $\vec{\tau} = (\tau_1, \dots, \tau_\ell)$ , let  $\vec{\mathcal{O}}(\vec{\tau}) = (\mathcal{O}(\tau_1), \dots, \mathcal{O}(\tau_\ell))$ . For  $\ell \in [L]$ , define  $\mathcal{T}_\ell = \mathcal{I}_1 \times \dots \times \mathcal{I}_\ell$ . We construct a FFS  $\Sigma$  with message space  $\mathcal{M}$  and tag space  $\mathcal{T} = \bigcup_{\ell \in [L]} \mathcal{T}_\ell$ , as follows.

- $\Sigma.\text{KeyGen}(1^\kappa)$ : output  $(vk, sk) \leftarrow \text{HIBS.Setup}(1^\kappa)$ .
- $\Sigma.\text{Sign}(sk, \vec{\tau} = (\tau_1, \dots, \tau_\ell), m)$ : let

$$sk_{\vec{\tau}} = \text{HIBS.KeyGen}^*(sk, 0, \vec{\tau}; \vec{\mathcal{O}}(\vec{\tau}))$$

and output  $\sigma \leftarrow \text{HIBS.Sign}(sk_{\vec{\tau}}, m)$ .

- $\Sigma.\text{Verify}(vk, \vec{\tau}, m, \sigma)$ : output  $b \leftarrow \text{HIBS.Verify}(vk, \vec{\tau}, m, \sigma)$ .
- $\Sigma.\text{Expire}(sk, T)$ : let  $T' = \text{Compress}(\mathcal{I}, T)$ ; output

$$\eta = \left\{ (\vec{\tau}, sk_{\vec{\tau}}) : sk_{\vec{\tau}} = \text{HIBS.KeyGen}^*(sk, 0, \vec{\tau}; \vec{\mathcal{O}}(\vec{\tau})) \right\}_{\vec{\tau} \in T'}$$

- $\Sigma.\text{Forge}(\eta, \tau, m)$ : if there exists  $sk_{\tau'}$  such that  $(\tau', sk_{\tau'}) \in \eta$  and  $\tau'$  is a prefix of  $\tau$ , let  $\ell$  be the length of  $\tau'$ , let

$$sk_{\vec{\tau}} = \text{HIBS.KeyGen}^*(sk_{\tau'}, \ell, \vec{\tau}; \vec{\mathcal{O}}(\vec{\tau}))$$

and output  $\sigma \leftarrow \text{HIBS.Sign}(sk_{\vec{\tau}}, m)$ ; otherwise, output  $\perp$ .

**Theorem 1.** If HIBS is a secure HIBS, Construction 1 instantiated with HIBS is a forward-forgeable signature scheme.

<sup>10</sup>The depth need not be finite, but we consider finite  $L$  for simplicity.

<sup>11</sup>The construction is presented in the random oracle model for simplicity, but does not *require* a random oracle: the random oracle can be replaced straightforwardly by a pseudorandom function (PRF) where the PRF key is made part of the HIBS secret key.

*Proof.* Correctness and unforgeability of Construction 1 follow directly from correctness and unforgeability of the underlying HIBS. The FFS requirement of *forgeability on expiry* moreover follows from the correctness requirement of the HIBS: the Forge algorithm of Construction 1 invokes HIBS.Sign using a secret key  $sk_\tau$  which is guaranteed (by construction of HIBS.KeyGen<sup>\*</sup>) to be the secret key corresponding to identity tuple  $\tau$ . The validity of HIBS.Sign invoked on a valid secret key corresponding to identity tuple  $\tau$  is guaranteed by the correctness of the HIBS.  $\square$

**Lemma 1** (Logarithmic succinctness). *Let HIBS be a depth- $L$  HIBS with message space  $\mathcal{M}$  and identity space  $\mathcal{I} = \{\mathcal{I}_\ell\}_{\ell \in [L]}$ . For each  $\ell \in L$ , let  $\preceq_\ell$  be a total order on  $\mathcal{I}_\ell$ . Let  $\mathcal{T}_L = \mathcal{I}_1 \times \cdots \times \mathcal{I}_L$ . For  $i \in |\mathcal{T}_L|$ , let  $\tau_i$  denote the  $i$ th element of  $\mathcal{T}_L$  in the lexicographic order induced by  $\{\preceq_\ell\}_{\ell \in L}$ . Construction 1 instantiated with HIBS is  $(S, 2z)$ -succinct and also  $(S_1, z)$ -succinct, where*

$$S = \{\{\tau_i\}_{j' \leq i \leq j} : j, j' \in [|\mathcal{T}_L|]\} \quad \text{and} \quad S_1 = \{\{\tau_i\}_{1 \leq i \leq j} : j \in [|\mathcal{T}_L|]\} .$$

and  $z(\cdot) = B \cdot \log_B(\cdot)$  where  $B = \max_{\ell \in L} \{|\mathcal{I}_\ell|\}$ . We assume  $B$  is constant.

*Proof.* Fix any  $j, j' \in [|\mathcal{T}_L|]$  and any set  $T = \{\tau_i\}_{j \leq i \leq j'}$ . By the definition of succinctness, it suffices to show that the output of Compress on  $T$ , is a set of nodes of size at most  $2B \cdot \log(|T|)$ .

For any identity tuple  $\iota \in \mathcal{I}$ , let  $\text{Sub}_\iota$  be the set of all level- $L$  identities of which  $\iota$  is a prefix. We say that  $T$  covers  $\iota$  if  $\text{Sub}_\iota \subseteq T$ . Let  $\text{Cover}_T$  be the set of all identities covered by  $T$ . We say  $T$  subsumes an identity  $\iota$  if  $\iota$  is a descendant of some  $\iota' \in \text{Cover}_T$  such that  $\iota' \neq \iota$ . By construction of Algorithm 2 (Compress), any identity subsumed by  $T$  will not be in the output set of Compress( $T$ ) (specifically, it will be removed in the innermost for-loop of Algorithm 2).

For any  $\ell \in [L]$ , consider any consecutive sequence of  $s$  level- $\ell$  identities covered by  $T$ . By definition of lexicographic ordering, there are fewer than  $2B$  level- $\ell$  identities in the sequence that are not subsumed by  $T$  (these identities will be at the beginning and/or end of the sequence). Moreover, if the sequence begins at the smallest level- $\ell$  identity in the lexicographic order, then there are fewer than  $B$  identities in the sequence that are not subsumed by  $T$  (these identities will be at the end of the sequence).

Thus, for each  $\ell \in [L]$ , there are fewer than  $2B$  level- $\ell$  identities in the output set of Compress( $T$ ). Moreover,  $L \leq \log_B(|\mathcal{T}_L|) \leq \log_B(T)$ . Therefore, the output set size is at most  $2B \log_B(T)$ . Moreover, if  $T \in S_1$ , then the output set size is at most  $B \log_B(T)$ .  $\square$

**Discussion of alternative approaches** As discussed above, forward-secure signatures (FSS) and FFS are different primitives with distinct requirements. One could build a FFS from a FSS by computing a long list of secret keys and then using them *in backwards order*. Using techniques of [36, 24], a sequence of keys could moreover be stored with logarithmic storage and computation to access a key. However, this optimization is only designed for contiguous sequences of keys; HIBS-based schemes allow for some succinct non-sequential key release and thus support more nuanced tag structures. Still, for certain applications, e.g., postquantum sequential key release, an FFS based on a FSS such as XMSS [13] could be useful.

The requirements of FFS also have some similarity to *timed authentication*. The TESLA timed authentication protocol [49, 50] considers releasing authentication (MAC) keys following a delay after sending the payload, in the broadcast authentication context. Such delayed verification is untenable for email for several reasons, even beyond the inconvenience of waiting 15 minutes for email delivery. Email's store-and-forward nature (see §2.2) means multiple MTAs may need to verify emails before forwarding (e.g., for spam filtering): if the first MTA waits to verify before forwarding,

the next MTA will be unable to verify because the delay has rendered the authentication forgeable. Also, the inability to discard incoming spam before a time delay may increase denial-of-service vulnerability, especially for smaller email providers.

## 5 Our Protocol Proposals

Sections 5.1 and 5.2 respectively describe our two proposed systems KeyForge and TimeForge.

### 5.1 KeyForge

KeyForge consists of two components: (1) replace the digital signature scheme used in DKIM with a succinct FFS; and (2) email servers periodically publish expiry information.

#### 5.1.1 FFS configuration for KeyForge

Figure 4 illustrates KeyForge’s key hierarchy. KeyForge is based on an  $L$ -level tag structure, corresponding to identity space  $\mathcal{I} = \{\mathcal{I}_\ell\}_{\ell \in [L]}$  where the level- $L$  identities represent 15-minute time chunks spanning a 2-year period. We use the following intuitive 4-level configuration for ease of exposition, but as discussed in §6, it is preferable for efficiency to keep  $|\mathcal{I}_\ell|$  equal for all  $\ell \in [L]$ .

$\mathcal{I}_1 = \{1, 2\}$	representing a 2-year time span
$\mathcal{I}_2 = \{1, \dots, 12\}$	representing months in a year
$\mathcal{I}_3 = \{1, \dots, 31\}$	representing days in a month
$\mathcal{I}_4 = \{1, \dots, 96\}$	representing 15-minute chunks of a day

A tag  $\tau = (y, m, d, c) \in \mathcal{I}_1 \times \mathcal{I}_2 \times \mathcal{I}_3 \times \mathcal{I}_4$  corresponds to a 15-minute chunk of time. The 15-minute chunks are contiguous, consecutive, and disjoint, so that any given timestamp is contained in exactly one chunk.  $\tau(t)$  denotes the unique 4-tuple tag  $(y, m, d, c)$  that represents the chunk of time containing a timestamp  $t$ , and  $t \sqsubset \tau$  denotes that  $\tau$  represents a chunk of time containing timestamp  $t$ .

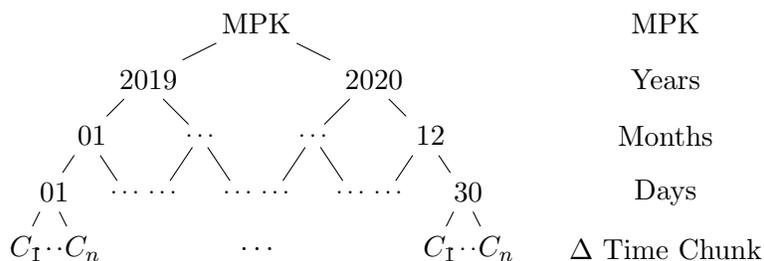


Figure 4: KeyForge Hierarchy Layout

KeyForge requires each signature at time  $t$  to be with respect to a tag (timestamp)  $t + \hat{\Delta}$ . The tag is sent in the email’s header, and used for verification at the receiving server. Algorithm 3 specifies KeyForge’s signing and verification.

**Efficient tree regeneration from private keys.** A key feature of our FFS construction is that the private keys from children (e.g., day-keys) are easy to generate from parent keys (e.g., the

---

**Algorithm 3** KeyForge.Sign and KeyForge.Verify

---

```
t = CurrentTime()
function KeyForge.Sign(sk, m, Δ)
    return FFS.Sign(sk, τ(t + Δ), m)
function KeyForge.Verify(vk, τ, m, σ)
    return t ⊆ τ AND FFS.Verify(vk, τ, m, σ)
```

---

MSK). This is not implied by the definition of HIBS,<sup>12</sup> and is essential for succinct expiry of entire portions of the tree (e.g., a year) by disseminating a single key. Further, regeneration can enhance security and availability: to limit key exposure, organizations could store the MSK in an HSM disconnected from the Internet, and keep only a child key pair in the MSA, thereby mitigating damage in case of compromise and allowing recovery from failure.

**Where to publish expiry information?** Regeneration allows KeyForge to have succinct expiry information; the number of private keys necessary to represent all expired chunks depends on the tree’s structure (see §6), but amounts to less than 4 KB for reasonable configurations.

Small expiry information means ease of distribution. While our implementation uses a simple public-facing webserver, one could imagine posting via DNS TXT records, public blockchains, or in outgoing email headers. Slow but permanent techniques (e.g., a blockchain) for keys higher up in the hierarchy (e.g., a year) could ensure that such keys are permanently available.

**When to publish expiry information?** KeyForge requires email servers to publish expiry information at regular intervals. A natural option is to publish expiry information every 15 minutes; to publish the expiry information corresponding to each chunk  $c$  at the end of the time period that  $c$  represents.

Publishing every 15 minutes yields the finest granularity of expiry possible under the basic four-level tag structure. Based on a server’s preference, it could release information at longer intervals (e.g., days) or shorter ones. In case of an attack, an adversary would be able to convince third parties of the authenticity of all emails in the current interval (e.g., the current day), so risk aversion prefers shorter intervals.

Server misconfiguration and clock skew may cause minor clock discrepancies between the MSA and MDA. To account for this, we delay publishing “expired” keys by 5 minutes. Although in practice most emails are received very quickly, the SMTP RFC [38] has a very lax give-up time of 4 days. To get a rough idea of how quickly emails tend to be delivered, we computed the time differences from the first Received header to the last in the Podesta email corpus [65],<sup>13</sup> and found that, of the 48,246 messages with parseable Received timestamps, over 99% (47,349) took less than 12 minutes.

While expiry time is a configurable parameter of KeyForge (e.g., by administrators), keeping it short is advisable to minimize time until universal forgeability. We leave a detailed study of email delivery times in practice to future work, while noting that such a study might support considerably reducing our conservative 15 minutes, and/or tailoring our approach to specific delay-prone situations. For example, delays are often caused by expected receiving-server outages (e.g.,

---

<sup>12</sup>The definition of HIBS is compatible with this property, but does not require it. Constructions typically have randomized subkey generation processes so do not have reproducible child keys.

<sup>13</sup>Beyond the irony, we chose the Podesta email corpus as it was distributed intact with attachments, and thus arguably more representative of a realistic user’s email distribution than other public datasets.

for server updates), which might be resolved by using a DMARC-like DNS record to signal to the sender to hold messages until later. Anti-spam techniques such as greylisting can delay email by 15 minutes more than usual; to address this, we can add 17 minutes’ leeway when first sending to a new domain.

We do not fully detail remediation procedures for timeouts, but note that similar authentication failures happen under DKIM and are commonly resolved via feedback loops such as Authentication Failure Reporting [28]. Shortening our expiry time is tricky given potentially adversarial routing delays: providing TCP-like flow control would be systematically possible, but we should also account for malicious MTAs trying to prolong messages’ unforgeability. A hard-cutoff maximum would likely be advisable.

**Why 15-minute chunks?** The time period associated with each leaf node is the maximum granularity of expiry information release.  $\hat{\Delta}$  is a lower bound on chunk size: since  $\hat{\Delta}$  represents email delivery time, publishing expiry information more often does not make sense.

**Why a 2-year public key lifetime?** Rotating keys is good practice; for operational reasons, the Messaging, Malware, and Mobile Anti-Abuse Working Group (M3AAWG) recommends DKIM key rotation every 6 months [40]. However, recognizing that, realistically, DKIM keys often last more than 6 months, our evaluation assumes a 2-year period.

**How many levels?** The optimal  $L$  depends on a trade-off between computation time and expiry succinctness; see §6.

**Flexible expiry policies** The basic tag structure described above is customizable: e.g., an extra level  $\mathcal{I}_\ell$  might represent an email’s “sensitivity,” allowing sensitive emails to expire faster. Alternatively, one might want certain emails to expire more slowly or never (e.g., bank/employer emails or contracts). In KeyForge, sensitivity can be expressed as the desired delay until expiry. KeyForge is highly configurable: after the first four levels, different email servers’ policies need not be consistent. Verification refers only to the first four levels of the tag (when checking  $t \sqsubseteq \tau$ ), so a sending server can add more levels beyond the basic four without sacrificing compatibility, to match its desired expiry policy.

## 5.2 TimeForge

KeyForge’s main limitation is that it requires signers to continuously release key material. Wide distribution can pose a practical challenge; users must depend on their provider to perform this task reliably. Unreliable distribution would limit a system’s realistic deniability.

TimeForge takes a different approach that eliminates reliance on follow-up action by signers. TimeForge uses a service that we term a *publicly verifiable timekeeper* (PVTK), which could be realized using various extant Internet systems. A PVTK is a global service that maintains a monotonically increasing clock. At any clock time  $t$ , any party can query the PVTK to obtain a publicly verifiable proof  $\pi_t$  that the current time is at least  $t$ . Simultaneously, the system ensures that attackers cannot forge such proofs at an earlier time.

The intuition behind TimeForge is straightforward. Let  $M$  be an email message sent at time period  $t$ . The sender first signs each message using a standard SUF-CMA signature scheme to produce a signature  $\sigma$ . She then authenticates the message, not directly using  $\sigma$ , but rather using a witness indistinguishable and non-interactive proof-of-knowledge (PoK) of the (informal) statement:

I know a valid sender signature  $\sigma$  on  $M$   
OR

I know a valid PVTk proof  $\pi_{t+d}$ , for some  $d \geq \Delta$ .

Assuming a trustworthy PVTk service, this proof authenticates the message during any time period prior to  $t + \Delta$ . Once a PVTk proof  $\pi_{t+\Delta}$  becomes public, the PoK becomes trivial for any party to generate. Witness indistinguishability ensures that a signer’s valid proof is indistinguishable from a “forgery” later computed using a revealed PVTk proof.

**Publicly verifiable timekeeping.** A PVTk scheme comprises three algorithms.

- $\text{TK.Setup}(1^\lambda)$  takes a security parameter  $\lambda$  and outputs a set of public parameters  $params$  and a trapdoor  $sk$ .
- $\text{TK.Prove}(sk, t)$  takes as input  $sk$  and the current time epoch  $t$ , and outputs a proof  $\pi_t$ .
- $\text{TK.Verify}(params, t, \pi_t)$  on input  $params$ , a time period  $t$ , and the proof  $\pi_t$ , outputs whether  $\pi_t$  is valid.

*Correctness and Security.* Correctness is straightforward.  $\Delta$ -PVTk security requires that an adversary with a PVTk oracle (which provides proofs for arbitrary time periods  $t$ ) must not be able to produce a valid proof for some time period  $t_{\max} + \Delta$  (except with negligible probability) where  $t_{\max}$  is the largest oracle query, and  $\Delta > 0$  is a constant parameter.

**Realizing a PVTk service.** A simple PVTk system can be constructed using a single server that maintains a clock, and periodically signs the current time using an SUF-CMA signature (our implementation does this).

While conceptually simple, deploying this solution at scale is likely to be costly, and may suffer denial-of-service and network-based attacks. A better approach might construct a PVTk from *existing* Internet services. Next, we outline several proposals.

**OCSF servers.** The Online Certificate Status Protocol, in its “stapling” configuration [51] allows TLS servers to obtain a standalone, signed, and timestamped certificate validity message from a Certificate Authority (CA); this can be viewed as an organic implementation of a PVTk server. To avoid reliance on a single CA, users can define the proof  $\pi_t$  to comprise *multiple* valid staples, *e.g.*, one from any  $k$  out of  $N$  chosen CAs. These parameters, as well as the CA identities, can be selected as part of the setup algorithm.

**CT and randomness beacons.** The Certificate Transparency protocol consists of a centrally-managed and publicly verifiable log for recording the issuance of certificates [41]. While CT is not intended as a timestamping protocol, each CT log entry is signed by the log operator (*e.g.*, Google), contains a timestamp, and may be re-purposed to implement a centralized PVTk service. Similarly, NIST operates a randomness beacon [45] that signs and distributes timestamps. While any single centralized service may be unreliable or subject to attack, a (fault-tolerant) combination of these extant services can be used to construct a “composite” PVTk system.

**Proof of work blockchains.** A number of cryptocurrencies use *proof of work* blockchains to construct an ordered transaction ledger [43, 66] which generates new blocks at a rate using intentionally chosen parameters. These ledgers can be used as a form of PVTk system, in which  $params$  comprises some initial block header  $B_s$ , and  $\pi_t$  comprises a list of block headers  $\{B_{s+1}, \dots, B_t\}$  drawn from the blockchain. Although this approach does not produce an exact timekeeping service (block intervals are probabilistic), nor does it guarantee cryptographic unforgeability (as chains can be forged given control of a substantial fraction of the network’s hash power), this is likely not a problem for the short intervals used in TimeForge.

**VDFs and puzzles.** Cryptographic “puzzles” are mathematical problems that require a known (or statistically predictable) number of computational operations to solve. Examples include cryptocurrency proof of work systems and timelock encryption schemes [53]. A related primitive, the *Verifiable Delay Function* [9, 52] creates a *sequential* puzzle that requires a precisely-known amount of work to solve, and allows the solver to produce a proof of the solution’s correctness. While puzzles and VDFs do not directly allow for the creation of a PVTk system, they enable a related primitive: at time  $t$  a sender may generate a puzzle challenge  $M$  (e.g., the contents of an email message) such that a proof  $\pi_{t+\Delta}$  can be found by applying a computational process to  $M$  in expected time approximately  $\Delta$ .

**A basic TimeForge signature scheme.** The TimeForge scheme consists of four algorithms: TF.Keygen, TF.Sign and TF.Verify, and TF.Forge. We assume a PVTk scheme with parameters  $params$  and an SUF-CMA signing algorithm Sig.

- TF.Keygen( $1^\lambda, params$ ). Run Sig.Keygen( $1^\lambda$ ) to generate  $(pk, sk)$  and output  $PK = (pk, params)$ , and  $SK = sk$ .
- TF.Sign( $PK, SK, M, t, \Delta$ ). Parse  $PK = (pk, params)$ . On input a message  $M$  and a time period  $t$ , compute  $\sigma \leftarrow \text{Sig.Sign}(SK, M || t || \Delta)$  and the following witness-indistinguishable (WI) non-interactive PoK:<sup>14</sup>

$$\Pi = \text{NIPoK}\{(\sigma, s, \pi) : \text{Sig.Verify}(pk, \sigma, M || t || \Delta) = 1 \vee (\text{TK.Verify}(params, \pi, s) = 1 \wedge s \geq t + \Delta)\}$$

Note that the prover can produce this proof using  $(\sigma, \perp, \perp)$  as the witness. Output  $\sigma_{\text{tf}} = (\Pi, t, \Delta)$ .

- TF.Verify( $PK, M, \sigma_{\text{tf}}$ ). Parse  $PK = (pk, params)$  and  $\sigma_{\text{tf}} = (\Pi, t, \Delta)$ , verify the proof  $\Pi$  with respect to the public values  $t, \Delta, pk, M$ , and output the verification result.

TF.Forge takes as input a PVTk proof  $\pi_s$  for some time period  $s \geq t + \Delta$ .

- TF.Forge( $PK, M, t, s, \Delta, \pi_s$ ). parse  $PK = (pk, params)$  and compute the NIPoK  $\Pi$  described in the TF.Sign algorithm, using the witness  $(\perp, s, \pi_s)$ . Output  $\sigma_{\text{tf}} = (\Pi, t, \Delta)$ .

*Defining Security.* Security for TimeForge is defined according to the following experiment. This experiment can be considered a variant of (weak) UF-CMA security definition for a signature scheme: an attacker must attempt to forge a TimeForge proof over a message  $M$  that she has not previously queried to a signing oracle. To assist in this, the attacker is given access to not one, but two oracles. The first is a signing oracle for the TimeForge signature scheme, and produces valid signatures for tuples of the form  $(M', t', \Delta')$ . The main difference from the standard UF-CMA experiment is the existence of a second oracle that models the PVTk service. To model this, the attacker additionally obtains PVTk parameters  $params$  at the start of the experiment, and may repeatedly query the PVTk oracle on chosen epoch numbers  $s$  to obtain PVTk proofs of the form  $\pi_s$ . Let  $s_{\max}$  be the largest time period queried to the PVTk oracle at the conclusion of the experiment. We say the attacker *wins* iff she outputs a message  $M$  and valid TimeForge proof  $\sigma_{\text{tf}} = (\Pi, t, \Delta)$  such that  $s_{\max} < t + \Delta$ . where  $(M, t, \Delta)$  was not queried to the signing oracle. We say that a TimeForge scheme is unforgeable under *chosen timestamp attacks* if  $\forall$  p.p.t. attackers  $\mathcal{A}$ , the adversary has at most a negligible advantage in succeeding at the above experiment.

<sup>14</sup>Here we use Camenisch-Stadler notation, where the witness values are in parentheses () and any remaining values are assumed to be public.

*Remark 2.* The definition above does not prevent forgeries that are “outside the expiration period.” Specifically, an intermediary can intercept a message embedding  $(M, t, \Delta)$  where  $t + \Delta$  is in the future, and author a new message  $M', t', \Delta'$  where  $t + \Delta$  has already been proved by the PVTk oracle. This is explicitly allowed by TimeForge; indeed, it is a goal of the system.

**Theorem 2.** *If (1) the PVTk service uses an SUF-CMA signature scheme, (2) the WI proof system is sound (extractable), and (3) the underlying signature scheme used by TimeForge is SUF-CMA, then the basic TimeForge scheme is secure under chosen timestamp attacks.*

*Proof.* Our proof is by contradiction. Let  $\mathcal{A}$  be an attacker that succeeds with non-negligible advantage in the chosen timestamp experiment. We construct a pair of algorithms  $\mathcal{B}_1, \mathcal{B}_2$  such that that one of the two algorithms (respectively) succeeds with non-negligible advantage in the SUF-CMA game against (1) the signature scheme  $\text{Sig}$ , or (2) the underlying PVTk signature scheme. We now describe the operation of each algorithm.

*An attack on the signature scheme Sig.* In this strategy we construct  $\mathcal{B}_1$ , which conducts the SUF-CMA experiment for the underlying signature scheme  $\text{Sig}$ .  $\mathcal{B}_1$  first obtains a public key  $pk$  from the SUF-CMA challenger. It next uses the PVTk signature scheme’s key generation algorithm to produce a keypair  $(params, sk_{\text{PVTk}})$  for the PVTk service and sends  $PK = (pk, params)$  to  $\mathcal{A}$ .

Each time  $\mathcal{A}$  queries the PVTk oracle on some timestamp  $s$ ,  $\mathcal{B}_1$  implements  $\text{TK.Prove}$  by using  $sk_{\text{PVTk}}$  to sign  $s$  and return the signature  $\pi_s$ . Whenever  $\mathcal{A}$  queries the TimeForge signing oracle on some tuple  $(M, t, \Delta)$ ,  $\mathcal{B}_1$  first queries the SUF-CMA signing oracle to obtain a signature  $\sigma$  on  $M||t||\Delta$ . It then constructs a TimeForge signature by constructing the proof described in the  $\text{TF.Sign}$  algorithm using  $\sigma$  as the satisfying witness, and returns  $\sigma_{\text{tf}}$  to  $\mathcal{A}$ .

When  $\mathcal{A}$  outputs a pair  $(M^*, \sigma_{\text{tf}}^*)$  that satisfies the win conditions of the experiment,  $\mathcal{B}_1$  parses  $\sigma_{\text{tf}}^*$  to obtain  $(\Pi^*, t^*, \Delta^*)$  and runs the extractor on  $\Pi^*$  to obtain the witness  $(\sigma^*, s^*, \pi^*)$ . (If the extractor fails,  $\mathcal{B}_1$  aborts.) If  $\sigma^*$  is a valid signature on  $M^*||t^*||\Delta^*$ , it outputs the pair  $(\sigma^*, M^*)$  as an SUF-CMA forgery.<sup>15</sup>

*An attack on the PVTk scheme.* In this strategy we construct  $\mathcal{B}_2$ , which conducts the SUF-CMA experiment against the PVTk signature scheme. This algorithm proceeds as in Strategy 1, except that here we set  $params$  to be the public key obtained from the SUF-CMA challenger, and generate the keypair  $(pk, sk) = \text{Sig.Keygen}(1^\lambda)$ . Queries to the PVTk oracle are answered by forwarding  $s$  to the SUF-CMA oracle and returning the resulting signature as  $\pi_s$ , and queries to the TimeForge oracle are answered honestly by running  $\text{TF.Sign}$  with  $sk$  as an input. As in the previous strategy, if  $\mathcal{A}$  succeeds in the experiment we run the extractor to obtain  $\sigma^*, s^*, \pi^*$ . Then  $\mathcal{B}_2$  verifies that  $\pi^*$  is a valid signature on  $s^*$  and, if so, outputs  $(\pi^*, s^*)$  as a forgery for the SUF-CMA experiment.

**Analysis.**  $\mathcal{A}$ ’s view is distributed identically when it interacts with  $\mathcal{B}_1$  or  $\mathcal{B}_2$ , so  $\mathcal{A}$ ’s advantage must be identical w.r.t.  $\mathcal{B}_1$  and  $\mathcal{B}_2$ . By soundness, the WI knowledge extractor fails with probability at most negligible in the security parameter. Thus both adversaries will abort with at most negligible probability due to extraction error.

It remains only to show that at least one of the two algorithms above must succeed with non-negligible advantage in the SUF-CMA experiment. This is true because the attacker’s output  $(\Pi^*, t^*, \Delta^*)$  and the extracted witness  $(\sigma^*, s^*, \pi^*)$  must satisfy the conditions that (1) for all PVTk queries  $s$  made during the experiment,  $s < t^* + \Delta^*$  (by the requirements of the experiment), (2) the message  $M^*||t^*||\Delta^*$  has not been queried to the SUF-CMA signing oracle (by the requirements of

<sup>15</sup>By the restrictions on  $\mathcal{A}$ ,  $M^*$  must represent a message that has not previously been queried to the SUF-CMA oracle.

the experiment), and (3) the following conditions are true (by the soundness of the proof system):

$$\begin{aligned} \text{Sig.Verify}(pk, \sigma^*, M^*) &= 1 \vee (\text{TK.Verify}(params, \pi^*, s^*)) \\ &= 1 \wedge s^* \geq t^* + \Delta^* \end{aligned}$$

If  $\mathcal{A}$  succeeds in the TimeForge experiment with non-negligible advantage, then the extracted witness must contain a valid signature  $\sigma^*$  on a message  $M^*||t||\Delta$  not queried to the signing oracle, **or** it must contain a valid signature  $\pi^*$  on some epoch  $s^*$  not queried to the PVTk oracle. Hence, one of  $\mathcal{B}_1$  or  $\mathcal{B}_2$  succeeds with non-negligible advantage.  $\square$

### 5.3 Realizing the TimeForge proof system

TimeForge can be realized using a variety of WI and ZK proof systems, combined with efficient SUF-CMA signature schemes. For example, a number of pairing-based signature schemes [7, 21, 6] admit efficient proofs of knowledge of a signature using simple Schnorr-style proofs [2]. More recent proving systems, e.g., Bulletproofs [14] and zkSNARKs (e.g., [48, 32]), admit succinct proofs of statements involving arbitrary arithmetic circuits and discrete-log relationships. Using the latter schemes ensures short proofs, in the hundreds of bytes, in some cases with a small, constant verification cost. Thus, even complex PVTk proofs such as block header sequences, can potentially be reduced to a succinct TimeForge signature.

**A concrete implementation.** For our basic implementation, which signs a timestamp, we considered several proof systems. For the relatively simple proof statement used in this scheme, Bulletproofs are not appropriate for two reasons: (1) the proof sizes that result exceed 1000 bytes, and (2) these proofs do not natively support efficient signatures. zkSNARKs produce bandwidth-efficient signatures, but at a significant cost due to the need to generate a trusted setup embedding the signature verification circuit. Based on these considerations, we propose and evaluate one concrete implementation based on Schnorr-style proving techniques, made non-interactive using the Fiat-Shamir heuristic. Our approach implements TimeForge using a dedicated server that produces (weak) Boneh-Boyen signatures [10] over the current time period  $t$ , which is encoded as an integer in  $\mathbb{Z}_q$ . Let  $g_1, g_2$  be generators of a pair of bilinear groups  $\mathbb{G}_1, \mathbb{G}_2$  of order  $q$ . Briefly, a Boneh-Boyen signature on a time period  $t$  comprises a single group element  $\sigma = g_1^{1/x+t}$ , where  $x$  represents the signing key, and the server’s public key is  $g_2^x$ . Verification is conducted by checking the following pairing equality:  $e(g_1, g_2) = e(\sigma, g_2^x g_2^t)$ .

Our proposed TimeForge proof of knowledge requires the following components. First, the prover provides a Pedersen commitment  $B$  to the current time period  $T_{\text{current}}$  using randomness  $r$ . The proof also reveals the (alleged) true signing time period  $T_{\text{signing}}$  in cleartext (in case it is different) and attaches  $\delta$ . Using these values, the prover employs the homomorphic property of Pedersen commitments to derive an implicit commitment  $C = g_1^{\gamma = T_{\text{current}} - T_{\text{signing}} - \delta} h^{r'}$ , and then uses a range proof to prove that it knows a value  $\gamma$  that is in the range  $[1, 2^{32}]$ . We use a range proof due to Camenisch, Chaabouni, and Shelat [20]. Alternatively, this proof could be implemented using a Bulletproof, due to Bootle *et al.* [14].

In addition to this commitment proof and range proof, we provide two separate Schnorr-style proofs in an “OR” construction:

1. *A standard Schnorr signature on the message.* This comprises an interactive proof of knowledge of a value  $sk \in \mathbb{Z}_q$  such that  $PK = g^{sk}$ , flattened into a signature of knowledge on

the signed message, using the Fiat-Shamir heuristic. (This represents the genuine signer’s signature on the message.)

2. A proof of knowledge of a Boneh-Boyer signature on the TimeForge time period  $T_{\text{current}}$ , signed using the TimeForge server secret key. For this construction we use a interactive zero-knowledge protocol given by Boneh, Boyen and Shacham [11, Protocol 1], flattened using the same Fiat-Shamir hash function.

## 5.4 KeyForge<sup>+</sup> and TimeForge<sup>+</sup>

KeyForge<sup>+</sup> (resp. TimeForge<sup>+</sup>) consists of KeyForge (resp. TimeForge) with two modifications: a *forge-on-request protocol* and *per-recipient-domain signatures*, described next.

**1. Forge-on-request protocol.** We add a protocol (detailed in Algorithm 4) by which email servers accept real-time requests for specified email content to be sent to the requester (and nobody else). The forge-on-request protocol ensures that all users have the capability to forge emails to themselves in real time, directly achieving *immediate recipient forgeability*. The requirement that the recipient be the requester is crucial: each requester is enabled to forge emails *only to herself*.

The requester’s email server attests to the requesting client’s identity (similarly to DKIM). We note that a malicious server could unauthorizedly sign requests for any client account it controls. This is outside our threat model, and such behavior is equally possible under DKIM (see also “Client-server trust” under §3.2): that is, today’s email ecosystem already relies on servers to attest honestly to their clients’ identities, and allows servers to spam their own clients (a behavior that might not keep them many clients).

The scheme as described so far already works for single-recipient emails, but can be problematic for multi-recipient<sup>16</sup> emails due to spam/spoofing attacks between co-recipients: as long as all recipients of an email receive exactly the same information, an attacker can always request a forged email addressed to a group including herself, and quickly forward it to trick the others into treating the email as legitimate. To remedy this, we require one more change to KeyForge.

**2. Per-recipient-domain signatures.** In DKIM, KeyForge, and TimeForge, the sending server signs each outgoing email once. In KeyForge<sup>+</sup> and TimeForge<sup>+</sup>, instead, the sending server signs each outgoing email *once per recipient domain*. That is, the sending server produces a signature  $\sigma_D = \text{Sign}(sk, (D, m))$  for each recipient domain  $D$ , where  $sk$  is the signing key and  $m$  is the email information that the sending server would have signed under DKIM (or KeyForge or TimeForge). The sending server sends each recipient domain  $D$  the email and *only* the signature  $\sigma_D$ .

Per-recipient-domain signatures prevent attackers from using the forge-on-request protocol to send spam/spoofing emails to co-recipients on forged emails. Adida *et al.* [1] previously proposed per-recipient signatures in a very similar context.

**Notation**  $\text{Email}_s(S, R, m, \mu, t)$  is as defined in §3.3, additionally taking into account that signatures in KeyForge<sup>+</sup> and TimeForge<sup>+</sup> are per recipient domain. **FReq** denotes a special message to betoken forge requests. For an email address  $a$ , let  $a.\text{dom}$  denote its domain.

Note that a forge-on-request protocol achieves a stronger guarantee than the use of ring signatures (i.e., signing with respect to both the sending and receiving servers’ public keys). A forge-on-request protocol enables any recipient with the ability to send mail from an email server to forge mail from any sender to herself. The ring signature approach enables her to do this only if she has the ability to sign fraudulent mail with the receiving server’s secret key.

<sup>16</sup>Here, “recipients” means any recipients whether via **to**, **cc**, or **bcc**.

---

**Algorithm 4** Forge-on-request

---

**Requester**

To request an email with message  $m$  and metadata  $\mu$  from `alice@foo.com`:

- Send  $(\text{FReq}, m, \mu, \text{alice@foo.com})$  to client's (i.e., its own) email server.

**Email server** (say, `bar.com`, with secret key  $s$ )

On receiving request  $(\text{FReq}, m, \mu, a)$  from own client `bob`:<sup>17</sup>

- If  $a.\text{dom} = \text{bar.com}$ :<sup>18</sup> Let  $t$  be the current time. Deliver  $e$  to `bob`, where  $e = \text{Email}_s(a, \text{bob@bar.com}, m, \mu, t)$ .

- Else: Let  $\sigma = \text{Sign}(\text{FReq}, m, \mu, a, \text{bob})$ .<sup>19</sup> Send  $(\text{FReq}, m, \mu, a, \text{bob}, \sigma)$  to server  $a.\text{dom}$ .

On receiving request  $(\text{FReq}, m, \mu, a, b, \sigma)$  from server  $b.\text{dom}$ :

- $v \leftarrow \text{Verify}(vk, (\text{FReq}, m, \mu, a, b), \sigma)$ , where  $pk$  is  $b.\text{dom}$ 's public key in DNS.

- If  $v = 0$ : Do not respond.

- Else (i.e.,  $v = 1$ ): Let  $t$  be the current time. Send  $e, e'$  to  $b.\text{dom}$ , where  $e = \text{Email}_s(a, b, m, \mu, t)$  and  $e' = \text{Email}_s(a, b, m, \mu, t - \hat{\Delta})$ .
- 

**Theorem 3.** *KeyForge<sup>+</sup> is non-attributable for recipients (Definition 1) and  $\hat{\Delta}$ -universally non-attributable (Definition 2).*

*Proof.* Follows directly from Lemmata 2 and 3. □

**Lemma 2.** *KeyForge<sup>+</sup> is non-attributable for recipients (Def. 1).*

*Proof.* Recall from Definition 1 that we must show that there is a PPT simulator  $\mathcal{S}$  such that for any sender  $S$  and recipient  $R$ , for any email message  $m$  and metadata  $\mu$ ,

$$\text{KeyForge}_{sk}(S, R, m, \mu, t) \approx_c \mathcal{S}^R(S, m, \mu), \quad (2)$$

where  $sk$  is the (master) secret key of  $S$ ,  $t$  is the time at which  $\mathcal{S}$  is invoked, and the superscript  $R$  denotes that  $\mathcal{S}$  has access to the recipient's email server. We construct  $\mathcal{S}$  in Algorithm 5.

---

**Algorithm 5** Simulator  $\mathcal{S}$  for recipient non-attributability

---

**Input:**  $S, m, \mu$

$t = \text{CurrentTime}()$

**send** forge request  $(m, \mu)$  to  $S$

**receive** answer  $\{e_0, e_1\}$

**parse**  $e_0, e_1$  as emails w.r.t. tags  $\tau_0, \tau_1$  respectively

**if**  $\tau(t) = \tau_0$  **then return**  $e_0$

**else if**  $\tau(t) = \tau_1$  **then return**  $e_1$

---

By definition of  $\hat{\Delta}$ , we know  $S$  received the request at some time  $t' \leq t + \hat{\Delta}$ . Thus, by construction of KeyForge, the emails  $e, e'$  must be signed with respect to the tags  $\tau(t' - \hat{\Delta}), \tau(t')$  (say, respectively). It follows that  $\tau(t) = \tau(t' - \hat{\Delta})$  or  $\tau(t) = \tau(t')$ . Therefore, at least one of the if-conditions in Algorithm 5 must be satisfied, and  $\mathcal{S}$  always produces an output. By construction of the if-statements and the forge-on-request protocol, the output of  $\mathcal{S}$  is an email signed by  $S$  for a tag corresponding to timestamp  $t$ , as (2) requires. Indeed, we achieve equality of distributions (not just indistinguishability). □

**Lemma 3.** *KeyForge<sup>+</sup> is  $\hat{\Delta}$ -universally non-attributable (Def. 2).*

*Proof.* Recall from Definition 2 that we must show that there is a PPT simulator  $\mathcal{S}$  such that for any sender  $S$  (with secret key  $sk$ ) and recipient  $R$ , for any email message  $m$ , metadata  $\mu$ , and timestamp  $t$ , the following holds at any time  $\geq t + \hat{\Delta}$ :

$$\text{KeyForge}_{sk}(S, R, m, \mu, t) \approx_c \mathcal{S}(S, R, m, \mu, t) . \quad (3)$$

Let  $\text{KeyForge}_{sk^*}^*$  be identical to  $\text{KeyForge}_{sk}$  except that whenever  $\text{KeyForge}_{sk}$  invokes  $\text{KeyForge.Sign}(sk, \cdot)$ ,  $\text{KeyForge}_{sk^*}^*$  instead invokes  $\text{HIBS.Sign}(sk^*, \cdot)$ . Next, we construct  $\mathcal{S}$  in Algorithm 6.

---

**Algorithm 6** Simulator  $\mathcal{S}$  for  $\hat{\Delta}$ -strong non-attributability

---

**Input:**  $S, R, m, \mu, t$

**retrieve** published expiry information  $\eta$  for  $S$

**for all**  $(\vec{\tau}, sk_{\vec{\tau}}) \in \eta$  **do**

**if**  $t \sqsubseteq \vec{\tau}$  **then**

$sk^* \leftarrow \text{HIBS.KeyGen}^*(sk_{\vec{\tau}}, \ell, \tau(t); \vec{O}(\vec{\tau}))$

**return**  $\text{KeyForge}_{sk^*}^*(S, R, m, \mu, t)$

---

Since  $\mathcal{S}$  is invoked at time  $\geq t + \hat{\Delta}$ , and  $\text{KeyForge}$  prescribes publication of expiry information at the end of each chunk of duration  $\hat{\Delta}$ , the expiry information  $\eta$  retrieved by  $\mathcal{S}$  includes expiry information with respect to time  $t$ . Therefore, the if-condition in Algorithm 6 is satisfied for at least one element of  $\eta$ .<sup>20</sup>

Recall that  $\text{KeyForge.Sign}$  invokes  $\text{FFS.Sign}$ , which invokes  $\text{HIBS.Sign}$ . By definition, if  $t \sqsubseteq \vec{\tau}$  and  $\eta$  is expiry information with respect to  $sk$ , then  $sk^*$  as computed in Algorithm 6 is the same key used to invoke  $\text{HIBS.Sign}$  (within  $\text{FFS.Sign}$ , which is within  $\text{KeyForge.Sign}$ ) at time  $t$ . Therefore, the output distributions of  $\text{KeyForge}_{sk^*}^*$  and  $\text{KeyForge}_{sk}$  are identical. It follows that  $\mathcal{S}$  satisfies (3).  $\square$

**On the efficiency of KeyForge<sup>+</sup>/TimeForge<sup>+</sup>** Per-recipient-domain signatures add sender-side overhead compared to single-signature schemes like DKIM, KeyForge, or TimeForge: the sending server must compute one signature per recipient domain per email, whereas the receiving server verifies just one signature, just like before. While this additional computation is unlikely to be prohibitive given the efficiency of signing, it must be taken into account when evaluating KeyForge<sup>+</sup> and TimeForge<sup>+</sup>'s efficiency, as discussed further in Section 6.

Implementing forge-on-request and per-recipient-domain signatures would entail more complexity and significant changes to the existing email infrastructure, than the base protocols. While immediate recipient forgeability is desirable for added protection against real-time attacks (see Threat Model 2), KeyForge is a more realistic candidate for near-term deployment as it is realizable with lighter-weight changes to the existing system: namely, replacing DKIM's signature scheme, and unilateral server publication of small amounts of data.

## 6 Implementation and Evaluation

We implemented prototypes of KeyForge and TimeForge and integrated them into Postfix, a common MDA/MSA. The entire project consists of roughly 2,000 lines of Go, C, and Python, and

---

<sup>20</sup>In fact, it will be satisfied for exactly one element of  $\eta$ , by construction of **Compress** which ensures that no timestamp is represented by more than one element.

Monthly KeyForge <sub>B</sub>	KeyForge <sub>B</sub> $\sigma$	Monthly KeyForge <sub>C</sub>	KeyForge <sub>C</sub> $\sigma$	DKIM RSA2048 $\sigma$	TimeForge $\sigma$
$30 \times 65 = 1950$	98	$30 \times (64 + 32) = 2880$	$64 \times 2 + 32 = 160$	256	841

Table 1: Bandwidth costs (in bytes) of KeyForge<sub>B</sub>, KeyForge<sub>C</sub>, and DKIM with RSA.  $\sigma$  denotes a signature.

is available open source.<sup>21</sup> Full details of cryptographic primitives and curve parameters are in Appendix D.

We evaluate two versions of KeyForge instantiated with different HIBS schemes: (1) KeyForge<sub>B</sub>, which uses Gentry-Silverberg’s “BasicHIDE” bilinear map based scheme [30] using a BN254 curve and (2) KeyForge<sub>C</sub>, which uses certificate chains on public keys using non-identity-based signatures, instantiated with Ed25519.<sup>22</sup> We also implemented a prototype of TimeForge (see Appendix 5.3), which is less efficient; it is intended as a proof of concept whose practicality will improve with advances in the underlying proof primitives (an active area of research). The two KeyForge implementations share the following bandwidth optimization.

**KeyForge bandwidth optimization.** HIBS schemes tend to have relatively large signatures. In KeyForge<sub>B</sub>, a signature must include public parameters for each node on the path to the current chunk. A public parameter in this configuration is 65B, yielding a bandwidth of 260B for a four-level Y/M/D/Chunk tree, resulting in a total of 293B per signature. KeyForge<sub>C</sub> similarly requires an Ed25519 signature between each node in the hierarchy, and has total signature size of 448B (four 64B path signatures, four 32B public keys, and the message signature). We optimize bandwidth by precomputing all path parameters except for the last chunk and store them in the DNS, along with the MPK. When verifying from a new server, KeyForge performs a DNS lookup and caches the result at a cost of 2-3KB per month (see Table 1).

Two components, the keyserver and mail filter, are shared between all implementations. They are described next.

**Mail Filter.** The filter ensures that sent emails are properly formatted, verifies incoming emails, and communicates the results to the MDA/MTA. The filter works by intercepting SMTP requests, adding necessary metadata to outgoing email headers, and requesting cryptographic operations from the keyserver. When sending a message, the filter attaches an expiry time (and other verification information) to the email’s header, hashes the metadata and message content, forwards the hash to the keyserver to sign, and finally adds this signature to the header. On receipt, the filter confirms that the signature’s hash matches the message and metadata, and forwards the signature, sending domain, and expiry timestamp to the keyserver for verification. If verification fails, the filter alerts PostFix and the message is dropped.

**Keyserver.** The keyserver performs signing and verification, communicates with the mail filter over RPC, and publishes expired keys (for KeyForge) via a simple webserver.

## 6.1 Evaluation

We evaluate messaging bandwidth, expiry data bandwidth, and speed. Our primary focus is on comparison with RSA-2048: it is the signature scheme commonly used in DKIM, and so a natural benchmark for practicality in the current email ecosystem. Although more bandwidth-efficient algorithms were approved for DKIM use some months ago, (e.g., Ed25519 with a 64 B signature [35]),

<sup>21</sup><https://github.com/keyforgery/KeyForge>

<sup>22</sup>The certificate-based approach has been attributed to folklore.

	Sign(ms)	Sign/s	Verify(ms)	Verify/s
TimeForge	24.58	49.68	23.24	43
KeyForge <sub>B</sub>	0.34	2,932	3.36	298
KeyForge <sub>C</sub>	0.13	17,197	0.13	7,541
RSA2048	0.93	1,075	0.05	19,966
Ed25519	0.03	27,001	0.10	9,781

Table 2: Time required for a single operation in milliseconds, and the equivalent number of operations per second. KeyForge times are for a 4-level tree, RSA is from OpenSSL benchmarks.

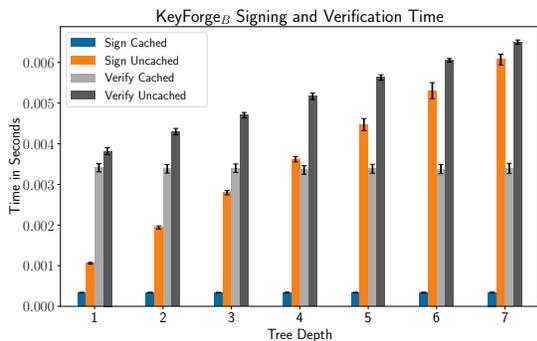


Figure 5: KeyForge<sub>B</sub> timings

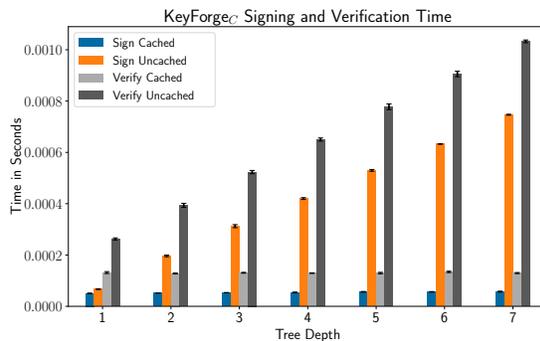


Figure 6: KeyForge<sub>C</sub> timings

these schemes appear to have had limited deployment to date.<sup>23</sup> Nonetheless, for completeness, this section also considers Ed25519 performance.

**Bandwidth.** Table 1 shows bandwidth costs for various configurations of KeyForge and TimeForge. Both KeyForge implementations have a bandwidth per email that is 42% *smaller* than a DKIM RSA-2048 signature.

**Speed.** To capture the range of KeyForge’s possible performance, we considered two cases: (1) where the public key path is verified from scratch (e.g., in setting up a new server, or verifying messages from a new domain) and (2) where path parameters are pre-verified and cached. Figures 5 and 6 show the results. Signing is largely unaffected by tree depth when caching.

Table 2 provides efficiency microbenchmarks for KeyForge, TimeForge, and Ed25519 and DKIM’s RSA-2048 via the OpenSSL suite’s benchmark. All KeyForge benchmarks are for a 4-level tree with caching. Our experiments were run on a laptop with power lower than a common server (see Appendix D), so our timings may be seen as upper bounds. Performance scales linearly with the number of cores; our measurements are for a single core.

**Optimizing for KeyForge expiration bandwidth.** While the Y/M/D/Chunk configuration is easy to intuit, an equal branching factor across tree levels yields a large gain in succinctness. Table 4 (in Appendix E) shows the average and maximum size of expiry info of various depth trees with an equal branching factor: e.g., the average expiry size for a 2-year period is 4.5MB, 4KB, or 1.8KB for depths 1, 4, and 7.

**Discussion and analysis.** We find that KeyForge, especially KeyForge<sub>C</sub>, is likely practical when

<sup>23</sup>E.g., as of October 2019, Gmail and Exchange use only RSA-2048.

using DKIM’s RSA-2048 as a benchmark. In both implementations, KeyForge’s signing time is *better* than RSA: KeyForge<sub>B</sub> and KeyForge<sub>C</sub> sign 2.7 and 16 times faster than RSA, respectively. KeyForge further beats RSA on signature bandwidth per email, at just 63% or less of RSA signature size in the worst case. RSA outperforms KeyForge only on verification time: KeyForge<sub>C</sub> is still eminently practical, with verification a factor of two slower than RSA, whereas KeyForge<sub>B</sub> is an order of magnitude slower.

Verification time is unlikely to affect KeyForge’s viability, as other factors such as hashing, I/O, and network latency are likely to dominate. Any hash-and-sign scheme must read the message into memory and perform a hash, so to provide a ballpark measurement of I/O and hashing, we timed OpenSSL’s SHA256 on the Podesta corpus [65], stored on-disk. The average time required was 10.2ms (2.689ms std),<sup>24</sup> indicating that hashing and I/O is surprisingly impactful. Network latency is significant as well — SMTP requires that a sending MTA perform a *minimum* of four round trips per email.<sup>25</sup> A highly optimistic round-trip time of 5ms would yield of 20ms per email, not including time to send message content.

The choice between KeyForge<sub>B</sub> and KeyForge<sub>C</sub> is likely implementation dependent: while KeyForge<sub>B</sub> requires less bandwidth, its drawbacks are speed and use of non-IETF-standardized curves (unlike KeyForge<sub>C</sub>).

**A note on adoption.** With an ecosystem as unwieldy as email, a reasonable concern might be that any large-scale update would be difficult. That said, now is an opportune time to propose such changes: the IETF has recently approved a new standard that will encourage MTAs to begin updating their DKIM signing and verification algorithms [35]. Further, if the community were to endorse a new standard, one could imagine large email providers (e.g., Google) displaying favorable security indicators akin to Gmail’s TLS indicators[31]. Such tactics have been successful in the context of HTTPS.

We have consulted members of the IETF, W3C, and the Gmail Security team, and optimized and evaluated our prototypes with their performance priorities and concerns in mind.

## Acknowledgements

We are grateful to Jon Callas for helpful discussions about motivations for email non-attributability and our scheme’s applicability to DKIM, and to Dan Boneh, Daniel J. Weitzner, John Hess, Bradley Sturt, Stuart Babcock, and Ran Canetti for their feedback on earlier versions of this work. This work was supported in part by the William and Flora Hewlett Foundation grant 2014-1601, and by the MIT Media Lab’s Digital Currency Initiative and its funders. We would like to acknowledge support from the National Science Foundation under awards CNS-1653110 and CNS-1801479, and a Google Security & Privacy Award.

## References

- [1] Ben Adida, David Chau, Susan Hohenberger, and Ronald L. Rivest. Lightweight email signatures. In *International Conference on Security and Cryptography for Networks*, pages 288–302.

---

<sup>24</sup>Email size is often pushed up by HTML formatting, embedded media, and attachments. Average email size in our corpus is 98 KB (691 KB std).

<sup>25</sup>SMTP messages require a round trip per command, and each email requires a MAIL, RCPT, and two DATA commands.

- Springer, 2006.
- [2] Ben Adida, Susan Hohenberger, and Ronald L. Rivest. Lightweight encryption for email. In *Steps to Reducing Unwanted Traffic on the Internet Workshop, SRUTI'05*. USENIX Association, 2005.
  - [3] D. F. Aranha and C. P. L. Gouvêa. RELIC is an Efficient Library for Cryptography. <https://github.com/relic-toolkit/relic>.
  - [4] Diego F. Aranha, Laura Fuentes-Castañeda, Edward Knapp, Alfred Menezes, and Francisco Rodríguez-Henríquez. Implementing pairings at the 192-bit security level. In *International Conference on Pairing-Based Cryptography*, pages 177–195. Springer, 2012.
  - [5] Associated Press. DKIM verification script. Available at <https://github.com/associatedpress/verify-dkim>.
  - [6] Mira Belenkiy, Melissa Chase, Markulf Kohlweiss, and Anna Lysyanskaya. P-signatures and noninteractive anonymous credentials. In *TCC 2008*, 2008.
  - [7] Mira Belenkiy, Melissa Chase, Markulf Kohlweiss, and Anna Lysyanskaya. Compact E-Cash and Simulatable VRFs Revisited. In *Pairing-Based Cryptography '09*, 2009.
  - [8] Steven Michael Bellovin. Spamming, phishing, authentication, and privacy. 2004.
  - [9] Dan Boneh, Joseph Bonneau, Benedikt Bünz, and Ben Fisch. Verifiable delay functions. In *Advances in Cryptology - CRYPTO 2018 - 38th Annual International Cryptology Conference, Santa Barbara, CA, USA, August 19-23, 2018, Proceedings, Part I*, pages 757–788, 2018.
  - [10] Dan Boneh and Xavier Boyen. Short signatures without random oracles. In Christian Cachin and Jan L. Camenisch, editors, *Advances in Cryptology - EUROCRYPT 2004*, pages 56–73, Berlin, Heidelberg, 2004. Springer Berlin Heidelberg.
  - [11] Dan Boneh, Xavier Boyen, and Hovav Shacham. Short group signatures. In Matt Franklin, editor, *Advances in Cryptology - CRYPTO 2004*, pages 41–55, Berlin, Heidelberg, 2004. Springer Berlin Heidelberg.
  - [12] Nikita Borisov, Ian Goldberg, and Eric Brewer. Off-the-record communication, or, why not to use PGP. In *Proceedings of the 2004 ACM Workshop on Privacy in the Electronic Society, WPES '04*, pages 77–84, New York, NY, USA, 2004. ACM.
  - [13] Johannes A. Buchmann, Erik Dahmen, and Andreas Hülsing. XMSS - A practical forward secure signature scheme based on minimal security assumptions. In Bo-Yin Yang, editor, *Post-Quantum Cryptography - 4th International Workshop, PQCrypto 2011, Taipei, Taiwan, November 29 - December 2, 2011. Proceedings*, volume 7071 of *Lecture Notes in Computer Science*, pages 117–129. Springer, 2011.
  - [14] B. Bunz, J. Bootle, D. Boneh, A. Poelstra, P. Wuille, and G. Maxwell. Bulletproofs: Short proofs for confidential transactions and more. In *2018 IEEE Symposium on Security and Privacy (SP)*, volume 00, pages 319–338.

- [15] Jack Burbank, David Mills, and William Kasch. Network Time Protocol Version 4: Protocol and Algorithms Specification. <https://tools.ietf.org/html/rfc5905> [<https://perma.cc/428T-HN3Y>].
- [16] John Callas, Eric Allman, Miles Libbey, Michael Thomas, Mark Delany, and Jim Fenton. DomainKeys Identified Mail (DKIM) Signatures.
- [17] Jon Callas. [ietf-dkim] Thinking about DKIM and surveillance. Available online at: [https://mailarchive.ietf.org/arch/msg/ietf-dkim/eWKbWdYmkX\\_d2ki\\_lAbczVSj8qY](https://mailarchive.ietf.org/arch/msg/ietf-dkim/eWKbWdYmkX_d2ki_lAbczVSj8qY) [<https://perma.cc/DQF6-SQNZ>].
- [18] Jon Callas. [ietf-dkim] DKIM Key Sizes, October 2016. <http://mipassoc.org/pipermail/ietf-dkim/2016q4/017195.html> [<https://perma.cc/7NNX-QJUK>].
- [19] Jon Callas. [ietf-dkim] DKIM Key Sizes, October 2016. <http://mipassoc.org/pipermail/ietf-dkim/2016q4/017207.html> [<https://perma.cc/K8LM-KJS7>].
- [20] Jan Camenisch, Rafik Chaabouni, and abhi shelat. Efficient protocols for set membership and range proofs. In Josef Pieprzyk, editor, *Advances in Cryptology - ASIACRYPT 2008*, pages 234–252, Berlin, Heidelberg, 2008. Springer Berlin Heidelberg.
- [21] Jan Camenisch and Anna Lysyanskaya. Signature schemes and anonymous credentials from bilinear maps. In *Advances in Cryptology-CRYPTO 2004*, 2004.
- [22] Ran Canetti, Cynthia Dwork, Moni Naor, and Rafail Ostrovsky. Deniable encryption. In Burton S. Kaliski, editor, *Advances in Cryptology — CRYPTO '97*, pages 90–104, Berlin, Heidelberg, 1997. Springer Berlin Heidelberg.
- [23] Charles Cazabon. getmail version 5. <http://pyropus.ca/software/getmail>.
- [24] Don Coppersmith and Markus Jakobsson. Almost optimal hash sequence traversal. In Matt Blaze, editor, *Financial Cryptography, 6th International Conference, FC 2002, Southampton, Bermuda, March 11-14, 2002, Revised Papers*, volume 2357 of *Lecture Notes in Computer Science*, pages 102–119. Springer, 2002.
- [25] D. Crocker. Internet Mail Architecture, 2009. <https://tools.ietf.org/html/rfc5598>.
- [26] Mario Di Raimondo and Rosario Gennaro. New approaches for deniable authentication. *Journal of Cryptology*, 22(4):572–615, Oct 2009.
- [27] Cynthia Dwork, Moni Naor, and Amit Sahai. Concurrent zero-knowledge. In *Proceedings of the Thirtieth Annual ACM Symposium on Theory of Computing*, STOC '98, pages 409–418, New York, NY, USA, 1998. ACM.
- [28] Hilda L. Fontana. Authentication Failure Reporting Using the Abuse Reporting Format. <https://tools.ietf.org/html/rfc6591> [<https://perma.cc/5MTF-ZD8P>].
- [29] Center for Strategic and International Studies (CSIS). Significant cyber incidents, 2018. <https://www.csis.org/programs/cybersecurity-and-governance/technology-policy-program/other-projects-cybersecurity>.

- [30] Craig Gentry and Alice Silverberg. Hierarchical ID-based cryptography. In Yuliang Zheng, editor, *Proceedings of ASIACRYPT 2002*, volume 2501 of *Lecture Notes in Computer Science*, pages 548–566. Springer, 2002.
- [31] Google. Making email safer for you, February 2016.
- [32] Jens Groth. On the size of pairing-based non-interactive arguments. In Marc Fischlin and Jean-Sébastien Coron, editors, *Advances in Cryptology – EUROCRYPT 2016*, pages 305–326, Berlin, Heidelberg, 2016. Springer Berlin Heidelberg.
- [33] Lachlan J. Gunn, Ricardo Vieitez Parra, and N. Asokan. Circumventing cryptographic deniability with remote attestation, 2019.
- [34] HIPAA Journal. United hospital district phishing attack impacts 2,143 patients, 2019. <https://www.hipaajournal.com/united-hospital-district-phishing-attack-impacts-2143-patients/>.
- [35] J. Levine. RFC 8463 - A New Cryptographic Signature Method for DomainKeys Identified Mail (DKIM).
- [36] Markus Jakobsson. Fractal hash sequence representation and traversal. *Proceedings of the 2002 IEEE International Symposium on Information Theory (ISIT)*, pages 437–44, 2002.
- [37] Markus Jakobsson, Kazue Sako, and Russell Impagliazzo. Designated verifier proofs and their applications. In *Proceedings of the 15th Annual International Conference on Theory and Application of Cryptographic Techniques*, EUROCRYPT’96, pages 143–154, Berlin, Heidelberg, 1996. Springer-Verlag.
- [38] John Klensin. RFC5321: Simple Mail Transfer Protocol, 2008.
- [39] Murray Kucherawy and Elizabeth Zwicky. Domain-based Message Authentication, Reporting, and Conformance (DMARC). <https://tools.ietf.org/html/rfc7489>.
- [40] Kurt Andersen. M3aawg DKIM Key Rotation Best Common Practices | M3aawg, March 2019. <http://www.m3aawg.org/DKIMKeyRotation> [March 2019 version archived at: <https://perma.cc/4WY6-SH8K>].
- [41] Ben Laurie, Adam Langley, and Emilia Kasper. Certificate Transparency. RFC 6962, June 2013.
- [42] Jeremy B. Merrill. Authenticating Email Using DKIM and ARC, or How We Analyzed the Kasowitz Emails. *ProPublica*, July 2017.
- [43] Satoshi Nakamoto. Bitcoin: A peer-to-peer electronic cash system, 2009.
- [44] Moni Naor. Deniable ring authentication. In *Advances in Cryptology — CRYPTO 2002*, pages 481–498, Berlin, Heidelberg, 2002. Springer Berlin Heidelberg.
- [45] NIST. NIST randomness beacon. <https://www.nist.gov/programs-projects/nist-randomness-beacon>.

- [46] Emad Omara, Benjamin Beurdouche, Jon Millican, Raphael Robert, Katriel Cohn-Gordon, and Richard Barnes. The Messaging Layer Security (MLS) Protocol. <https://tools.ietf.org/html/draft-ietf-mls-protocol-07> [<https://perma.cc/YQ5X-36LB>].
- [47] Open Whisper Systems. Curve.java, 2018. <https://github.com/signalapp/libsignal-protocol-java/blob/master/java/src/main/java/org/whispersystems/libsignal/ecc/Curve.java> [<https://perma.cc/6CBK-745E>].
- [48] Bryan Parno, Jon Howell, Craig Gentry, and Mariana Raykova. Pinocchio: Nearly practical verifiable computation. In *Proceedings of the IEEE Symposium on Security and Privacy*. IEEE, May 2013.
- [49] Adrian Perrig, Ran Canetti, J.D. Tygar, and Dawn Song. The TESLA broadcast authentication protocol. *RSA CryptoBytes*, 5:2–13, 2002. Available at: [https://people.eecs.berkeley.edu/~tygar/papers/TESLA\\_broadcast\\_authentication\\_protocol.pdf](https://people.eecs.berkeley.edu/~tygar/papers/TESLA_broadcast_authentication_protocol.pdf).
- [50] Adrian Perrig, Dawn Song, Ran Canetti, J. D. Tygar, and Bob Briscoe. Timed efficient stream loss-tolerant authentication (TESLA): multicast source authentication transform introduction. *RFC*, 4082:1–22, 2005.
- [51] Yngve N. Pettersen. The Transport Layer Security (TLS) Multiple Certificate Status Request Extension. RFC 6961, June 2013.
- [52] Krzysztof Pietrzak. Simple verifiable delay functions. *IACR Cryptology ePrint Archive*, 2018:627, 2018.
- [53] R. L. Rivest, A. Shamir, and D. A. Wagner. Time-lock puzzles and timed-release crypto. Technical report, Cambridge, MA, USA, 1996.
- [54] Ronald L. Rivest, Adi Shamir, and Yael Tauman. How to leak a secret. In *ASIACRYPT*, 2001.
- [55] Raphael Satter. Emails: Lawyer who met Trump Jr. tied to Russian officials. *The Associated Press*, July 2018.
- [56] Adi Shamir. Identity-based cryptosystems and signature schemes. In G. R. Blakley and David Chaum, editors, *Proceedings of CRYPTO '84*, volume 196 of *Lecture Notes in Computer Science*, pages 47–53. Springer, 1984.
- [57] Jonathan Stempel and Jim Finkle. Yahoo says all three billion accounts hacked in 2013 data theft, 2017. <https://www.reuters.com/article/us-yahoo-cyber/yahoo-says-all-three-billion-accounts-hacked-in-2013-data-theft-idUSKCN1C8201>.
- [58] Michael Thomas. Requirements for a DomainKeys Identified Mail (DKIM) Signing Practices Protocol. <https://tools.ietf.org/html/rfc5016>.
- [59] Nik Unger and Ian Goldberg. Deniable key exchanges for secure messaging. In *Proceedings of the 22Nd ACM SIGSAC Conference on Computer and Communications Security, CCS '15*, pages 1211–1223, New York, NY, USA, 2015. ACM.
- [60] Nik Unger and Ian Goldberg. Improved strongly deniable authenticated key exchanges for secure messaging. *PoPETs*, 2018(1):21–66, 2018.

- [61] WhatsApp. WhatsApp encryption overview: Technical white paper. <https://www.whatsapp.com/security/WhatsApp-Security-Whitepaper.pdf> [<https://perma.cc/6L5W-C8S5>].
- [62] Wikileaks. Sony Email Leak. Available at <https://wikileaks.org/sony/emails/>, 2012.
- [63] Wikileaks. The Global Intelligence Files: STRATFOR email leak. Available at [https://wikileaks.org/gifiles/docs/13/1328496\\_stratfor-.html](https://wikileaks.org/gifiles/docs/13/1328496_stratfor-.html), 2012.
- [64] Wikileaks. Search the DNC Database, July 2016. <https://wikileaks.org/dnc-emails/>].
- [65] Wikileaks. WikiLeaks: DKIM Verification, nov 2016. <https://wikileaks.org/DKIM-Verification.html> [<https://perma.cc/H3SR-YB44>].
- [66] Gavin Wood. Ethereum: A secure decentralised generalised transaction ledger. *Ethereum project yellow paper*, 151:1–32, 2014.

## A Comparison Table

Table 3: Comparison of Messaging Schemes

	Messaging			Email	
	OTR v3	Signal	MLS	DKIM	KeyForge / TimeForge
ATTRIBUTABLE?		non-attributable		attributable	non-attributable
INTERACTIVE?	yes (see Remark 3)				no
KNOWN ENDPOINTS?	yes	yes	yes		no
SYNCHRONOUS?	yes	no	yes	no ( $\Leftarrow$ unknown endpoints)	
ROUTING	central server(s)	central server(s)	decentralized w/ known endpoints	fully decentralized w/ unknown endpoints	
FEDERATED?	no ( $\Leftarrow$ central server)	no ( $\Leftarrow$ central server)	yes	yes	
AUTHENTICATION BANDWIDTH	464B*	64B (Ed25519) <sup>†</sup>	64B [46]	256B (RSA2048)	KeyForge: 146B TimeForge: 1096B
BANDWIDTH PER USER $n = \#$ in group chat	$N/A^{\ddagger}$	$O(n^2)^{\S}$	$O(n^2)$	$O(1)$	

\* This is the size of the initial Diffie-Hellman handshake (which would amortize over multiple messages), not including any per-message costs.

<sup>†</sup> The exact signature scheme is implementation-dependent. This data reflects both the current Signal implementation and the best bandwidth among implementations of which we are aware. [47]

<sup>‡</sup> OTR does not support group chat.

<sup>§</sup> WhatsApp has an optimization that reduces message traffic; published documentation leaves unclear if this affects authentication handshake bandwidth. [61]

Table 3 lays out characteristics of KeyForge, DKIM-based email and other types of non-attributable messaging schemes. As also explained in Section 2, non-attributable email presents a very different set of constraints from other types of non-email messaging schemes, so direct performance comparisons between these settings are not meaningful. Table 3 highlights in more detail how the email setting differs from other types of messaging.

*Remark 3.* All other non-attributable messaging schemes of which we are aware require some sort of interactivity — either by having the sender know the endpoint a priori or by having the receiver communicate back to the sender before a message is sent. While Signal’s updates to OTR allow it to be more asynchronous than some other protocols, Signal is still interactive in that it requires communication of key material to known endpoints. Such interactivity makes these messaging schemes inappropriate for use in email.

## B ARC

Authenticated Received Chain (ARC) is an experimental and (at the time of writing) largely unimplemented IETF standard that aims to address the issues caused by indirect email flows. Legitimate modification of messages in transit may occur in a number of circumstances. For example, an MTA that is also a virus scanner may remove malicious attachments, and mailing lists may prepend a name to the subject of an email. Unfortunately, any alteration of the body or headers invalidates the original signature. ARC acts as an attestation by an intermediary a subset of DKIM, SPF, DMARC, or previous ARC signatures were verifiable before content was modified. To do so, the intermediary adds its own signature to the header of the message, along with metadata about what was verified.

**Complications from ARC.** Arc would pose minor complications for non-attributability if widely adopted, as third-party MTAs sign emails in transit. Though here presented as a modification of DKIM, KeyForge and TimeForge can be easily extended to accommodate ARC; the third-party signer uses an FFS for signing, publishes expiry information, and offers a forge-on-request protocol.

## C Preventing real-time attacks requires interaction

**Claim 1.** *In the store-and-forward model (see §2.2), where final-destination addresses and keys may be unknown to senders, any email protocol that (1) proves (with soundness) the sending domain’s identity to recipients and (2) is secure against real-time attacks must involve interaction, i.e., recipient communication to the sending server after receipt of an initial message.*<sup>26</sup>

*Proof (sketch).* By assumption, there is a verification procedure  $V$  that the receiver  $R$  may run on the single-message protocol transcript, and  $V$  will (with overwhelming probability) output 1 if the claimed sending domain is correct and 0 otherwise. Since the final destination may be unknown to the sender,  $V$ ’s behavior cannot depend on any information associated specifically with the receiver, such as key material. That is,  $V$ ’s inputs must consist only of the protocol transcript and public information. It follows that  $R$  may share the received message(s) with any third party  $T$ , who then becomes able to run  $V$  for itself on the same inputs as  $R$  would use. By soundness,  $V$ ’s output on these inputs must suffice to convince  $T$  of the sending domain’s identity *unless*  $R$  has the ability to non-interactively forge protocol transcripts that appear to be from this sending domain, in real time — say, using an algorithm  $F$ . But because the final destination may be unknown to the sender, the behavior of any such  $F$  algorithm cannot depend on any information associated specifically with the receiver, so the  $F$  algorithm could also be run by other parties to successfully forge protocol transcripts that appear to be from this sending domain. This contradicts soundness (Condition (1) from the theorem statement).  $\square$

Among other things, this rules out immediate recipient forgeability by timed-authentication approaches like TESLA [49, 50] (also mentioned in §4 under “Discussion of alternative approaches”).

---

<sup>26</sup>In fact, a version of Claim 1 extends to any *unidirectional* protocol, which may involve multiple messages but all sent by the same party, as long as there is realistic nontrivial variance in network delays. For simplicity, the proof sketch is for single-message protocols.

## D Compilation and Evaluation parameters

We performed all benchmarks on a 2017 MacBook Pro, 15-inch, with an Intel 4-core 3.1GHz processor and 16GB of RAM. We use the RELIC toolkit’s [3] implementation of a BN-245 curve. This configuration conservatively yields keys with a 110-bit security level [4], which is on par with the standard 2048-bit RSA. We chose RELIC due to its support for many pairing friendly curves and low overhead.

Openssl was compiled with the following flags:

```
clang -I. -I.. -I../include -fPIC -fno-common -DOPENSSL_PIC -DOPENSSL_THREADS -D_REENTRANT
-DHAVE_DLFCN_H -arch x86_64 -O3 -DL_ENDIAN -DGHASH_ASM -DECP_NISTZ256_ASM -DDSO_DLFCN
-Wall -DOPENSSL_IA32_SSE2 -DOPENSSL_BN_ASM_MONT -DOPENSSL_BN_ASM_MONT5 -DOPENSSL_BN_ASM_GF2m
-DSHA1_ASM -DSHA256_ASM -DSHA512_ASM -DMD5_ASM -DAES_ASM -DVPAES_ASM -DBSAES_ASM -DWHIRLPOOL_ASM
```

RELIC was compiled with the following flags:

```
cmake ../ -DALLOC=DYNAMIC -DFP_PRIME=381 -DARITH=gmp-sec -DWSIZE=64 -DEP_SUPER=off
-DFP_METHD="INTEG;INTEG;INTEG;MONTY;LOWER;SLIDE" -DCOMP="-O3 -mtune=native -march=native"
-DFP_PMERS=off -DFP_QNRES=on -DPP_METHD="LAZYR;OATEP" -DFPX_METHD="INTEG;INTEG;LAZYR"
```

## E Size of Expiry Information

Table 4 shows the size of expiry information.

Table 4: Expiry information size

$L$	Expiry info size (bytes)			
	1 year		2 years	
	Avg	Max	Avg	Max
1	1121248	2242496	1679814	4485056
2	12700	25344	16934	33792
3	3283	6464	3920	7744
4	1787	3520	2016	3968
5	1275	2496	1408	2752
6	1048	2048	1117	2176
7	859	1664	934	1792