

# Intrusion-Tolerant and Confidentiality-Preserving Publish/Subscribe Messaging

Sisi Duan, Chao Liu, Xin Wang, Yusen Wu, Shuai Xu, Yelena Yesha, Haibin Zhang  
University of Maryland, Baltimore County

**Abstract**—We present Chios, an intrusion-tolerant publish/subscribe system which protects against Byzantine failures. Chios is the first publish/subscribe system achieving decentralized confidentiality with fine-grained access control and strong publication order guarantees. This is in contrast to existing publish/subscribe systems achieving much weaker security and reliability properties.

Chios is flexible and modular, consisting of four fully-fledged publish/subscribe configurations (each designed to meet different goals). We have deployed and evaluated our system in multiple platforms (rack servers, Amazon EC2, IBM Cloud and IBM bare metal servers). We compare Chios with various publish/subscribe systems. Chios is as efficient as an unreplicated, single-broker publish/subscribe implementation, only marginally slower than Kafka and Kafka with passive replication, and at least an order of magnitude faster than all Hyperledger Fabric modules and publish/subscribe systems using Fabric.

## I. INTRODUCTION

Publish/Subscribe (pub/sub) is a popular messaging pattern allowing disseminating information from publishers to different subsets of interested subscribers via an overlay of brokers (servers). Publishers advertise information to the brokers and send publications as advertised. Subscribers express their interests for receiving a subset of publications by issuing subscriptions to brokers. Upon receiving publications from publishers matching the interests of subscribers, brokers send the corresponding publications to the interested subscribers.

One distinguishing feature of a pub/sub system is that it decouples publishers and subscribers in both time and space: publishers and subscribers do not need to know or synchronize with one another. This feature enables both system flexibility and scalability. Pub/Sub systems are widely used in practice, such as Amazon SNS [4], AMQP [54], Apache Kafka [7], FAYE [25], Google Cloud Pub/Sub [30], and MQTT [44]. Pub/Sub serves as the core middleware for numerous applications, e.g., data collection and analysis, Internet-of-Things (IoT), network management and monitoring, streaming services.

Despite their popularity, existing pub/sub systems (built in both industry and academia) suffer reliability and confidentiality problems. Let us illustrate the issues with a health record exchange pub/sub system [31], [24], where the actors include patients and providers (physicians, hospitals, pharmacists), both of which can be publishers and subscribers. Publications may be medical files (e.g., reports, X-ray images) sent from patients or providers to patients or providers. Publications may also be new drug information and updates about the availability of facilities sent from providers to patients. For

instance, an emergency unit receives a patient in critical conditions and disseminates the patient medical files as a publisher to various hospital units, while the hospital units may submit subscriptions (e.g., specialties, qualifications, schedule for patient admission and treatment sessions). As another example, a new-born is identified by a hospital for a rare dermatology disease. The hospital represents the new-born and the parents to send the medical images to some local expert dermatologists for timely treatment.

Consider another example of a market report notification system, where publishers are private sectors publishing paid market reports, and subscribers are investors who receive reports according to their interests (e.g., reports for certain categories, reports for specific periods). The brokers match publications with the interests and send the publications to interested (and paid) investors.

With the examples in mind, we now discuss the challenges of building intrusion-tolerant pub/sub systems.

**Confidentiality and fine-grained access control (Or: Two-way information control).** Publications in both examples (health records, private market reports) need to be confidentiality-protected. In fact, confidentiality in pub/sub is strongly tied to *access control*, a process by which subscribers are granted access to certain publications based upon certain rules. The middleware community has long been expecting pub/sub systems where publishers can define by whom and how their data can be accessed, preferably not just role-based but also attribute-based.

For instance, the “ideal” situation for health record exchange is that publishers (patients and providers on behalf of patients) can decide by whom, when, and how their health records can be viewed or used. Patients should be able to decide which doctors can see their records, either exactly (by name), or those that meet certain criteria (e.g., “D.C. doctors”, “more than 15-year practice in dermatology”, “no malpractice history”) [19]. For the market report example, publishers may enforce access control based on subscribers’ qualifications, attributes, and if subscribers paid for the service (to the brokers).

Confidentiality-preserving pub/sub with fine-grained access control *enhances* the conventional pub/sub systems with *two-way* information control. In conventional pub/sub systems, subscribers can filter the information via subscriptions, but publishers cannot control who can receive the publications. The one-way information control is undesirable for applications such as cross-domain pub/sub systems where publications need to be protected (as shown in the health exchange

and market report examples), most pub/sub systems in private corporate networks, and any IoT and big data applications where individual user data are sensitive.

Achieving the goal *securely*, however, is difficult. Existing pub/sub systems with confidentiality or access control either rely on non-cryptographic trusted domains (an overly strong assumption), centralized architectures, and/or violate the decoupling feature of pub/sub systems [31], [32], [33], [52], [8], [59], [40], [27], [56], [51]. Building a decentralized pub/sub system with fine-grained access control is deemed to be a major open problem [46]. First, the approach to encrypting publications using the keys of subscribers does not work, because, due to the decoupling feature of pub/sub systems, publishers do not know the identities or keys of subscribers. Second, publishers cannot encrypt the data using the keys of brokers either, as brokers would know the publications in plaintext. Even if a single broker is compromised, all historical publications will be leaked.

**Reliability.** Another challenge of building intrusion-tolerant pub/sub systems is reliability under Byzantine (arbitrary) failures. Existing reliable pub/sub systems [20], [36], [38], [39], [48] only achieve weak reliability notions. One particular reliability notion is *publication total order* which guarantees subscribers should receive relevant publications in the same order. For instance, in the stock market, seeing a high price followed by a low price means something very different from seeing a low price followed by a high price; it is vital to ensure that all subscribers receive the price information in the same order. Publication total order would be easy to achieve if brokers use Byzantine fault-tolerant (BFT) state machine replication to maintain a total order of publications and ask subscribers to deliver publications according to the total order.

Even so, due to the two-way control (publication filtering via subscriber interests and publisher access control), not *all* publications will be sent to all subscribers. Therefore, the approach that brokers maintain a total order fail to work. In particular, subscribers do not know if they should wait for or skip publications with certain sequence numbers, as subscribers do not know if the corresponding publications are on the way or will never arrive.

**Discussion.** With the rise of blockchains, two pub/sub systems using blockchains (Hyperpubsub [60] and Trinity [48]) were proposed to defend against (Byzantine) failures. Both systems make a black-box usage of existing pub/sub systems and blockchain systems. Hyperpubsub uses Apache Kafka and Hyperledger Fabric [6], while Trinity combines MQTT and one of the four blockchains (Fabric, Tendermint [53], and test networks for Ethereum [55] and IOTA [34]). The two systems, however, suffer from at least three problems. First, both systems are essentially auditing systems using blockchains. The overall systems are not Byzantine fault-tolerant, as neither Kafka (only partially crash fault-tolerant) or MQTT (not fault-tolerant) can defend against Byzantine failures. Both liveness and safety are violated if any brokers of the two systems are compromised. Second, both Hyperpub and Trinity leveraging fully-fledged blockchains have demon-

strated poor performance, because blockchains are essentially storage systems not designed for pub/sub systems, and many features of blockchains are not needed for pub/sub systems. Third, both Hyperpubsub and Trinity directly combine existing pub/sub and blockchains systems and therefore require a much larger number of nodes and resources than a blockchain system or a conventional pub/sub system.

Neither Hyperpubsub nor Trinity achieves confidentiality or publication total order, two goals we aim to address in this paper.

**Our contribution.** We design, implement, and evaluate Chios, a Byzantine fault-tolerant (BFT) pub/sub system with fine-grained access control and strong reliability, without sacrificing the decoupling property of pub/sub. Chios's security assumption is standard to BFT and threshold cryptography, i.e., an adversary cannot corrupt more than 1/3 of the total brokers. We summarize our contribution in the following:

- We formally define the properties of a BFT and confidentiality-preserving pub/sub system, covering strong access control and message ordering guarantees, in the sense of cryptography and reliable distributed systems.
- We demonstrate Chios is provably secure under our definitions by devising and extending cryptographic and reliable distributed system protocols (e.g., vector-label-input threshold encryption, broadcast encryption with decentralized key distribution). Chios is the first pub/sub system achieving decentralized and fine-grained access control as well as publication total order. We compare Chios with existing pub/sub systems in Table I.
- Chios is versatile and modular, supporting three additional and fully-fledged pub/sub instances designed to meet different goals (e.g., different performance metrics, different application scenarios). This includes an instance that combines threshold encryption and broadcast encryption to enable more efficient and dynamic access control. For the instance, we also provide an optimized instantiation that is more efficient than a trivial instantiation. Both the general protocol and the instantiation use a novel approach to maintain the decoupling property of pub/sub.
- We implement and evaluate Chios, showing that all its variants are nearly as efficient as its unreliable (unreplicated) counterpart and existing pub/sub systems (Kafka and Kafka with passive replication) and orders of magnitude faster than blockchain-based systems (Fabric, Trinity, Hyperpubsub). None of existing pub/sub systems or blockchain-based systems achieve decentralized confidentiality or strong order guarantees.

## II. RELATED WORK

**Fault-tolerant pub/sub.** Most of industry pub/sub systems (Apache Kafka [7], FAYE [25], Google Cloud Pub/Sub [30], and MQTT [44]) do not have strong fault tolerance guarantees. For instance, Kafka is crash fault-tolerant for its controller part. For its broker components, most Kafka implementations are not fault-tolerant, though Kafka can be configured to use passive replication for weak fault tolerance. Pub/Sub systems

| systems          | brief description   | Byzantine publisher | Byzantine broker | confidentiality and access control                           | publication total order | publication liveness |
|------------------|---|---------------------|------------------|--|-------------------------|----------------------|
| Chios            | BFT and confidentiality-preserving favor performance over reliability | ●                   | ●                | decentralized; attribute-based                               | ●                       | ●                    |
| Kafka [7]        |   | ○                   | ○                |  | ○                       | ○                    |
| AMQP [54]        | “pub/sub for business”  | ○                   | ○                | use trusted virtual host; password for access control        | ○                       | ○                    |
| Hyperpubsub [60] | auditing system for Kafka   | ○                   | ○                |  | ○                       | ○                    |
| Trinity [48]     | auditing system for MQTT  | ○                   | ○                |  | ○                       | ○                    |
| P2S [20]         | crash fault-tolerant  | ○                   | ○                |  | ●                       | ●                    |
| PubliyPrime [39] | Byzantine failure detection   | ○                   | ●                |  | ●                       | ●                    |
| JM [36]          | deconstructing BFT using a large number of nodes                      | ●                   | ●                |  | ○                       | ●                    |
| IRC [33]         | access control using ABE  | ●                   | ●                | centralized authority needed; expensive pairing-based crypto | ○                       | ○                    |
| EventGuard [51]  | use trusted components  | ○                   | ●                | trusted nodes for confidentiality                            | ○                       | ○                    |

Table I

CHARACTERISTICS OF REPRESENTATIVE PUB/SUB PROTOCOLS. ● DENOTES PARTIAL SUPPORT. P2S AND PUBIYPRIME ACHIEVE WEAKER ORDERING GUARANTEES THAN PUBLICATION TOTAL ORDER. (THE FORMAL DEFINITIONS OF PUBLICATION TOTAL ORDER AND PUBLICATION LIVENESS ARE IN SEC. III-A.)

with strong reliability have been mostly studied for the case of crash failures [38], [58], [20]. Only a handful of works consider a weaker subset of Byzantine failures [36], [39] and none of them achieve publication total order. Besides, PubliyPrime [39] does not handle Byzantine publishers or subscribers.

**Pub/Sub with payload confidentiality.** Confidentiality in pub/sub systems can be generally divided into two categories [46]: 1) confidentiality for publication headers and subscription constraints; 2) payload confidentiality (the ability to hide the payload of the publications, e.g., the patient health record). The confidentiality issue has become a major obstacle to wider adoption of pub/sub systems [46].

Chios addresses payload confidentiality but not confidentiality for publication headers or subscription constraints. Most prior pub/sub systems that handle payload confidentiality rely on overly strong “trusted domain” assumptions and do not maintain the decoupling feature of pub/sub systems that is essential to pub/sub system flexibility and scalability [8], [59], [40], [27], [56]. Srivatsa and Liu [51] devised EventGuard with many goals similar to ours. EventGuard, however, assumes a trusted service for confidentiality and authenticity.

**Pub/Sub with access control (but no fault tolerance).** While there are a number of pub/sub systems [31], [32], [33], [52] that use attribute-based encryption (ABE) [12] to achieve fine-grained access control, they all suffer from the following problems: 1) Efficient ABE schemes rely on relatively slow pairing-based cryptography. 2) All these systems use a trusted central authority which is a single point of failure. While the so-called *decentralized ABE* schemes exist [43], decentralization here actually means that *anyone* can serve as an ABE authority by creating a public key and issuing private keys to different users, but it does not mean that the keys are generated interactively among distributed nodes.

**Reliable distributed systems with confidentiality.** Several works achieve confidentiality in distributed file or storage systems that support *store* and *retrieve* operations [2], [28], [35],

[42], [17], [10], [47]. In these systems, clients apply encryption, or secret sharing, to the data before the data is uploaded to the system. Notably, Depspace [10] explores how to use publicly verifiable secret sharing and hash function to encrypt and locate client data, but it does not achieve linearizability. Besides, Depspace uses hash function to encrypt data, which makes it vulnerable to brute-force and dictionary attacks. Access control in Depspace is centralized, as the access control certificates are issued by a single administrator. AVID [17] suggests the use of threshold encryption to provide access control for Byzantine reliable broadcast and asynchronous verifiable information dispersal. AVID, however, considers a much simpler setting and does not have an implementation.

Yin et al. [57] built a BFT protocol which privately processes user data by separating agreement from execution and using threshold signatures. Assuming the same architecture, Duan and Zhang [23] provided a more efficient construction that uses only symmetric encryption. Both protocols require a lot more nodes than a conventional BFT protocol.

Many recent works [18], [14], [41] explore how to perform private computation on blockchains using trusted execution environments (TEEs), e.g., Intel SGX. These systems require trusting a single TEE vendor (e.g., Intel). Some cryptographic proposals use zkSNARKs [9] or multi-party computation [21] to achieve private computation. These approaches are limited in practice, as the cost to deal with *generic* operations is very high, and the throughput is low.

### III. SYSTEM AND THREAT MODEL

**Background on pub/sub systems.** Pub/Sub systems enable disseminating information from publishers (information sources) to subscribers (interested recipients) via an overlay of brokers (servers). Publishers advertise information to the brokers and send publications as advertised. Subscribers express their interests for receiving a subset of publications by issuing subscriptions. Brokers store subscriptions received from subscribers. Upon receiving matching publications from

publishers, brokers send the corresponding publications to the interested subscribers. Besides storing subscriptions, brokers may maintain routing tables to deliver subscribers information.

The communication between publishers and subscribers is decoupled both in time and space. In particular, publishers and subscribers do not need to know or synchronize with one another. Indeed, direct communication among end-customers may not be possible. The decoupling feature enables flexible and scalable information exchange and also avoids maintenance and charging difficulties for end-customers. Moreover, this allows anonymity between publishers and subscribers (assuming brokers are correct).

This paper considers *topic-based* pub/sub, which is dominant in industry pub/sub systems (e.g., Kafka, FAYE, MQTT, Amazon SNS, Google Cloud Pub/Sub). In topic-based pub/sub, a publication includes a *header* and a *payload*. The header contains the *topics* and their values (e.g., ID = “Alice”, county = “Orange”, price = “105”), while the payload contains the complete bulk data. Correspondingly, a subscription includes a set of *constraints* on the topics (e.g., ID = “Alice”, county = “Franklin” or “Orange”, price = 100). The brokers need to match publications against stored subscriptions according to the constraints of the topics (“equation,” “and,” “or” for topic-based pub/sub).

**BFT.** We consider BFT state machine replication (SMR) protocols, where  $f$  out of  $n$  replicas may fail arbitrarily (Byzantine failures) and a computationally bounded adversary can coordinate faulty replicas. A replica *delivers operations*, each *submitted* by some client. The client should be able to compute a final response to its submitted operation from the responses it receives from replicas. Correctness of a secure BFT protocol is specified as follows.

- **Agreement:** If any correct replica delivers an operation  $m$ , then every correct replica delivers  $m$ .
- **Total Order:** If a correct replica has delivered  $m_1, m_2, \dots, m_s$  and another correct replica has delivered  $m'_1, m'_2, \dots, m'_{s'}$ , then  $m_i = m'_i$  for  $1 \leq i \leq \min(s, s')$ .
- **Liveness:** If an operation  $m$  is submitted to  $n - f$  correct replicas, then all correct replicas will eventually deliver  $m$ .

#### A. Formalizing BFT Pub/Sub

**Syntax.** In our setting, publishers and subscribers are clients. Publishers can be subscribers and vice versa. We use brokers, servers, and replicas interchangeably. We consider an overlay network, where brokers are connected in a complete graph.

A BFT pub/sub system consists of the following (possibly interactive) operations (reg, advertise, sub, pub, notify, read). An interactive registration algorithm reg is run by clients and brokers. Through the reg algorithm, new clients can be registered in the system and brokers can verify and store client (access) attributes (e.g., ages, certificates) enabling them to have access to publications in the future. For instance, a publisher may want only clients with certain attributes to see its publications. Clients should be able to register independently, and in particular, potential publishers and subscribers need not know one another. A client may not need to decide at this stage

if the client would like to register as a publisher, a subscriber, or both, but rather may do this later via advertise and sub.

Publishers advertise to the replicas information that will be sent to all or a subset of clients. The advertise messages may be viewed as special publications. Subscribers send brokers subscriptions to express their interests via a sub operation. Brokers store subscriptions received from subscribers. Upon receiving matching publications from publishers via a pub operation, brokers send the corresponding publications to the interested subscribers via a notify operation. The read operation is similar to that of popular pub/sub systems (e.g., Kafka) and allows a client to read particular data of interests from brokers.

Operations (reg, advertise, sub, pub) change broker state and are collectively called write operations. Operations (notify, read) do not change broker state.

In our system, a publisher can send an encrypted publication together with access control rules ac to the system. We say a subscriber (a client) is *authorized* to see a publication  $m$ , if the publisher submitting  $m$  has listed the subscriber in its access control rules ac.

**Goals.** The goal of our secure BFT pub/sub system is to achieve CIA (confidentiality, integrity, availability) against malicious brokers, publishers, and subscribers. As in a BFT system, we assume a strong adversary that can passively corrupt  $f$  out of  $n$  replicas and adaptively corrupt an unbounded number of clients. We divide the goals into confidentiality and reliability goals.

*Confidentiality and access control goals.* We provide a unified definition of security covering all confidentiality aspects (access control as specified by data providers and confidentiality for non-subscribers and brokers). Specifically, given a BFT pub/sub system, we associate the following game to an adversary  $\mathcal{A}$  in Fig. 1.

- $\mathcal{A}$  chooses to corrupt a fixed set of  $f$  brokers.
- $\mathcal{A}$  interacts with honest parties arbitrarily and chooses to corrupt clients adaptively.
- $\mathcal{A}$  selects two messages  $m_0$  and  $m_1$ , an ac, and a unique tag tid that specifies an instance, and submits them to the encryption oracle for the system.  $\mathcal{A}$  cannot corrupt any clients specified by ac (otherwise,  $\mathcal{A}$  would have trivially won the game). The oracle randomly selects a bit  $b$  and computes an encryption  $c$  of  $m_b$  with ac and tid, and sends the ciphertext to  $\mathcal{A}$ .
- $\mathcal{A}$  interacts with honest parties arbitrarily subject only to the following two conditions that 1)  $\mathcal{A}$  cannot ask the decryption oracle for the ciphertext  $c$  with ac and tid, and 2)  $\mathcal{A}$  cannot corrupt any clients specified by ac.
- Finally,  $\mathcal{A}$  outputs a bit  $b'$ .

Figure 1. We define the advantage of the adversary  $\mathcal{A}$  to be the absolute difference between  $1/2$  and the probability that  $b' = b$ .

Note it is easy to have a unique tid for a client operation (e.g., using a concatenation of the client identity cid and the timestamp of the operation ts). We comment that we do

not need to additionally define decryption consistency (as in threshold encryption), as this is captured by Agreement 2 of the reliability goals (introduced below).

Our definition is easily shown to imply input causality (causal order) [49], which prevents the faulty replicas from creating an operation derived from a correct client’s but that is delivered (and so executed) before the operation from which it is derived. The problem of preserving input causality was introduced in BFT atomic broadcast protocols by Reiter and Birman [49], later refined by Cachin et al. [15], and recently generalized by Duan et al. [22]. Preserving causal order equally makes sense in BFT pub/sub systems.

We do not aim to achieve confidentiality on publication headers or subscription constraints, although they need to be protected for some applications.

*Reliability goals.* We have the following reliability goals:

- **Agreement 1:** If any correct replica delivers a write operation  $m$ , then every correct replica delivers  $m$ .
- **Agreement 2:** If any correct subscriber delivers a publication  $p$  matching its subscription  $T$ , then every correct subscriber who has the same subscription  $T$  and has access to  $p$  delivers  $p$ .
- **Total Order 1:** If a correct replica has delivered write operations  $m_1, m_2, \dots, m_s$  and another correct replica has delivered  $m'_1, m'_2, \dots, m'_{s'}$ , then  $m_i = m'_i$  for  $1 \leq i \leq \min(s, s')$ .
- **Total Order 2 (Publication total order):** If a correct subscriber has delivered  $p_1, p_2, \dots, p_s$  for a subscription  $T$  and another correct subscriber has delivered  $p'_1, p'_2, \dots, p'_{s'}$  for  $T$ , and if the two subscribers have the same access attributes, then  $p_i = p'_i$  for  $1 \leq i \leq \min(s, s')$ .
- **Liveness 1:** If a write operation  $m$  is submitted to  $n - f$  correct replicas, then all correct replicas will eventually deliver  $m$ .
- **Liveness 2 (Publication liveness):** If a publisher is correct and submits  $p$  matching a subscription  $T$ , then all correct subscribers that issued a subscription  $T$  and have access to  $p$  will eventually deliver  $p$ . If a subscriber issues a subscription  $T$ , then it will deliver all authorized publications matching  $T$ .
- **No Creation:** If a subscriber delivers a publication, then the publication was published by some publisher.
- **No Duplication:** A subscriber delivers no publications twice.

Agreement 1, Total Order 1, and Liveness 1 are properties for all write operations. The other properties are ones for pub/sub operations with respect to subscribers. We have considered access control when defining these properties. The properties can be easily simplified to work without considering access control.

Prior formalization on reliable pub/sub systems [58], [20], [48], [39] only consider a much smaller subset of properties we defined here. In particular, a weaker notion of publication total order was considered in several systems [58], [20], [48], where neither subscription restraints nor access control rules are considered, and total order is enforced among all publications

across all subscribers. The weaker notion is immediately implied by the total order property of brokers (Total Order 1), as subscribers can directly deliver publications in the sequence number order determined by brokers. Moreover, in [58], Total Order 1 is not required, because they did not use a state machine replication approach.

No Creation and No Duplication have been previously formalized by Jehl and Meling [36] but with different names (“authentication” and “uniqueness”).

#### IV. THE CHIOS SYSTEM

Chios addresses two important problems in pub/sub systems, achieving decentralized, privacy-preserving pub/sub with fine-grained access control, and ensuring publication total order even with the two-way information control.

We first review threshold encryption. Then, we describe a toy protocol achieving all security goals except publication total order. Finally, we show our core protocol (Chios) achieving publication total order.

##### A. Review of VIL Threshold Encryption

Conventional labeled threshold encryption takes a single string as the label. We extend the primitive to support a vector of strings  $L = (L_1, \dots, L_s) \in \{0, 1\}^{**}$  as labels. By a vector we mean a sequence of zero or more strings, and we let  $\{0, 1\}^{**}$  denote the space of all vectors. Our scheme supports an *arbitrary* number of vectors, each of which can be of *arbitrary* length.

Syntactically, a robust  $(t, n)$  VIL (variable-input-length) threshold encryption consists of the following algorithms. A probabilistic key generation algorithm TGen takes as input a security parameter  $l$ , the number  $n$  of total servers, and threshold parameter  $t$ , and outputs  $(pk, vk, sk)$ , where  $pk$  is the public key,  $vk$  is the verification key, and  $sk = (sk_1, \dots, sk_n)$  is a list of private keys. A probabilistic encryption algorithm TEnc takes as input a public key  $pk$ , a message  $m$ , and a vector label  $L$ , and outputs a ciphertext  $c$ . A probabilistic decryption share generation algorithm ShareDec takes as input a private key  $sk_i$ , a ciphertext  $c$ , and a label  $L$ , and outputs a decryption share  $\tau$ . A deterministic share verification algorithm Vrf takes as input the verification key  $vk$ , a ciphertext  $c$ , a label  $L$ , and a decryption share  $\tau$ , and outputs  $b \in \{0, 1\}$ . A deterministic combining algorithm Comb takes as input the verification key  $vk$ , a ciphertext  $c$ , a label  $L$ , a set of  $t$  decryption shares, and outputs a message  $m$ , or  $\perp$  (a distinguished symbol).

Our VIL threshold encryption scheme, TDH2-VIL, extends the TDH2 threshold encryption by Shoup and Gennaro [50] and is described in Appendix B.

##### B. A Toy Protocol: Chios without Publication Total Order

**System setup.** We assume that the number of brokers is  $n$ , and  $f$  out of  $n$  brokers can fail arbitrarily (Byzantine failures). We set up an  $(f + 1, n)$  VIL threshold encryption (TGen, TEnc, ShareDec, Vrf, Comb) so that a public key  $pk$  and verification keys  $vk$  are associated with the system, while a secret key is

shared among all brokers, with a broker  $i$  having a key  $sk_i$  for  $i \in [1..n]$ .

**Publisher and subscriber registration.** In Chios, communication among publishers and subscribers is decoupled both in time and space. Publishers and subscribers do not need to know or synchronize with one another. A client (a publisher or a subscriber) registers with brokers using their attributes. During the registration, the brokers collectively verify and store client attributes. Chios runs BFT to ensure the registration information is consistent among brokers. More specifically:

- A client sends its attributes and the corresponding proof to brokers as a special registration operation.
- Upon receiving a registration operation, brokers verify the correctness of client attributes. Brokers discard the operation if the verification fails. (Note the verification of the client attributes can be done offline or online, as in PKI registration.) Brokers run the BFT protocol to assign a sequence number to the registration operation and store the operation in sequence number order. Brokers send replies signaling the success of registration.
- Upon receiving  $f + 1$  matching replies, the client completes the registration.

**Advertisements and subscriptions.** During the advertisement process, publishers advertise to the system their publication scopes, and the brokers broadcast the type of events to all potential subscribers (who show an intent to receive subscriptions during the registration process or later via subscriptions). The advertise operation can be viewed as a special pub operation. During the subscription process, subscribers submit their subscriptions which are stored at the brokers. Advertisements and subscriptions are treated as BFT write operations that need to be ordered.

**Publishing (with confidentiality and fine-grained access control).** Let  $ts$ ,  $op$ ,  $o$ ,  $hr = [hr_1..hr_s]$ ,  $ac = [ac_1..ac_t]$ , and  $p$  be the timestamp, the operation type (pub), the executable operation  $o$  (which makes Chios stateful), the header, the access control policies, and the payload of a publication, respectively. The header  $hr$  consists of the topics of a publication and optionally additional associated-data that do not need to be privacy-protected. The approach provides fine-grained (per-publication) and attribute-based access control.

- A publisher  $cid$  takes as input  $ts$ ,  $op$ ,  $o$ ,  $hr$ ,  $p$ , and  $ac$ , and computes a threshold encryption ciphertext as follows. The vector of labels  $L$  for the client is of the form  $(cid, ts, op, hr, ac)$ . The client  $cid$  takes as input the threshold encryption public key  $pk$ ,  $L$ , and  $p$ , and outputs a labeled ciphertext  $(L, c) \xleftarrow{\$} TEnc(pk, p, L)$  using our vector-label-input threshold encryption. It sends brokers  $(L, c)$  as a BFT write operation.
- Upon receiving a client publication, brokers run the BFT protocol to order the publication (by assigning a sequence number to the publication), store the publication, and execute the associated operation  $o$  in sequence number order. The brokers send replies to the write request which may contain the executed result for the publisher.

- Upon receiving  $f + 1$  matching replies, the client completes the publish operation.

**Notify.** During the process, brokers enforce access control and send publications to authorized and interested subscribers. More specifically:

- Brokers decide authorized and interested subscribers for a publication  $(L, c)$  by matching publication topics with existing subscription constraints, checking access control policies associated with the publication, and checking global access control policies already installed in the brokers. For authorized and interested subscribers, each broker  $i \in [1..n]$  sends them its decryption share  $\tau_i \xleftarrow{\$} \text{ShareDec}_{sk_i}(L, c)$  and the sequence number  $sn$  assigned to the labeled ciphertext  $(L, c)$ .
- Upon receiving  $f + 1$  matching publications with *valid* decryption shares from the brokers with the same sequence number  $sn$ , a subscriber runs Comb to obtain the publication in plaintext and delivers it.

**Read.** As in Kafka, Chios can serve as a storage system and an authorized client can read stored data (publications) at brokers via engaging a protocol between the client and brokers.

- A client sends brokers a read request for a particular publication of the form  $(L, c)$ .
- Upon receiving a read request, brokers decide if the client is authorized by checking access control policies associated with the publication. If the client is allowed to have access to the publication, each broker  $i \in [1..n]$  sends the client its decryption share  $\tau_i \xleftarrow{\$} \text{ShareDec}_{sk_i}(L, c)$ .
- Upon receiving  $f + 1$  matching replies with *valid* decryption shares from the brokers with the same sequence number  $sn$ , the client runs Comb to obtain the publication in plaintext and delivers it.

The above system achieves all properties in Sec. III-A except publication total order.

### C. Chios with Publication Total Order

Intuitively, to achieve publication total order, each subscriber needs to maintain a log of valid publications received and deliver them according to the sequence number order assigned by brokers; however, due to access control and subscriber interests, *not all* publications will be sent to all subscribers. Therefore, subscribers do not know if they should wait for or skip publications with certain sequence numbers.

To tackle the issue, we first require servers to additionally maintain topic-based sequence numbers in addition to the global sequence numbers. Doing so, however, does not suffice, as even if two subscribers have the same subscriptions, they may not receive the same publications due to the access control rules. We thus also require that servers send empty messages with sequence numbers to subscribers who are not authorized to receive the corresponding publications. This way, subscribers can safely skip empty publications and go ahead to deliver publications with larger publication sequence numbers.

We now describe in more detail how Chios achieves publication total order. As illustrated in Fig. 2, we maintain two

| $sn$ | $cid$ | $ts$ | $op$  | $p$   | $ac$    | $tp$                            |
|------|-------|------|-------|-------|---------|---------------------------------|
| 0    | 1000  | 4    | pub   | $m_0$ | NULL    | price="105"                     |
| 1    | 1001  | 6    | write | $m_1$ | 101     | NULL                            |
| 2    | 1000  | 10   | pub   | $m_2$ | 100,101 | price="105",<br>county="Orange" |

Figure 2. Data blocks and the publication order indices.

| $tp$            | S-PS    |
|-----------------|---------|
| price="105"     | 0-0,2-1 |
| county="Orange" | 2-0     |

tables: a table for data blocks and a table for publication order indices. The data block table maintains all operations in the system, which are stored in the database. The publication order index table contains metadata of the data blocks and can be derived from the data block table. The index table is stored either in the database or in memory.

For each operation, we store the sequence number ( $sn$ ), the client id ( $cid$ ), the operation type ( $op$ ), the message payload ( $p$ ), timestamp ( $ts$ ), access control rules ( $ac$ ), and the publication topics ( $tp$ ). Certain fields in the data blocks can be NULL.

The publication order index table helps achieve topic-based total order (i.e., total order for the publications according to the topics). Specifically, for each topic, we maintain a simple data structure S-PS, where the S field consists of the sequence numbers of operations ( $sn$ , the same sequence numbers as in the data blocks table), and the PS field consists of the per-topic sequence numbers ( $ps$ ).

The PS field contains incremental sequence numbers for a specific topic, ensuring there is no gap in the sequence numbers for operations with the same topic. For instance, as shown in Fig. 2, in the data block table, operations with sequence number 0 and 2 are publications. There are two topics involved in the data block table: price = "105" and county = "Orange". Correspondingly, there are two topics in the publication index table. As both publications have the topic (price = "105"), the topic in the index table has two S-PS numbers: 0-0 and 2-1. The numbers 0 and 2 in the S field are the sequence numbers in the data block table, while the numbers 0 and 1 in the PS field are per-topic sequence numbers. Specifically, brokers distinguish three cases:

- For authorized and interested subscribers, each broker  $i \in [1..n]$  sends them  $(tp, ps, \tau_i)$ , where  $tp$  is the topic,  $ps$  is the topic sequence number, and  $\tau_i \xleftarrow{s} \text{ShareDec}_{sk_i}(L, c)$  is the decryption share for broker  $i$ .
- For unauthorized and interested subscribers, each broker  $i \in [1..n]$  sends them  $(tp, ps, \perp)$ , where  $\perp$  is a short distinguished symbol denoting an empty message payload (so that subscribers can safely skip the sequence numbers for a particular topic).
- For uninterested subscribers, brokers send nothing.

Each subscriber maintains a log of publications (either empty publications or publications in plaintext) for each topic  $tp$ . It delivers publications according to the  $ps$  order, and example of which is illustrated in Fig. 3. More specifically,

- Upon receiving  $f + 1$  matching publications of the form  $(tp, ps, \tau_i)$  from different brokers, a subscriber runs Comb to obtain a publication in plaintext  $p$  and stores  $p$  in  $\Delta$  in its  $ps$ 's position.

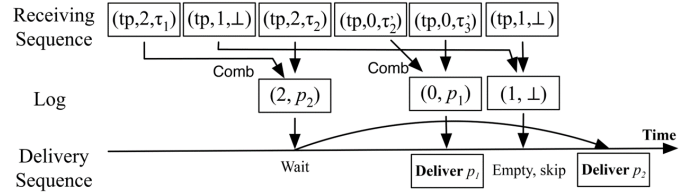


Figure 3. An example of how a subscriber delivers publications assuming  $f = 1$  and  $n = 4$ . The subscriber receives a sequence of messages from BFT brokers and stores them in its buffer. It first receives  $f + 1 = 2$  matching messages with  $ps = 2$ . It then runs Comb to obtain a publication  $p_2$  in plaintext and stores in its log. The subscriber has to wait until publications with smaller sequence numbers (i.e.,  $ps = 0, 1$ ) have been dealt with. After the subscriber receives 2 matching messages with  $ps = 0$ , it runs Comb to obtain  $p_1$  and delivers  $p_1$ . It then waits for messages with  $ps = 1$ . After the subscriber receives two empty messages for  $ps = 1$ , it directly skips the message and delivers message  $p_2$  stored.

- Upon receiving  $f + 1$  matching publications of the form  $(tp, ps, \perp)$  from different brokers, the subscriber directly skips the empty publication in the array  $\Delta$  in its  $ps$ 's position.
- The subscriber delivers a publication  $p \in \Delta$  with a sequence number  $ps$ , if all publications with sequence numbers smaller than  $ps$  are either delivered (for non-empty publications) or skipped (empty publications).

## V. A MODULAR FRAMEWORK

Chios provides a modular framework allowing trade-offs between functionality, security, and efficiency. Chios currently supports four modules, including an encryption-free module (Module 1), the module using threshold encryption (Module 2, the one we described in Sec. IV), a module using hybrid encryption, and one with broadcast encryption. The last two can be found in Appendix C.

## VI. IMPLEMENTATION

Chios consists of a Java library and a Python library with about 30,000 lines of new code. We utilize BFT-SMaRt [11] written in Java as the underlying consensus engine, as BFT-SMaRt is "the most advanced and most widely tested implementation of a BFT consensus protocol" [37]. We use LevelDB [29] as the database. We extend the BFT-SMaRt library and implement a key-value store service. The Java library serves as an ordering service, which assigns a sequence number to a client operation. Then, we wrap the library in Python and develop all the core functionalities.

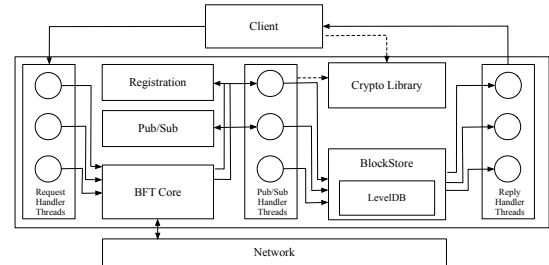


Figure 4. System architecture and message flow.

Fig. 4 illustrates the system architecture and the message flow. The client operations are first handled through a request

handler thread pool and the operations are then relayed to the BFT core. The BFT core batches concurrent client operations and assigns a sequence number to each operation. The ordered client operations are then processed by a pub/sub handler thread pool. Each thread processes a client operation at a time and outputs a reply according to the operation type. Chios uses a *batch-process, block-store* approach, where operations are batched according to a tunable parameter `BlockSize`, ordered, processed, and the results are stored in the database in blocks.

We use ECDSA for authentication and use SHA-256 as our hash function. We implement TDH2-VIL and threshold PRF [16] using the Charm Python library [3]. We use the NIST P-256 curve to provide 128-bit security. We use AES and CBC with ciphertext stealing as our blockcipher and encryption scheme, respectively, to implement the NNL scheme [45].

## VII. EVALUATION

**Settings.** We deployed Chios in multiple platforms (Amazon EC2, IBM Cloud and bare metal servers, and a local Dell/EMC cluster). In this section, we focus on Amazon EC2 where we utilize up to 66 nodes from five different regions in five continents. Each node, by default, is a compute-optimized *c5.2xlarge* type with 8 virtual CPUs (vCPUs) and 16GB memory. We also test the performance using a general-purpose *t2.medium* type with two vCPUs and 4GB memory to evaluate the performance on different hardware. We evaluate our protocols in both LAN and WAN settings, where the LAN nodes are selected from the same EC2 region, and the WAN nodes are uniformly selected from different regions.

We evaluate the protocols under different network sizes (number of replicas) and contention levels (number of concurrent clients). For each experiment, we use  $f$  to represent the network size, where  $3f + 1$  brokers are launched in total. We use P, C, and B to represent the encryption-free module (Module 1), the threshold encryption module (Module 2), and the broadcast encryption module (Module 4), respectively. We did not include the hybrid encryption module, as its performance is almost indistinguishable from the encryption-free module. Let  $\text{Mod} \in \{P, C, B\}$  and let  $\text{op}(\text{Mod})$  represent the operation  $\text{op}$  in the operation using the  $\text{Mod}$  module. For instance,  $\text{pub}(C)$  denotes pub operations for Module 2.

**Overview.** For the minimum one failure setting ( $f = 1$ ), the Chios protocol with all desirable features (pub/sub, decentralized confidentiality, and fine-grained access control), achieves throughput of 45 kops/s for pub operations in LAN. To rigorously demonstrate Chios’s performance, we first compare Chios Module P with five other pub/sub systems. Next, we evaluate the performance for different Chios Modules.

**Comparison with five other pub/sub systems.** We first compare Chios Module P with the following five systems, where Chios-Solo, Kafka, and Fabric-Solo are unreplicated systems, while Kafka-Rep and Fabric-Kafka are crash fault-tolerant systems:

- Chios-Solo. Unreplicated, single-node version of Chios.
- Kafka. As we summarized in Table I in Sec. I, Kafka favors performance over reliability and does not achieve any security

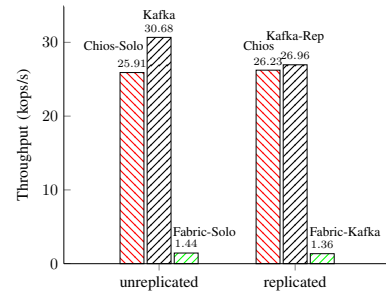


Figure 5. Throughput of Chios, Chios-Solo, Kafka, Kafka-Rep, Fabric-Kafka, and Fabric-Solo.

or reliability goals which we surveyed even in the crash failure model.

- Kafka-Rep. Kafka also supports passive (primary-backup) replication for its brokers with no total order guarantees.
- Fabric-Kafka [6]. Fabric is a popular permissioned blockchain system. Fabric currently does not protect against Byzantine failures. Fabric-Kafka uses the Zookeeper system in Kafka to achieve consensus and is thus only crash fault-tolerant. Hyperpubsub is pub/sub auditing system using Fabric (with Raft) and it is thus slower than Fabric-Kafka.
- Fabric-Solo [6] uses a single node for consensus and is thus not fault-tolerant. One Trinity instance [48] uses Fabric-Solo as its pub/sub auditing system and is slower than Fabric-Solo.

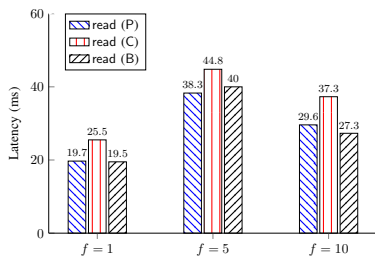
We implement a read/write smart contract for Fabric and use the write operation for the write throughput. Our evaluation for throughput is standard: publishers send brokers operations, and we increase the number of publishers to obtain the peak throughput. We first find out the number of publishers when each system reaches peak throughput. To ensure a fair comparison, we evaluate the systems under the same *total* workload. Namely, the total number of operations sent publishers is the same for all systems. We let the size of all operations be 1KB and we utilize network sizes that tolerate one failure, i.e., four for Chios and three for Fabric and Kafka.

We report the throughput in LAN using 200 clients of the six systems in Fig. 5. Chios Module P is as efficient as Chios-Solo and is only marginally less efficient than Kafka and Kafka-Rep. Chios is significantly more efficient than Fabric-Kafka and Fabric-Solo and thus even much more efficient than Hyperpubsub and Trinity.

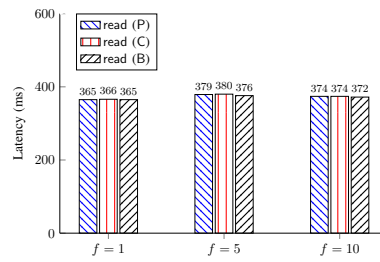
*It is unfair to compare Chios with Kafka with more nodes, as Kafka uses independent server instances for horizontal scalability.*

**Latency of Chios modules.** We assess the latency in both the LAN and WAN settings. We examine the average latency under no contention where only one client sends a single operation to the servers. We let the `BlockSize` be one to understand the latency caused by the protocol itself. In the LAN setting, the network latency is relatively small, so the overhead is more caused by the BFT agreement and execution of operations (e.g., verifying operation types, database interaction). In the WAN setting, the network latency causes more performance degradation than that in the LAN setting. For the threshold encryption and broadcast encryption modules, the latency

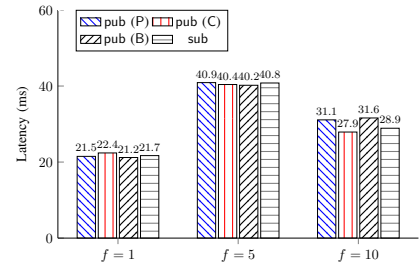




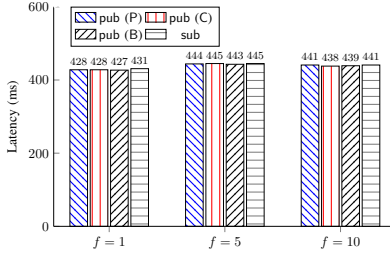
(a) Latency for read operations in the LAN setting with  $f = 1, 5$  and  $10$ .



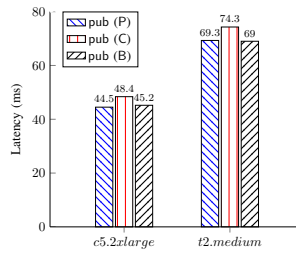
(b) Latency for read operations in the WAN setting with  $f = 1, 5$  and  $10$ .



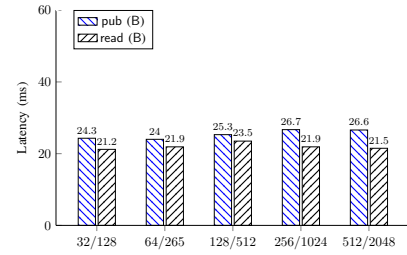
(c) Latency for pub operations in the LAN setting with  $f = 1, 5$  and  $10$ .



(d) Latency for pub operations in the WAN setting with  $f = 1, 5$  and  $10$ .

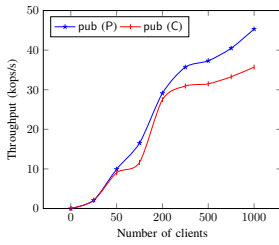


(e) Latency for pub operations in the LAN setting on different hardware with  $f = 1$ .

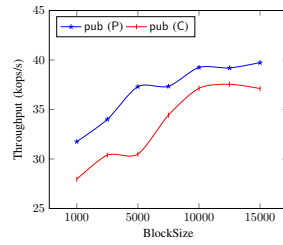


(f) Latency for read(B) and pub(B) operations in LAN with  $f = 1$ , where  $a/b$  represents that  $a$  out of  $b$  subscribers are revoked.

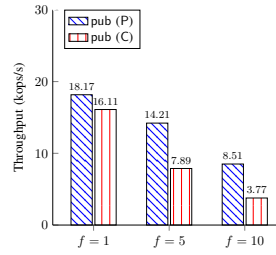
Figure 6. Latency of different Chios modules.



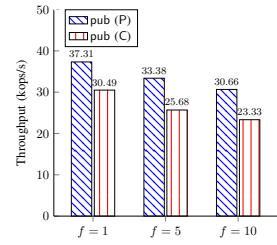
(a) Throughput in the LAN as the number of client threads increases.



(b) Throughput in the LAN as the size of block increases.



(c) Peak throughput in the WAN using 500 clients.



(d) Throughput in the LAN using 500 clients.

Figure 7. Throughput of different Chios modules.

evaluated includes the overhead of client-side encryption.

**For read operations.** We assess the latency for read operations in all three encryption modules as the network size increases. Fig. 6(a) reports the latency for the LAN setting. As read(P) involves no encryption, it has the lowest latency among all three modules. For read(C), replicas verify the ac rules, decrypt the ciphertext, and send decryption shares to the clients. Additional overhead is thus incurred. For read(B), we test the latency for the content distribution phase, as the key distribution phase needs to be done only once. The performance of read(B) is consistently better than that of read(C), as it uses symmetric cryptography only. In the WAN setting, the latency difference among the three modules is smaller, as shown in Fig. 6(b), mainly because network latency dominates the overhead.

**For pub and sub operations.** We report their latency in Fig. 6(c) and Fig. 6(d). We also report the latency of pub using different hardware in Fig. 6(e). We find that the latency for pub operations is higher than that of read operations. We also find that the latency difference for pub operations between the LAN and WAN settings is much higher than that for read operations. The findings are expected, as Chios implements the BFT read

optimization which reduces much communication overhead.

**Other operations.** In Fig. 6(f), we evaluate the performance of read(B) and write(B) operations as the number of subscribers increases. In all these experiments, we randomly revoke 1/4 of the total subscribers. We find for both operations, the latency is steady, regardless of the number of subscribers. The reason is that the broadcast encryption module uses symmetric cryptography only for the content distribution phase.

**Throughput of Chios modules.** We evaluate the throughput of Chios with varying BlockSize in the LAN setting when  $f = 1$ . Fig. 7(a) demonstrates the throughput when the BlockSize is the 5,000 and as the number of concurrent clients increases from 25 to 1,000. The system reaches peak throughput when the number of concurrent clients is larger than 800. The peak throughput that we observe for pub (P) is around 45 kops/s in LAN and 18 kops/s in WAN. We report the throughput when the total number of clients is 375 and as the BlockSize increases in Fig. 7(b). We observe that the throughput becomes larger when the BlockSize increases; however, after BlockSize is larger than 10k, the throughput ceases to increase. In all experiments, the throughput for pub (C) is lower than that of pub (P) due to the cryptographic overhead.

We report the throughput for pub (P) and pub (C) using up to 31 servers and 500 concurrent clients for the LAN setting and the WAN setting, in Fig. 7(d) and Fig. 7(c), respectively. For both the LAN and WAN settings, we find that the throughput for both modules degrade when the number of servers increases (**resembling that of BFT-SMaRt, the consensus engine for Chios**), and the throughput for pub (C) degrades more significantly due to the cryptographic overhead.

### VIII. CONCLUSION

We design and implement Chios, a highly efficient and intrusion-tolerant pub/sub system. Chios addresses two major challenges in pub/sub in terms of confidentiality and reliability: Chios achieves decentralized confidentiality with fine-grained and attribute-based access control and publication total order with two-way information control. Chios provides modular instances designed to meet different goals. Through extensive evaluation, we demonstrate Chios is efficient.

### REFERENCES

- [1] B. Adida, O. de Marneffe, O. Pereira, and J.-J. Quisquater. Electing a University President using open-audit voting: Analysis of real-world use of Helios. *EVT/WOTE 2009*.
- [2] A. Adya et al. FARSITE: Federated available and reliable storage for incompletely trusted environments. *OSDI '02*.
- [3] J. A. Akinyele et al. Charm: a framework for rapidly prototyping cryptosystems. *Journal of Cryptographic Engineering*, 3(2):111–128, 2013.
- [4] Amazon Simple Notification Service (SNS). <https://aws.amazon.com/sns/>
- [5] E. Androulaki, C. Cachin, A. D. Caro, and E. Kokoris-Kogias. Channels: Horizontal scaling and confidentiality on permissioned blockchains. *ESORICS 2018*.
- [6] E. Androulaki et al. Hyperledger Fabric: a distributed operating system for permissioned blockchains. *EuroSys 2018*.
- [7] Apache Kafka. <https://kafka.apache.org/>
- [8] J. Bacon, D. M. Eyers, J. Singh, and P. R. Pietzuch. Access control in publish/subscribe systems. *DEBS 2008*.
- [9] E. Ben-Sasson, A. Chiesa, A. D. Genkin, and E. Tromer, and M. Virza. SNARKs for C: Verifying program executions succinctly and in zero knowledge. *CRYPTO 2013*.
- [10] A. Bessani, E. Alchieri, M. Correia, and J. Fraga. DepSpace: A Byzantine fault-tolerant coordination service. *EuroSys 2008*.
- [11] A. Bessani, J. Sousa, and E. Alchieri. State Machine Replication for the Masses with BFT-SMART. *DSN 2014*.
- [12] J. Bethencourt, A. Sahai, and B. Waters. Ciphertext-policy attribute-based encryption. *IEEE S&P 2007*.
- [13] D. Boneh, C. Gentry, and B. Waters. Collusion resistant broadcast encryption with short ciphertexts and private keys. *CRYPTO 2005*.
- [14] M. Brandenburger, C. Cachin, R. Kapitza, and A. Sorniotti. Blockchain and trusted computing: Problems, pitfalls, and a solution for Hyperledger Fabric. *SRDS 2019*.
- [15] C. Cachin, K. Kursawe, F. Petzold, and V. Shoup. Secure and efficient asynchronous broadcast protocols (extended abstract). *CRYPTO 2001*.
- [16] C. Cachin, K. Kursawe, and V. Shoup. Random oracles in Constantino: Practical asynchronous Byzantine agreement using cryptography. *Journal of Cryptology* 18(3), 219–246.
- [17] C. Cachin and S. Tessaro. Asynchronous verifiable information dispersal. *SRDS 2005*.
- [18] R. Cheng et al. Ekiden: A platform for confidentiality-preserving, trustworthy, and performant smart contract execution. *EuroS&P, 2019*
- [19] CMS advances interoperability & patient access to health data through new proposals. <https://edit.cms.gov/newsroom/fact-sheets/cms-advances-interoperability-patient-access-health-data-through-new-proposals>
- [20] T. Chang, S. Duan, H. Meling, S. Peisert, and H. Zhang. P2S: A fault-tolerant publish/subscribe infrastructure. *DEBS 2014*.
- [21] R. Cramer, I. Damgard, and J. Nielsen. Secure multiparty computation and secret sharing. Cambridge University Press.
- [22] S. Duan, M. K. Reiter, and H. Zhang. Secure causal atomic broadcast, revisited. *DSN 2017*.
- [23] S. Duan and H. Zhang. Practical state machine replication with confidentiality. *SRDS*, 2016.
- [24] Benjamin Eze et al. Policy-based data integration for e-health monitoring processes in a B2B environment: experiences from Canada. *JTAER* 5, 1 (April 2010), 56–70.
- [25] FAYE, Simple pub/sub messaging for the web. <https://faye.jcoglan.com>
- [26] A. Fiat and M. Naor. Broadcast encryption. *CRYPTO 1993*.
- [27] L. Fiege et al. Security aspects in publish/subscribe systems. *DEBS'04*.
- [28] J. Garay, R. Gennaro, C. Jutla, and T. Rabin. Secure distributed storage and retrieval. *TCS*, 243 (1-2): 363–389, 2000.
- [29] S. Ghemawat and J. Dean. LevelDB, A fast and lightweight key/value database library by Google. 2014.
- [30] Google Cloud Pub/Sub. <https://cloud.google.com/pubsub/>
- [31] M. Ion, G. Russello, and B. Crispo. An implementation of event and filter confidentiality in pub/sub systems and its application to e-health. *CCS 2010*.
- [32] M. Ion, G. Russello, and B. Crispo. Supporting publication and subscription confidentiality in pub/sub networks. *SecureComm 2010*.
- [33] M. Ion, G. Russello, and B. Crispo. Design and implementation of a confidentiality and access control solution for publish/subscribe systems. *Computer Networks* 56, 7 (2012), 2014–2037.
- [34] IOTA. <https://www.iota.org/>
- [35] A. Iyengar, R. Cahn, C. Jutla, and J. Garay. Design and implementation of a secure distributed data repository. *IFIP ISC*, 1998.
- [36] L. Jehl and H. Meling. Towards Byzantine fault tolerant publish/subscribe: a state machine approach. *HotDep 2013*.
- [37] A. Kate, Y. Huang, and I. Goldberg. Distributed key generation in the wild. *IACR Cryptology ePrint Archive* 2012: 377 (2012).
- [38] R. Kazemzadeh and H. Jacobsen. Reliable and highly available distributed publish/subscribe service. *SRDS 2009*.
- [39] R. Kazemzadeh and H. Jacobsen. Publiprime: Exploiting overlay neighborhoods to defeat byzantine publish/subscribe brokers.
- [40] H. Khurana. Scalable security and accounting services for content-based publish/subscribe systems. *SAC 2005*.
- [41] A. E. Kosba et al. Hawk: The Blockchain Model of Cryptography and Privacy-Preserving Smart Contracts. *IEEE Security & Privacy* 2016.
- [42] J. Kubiatoiwicz et al. OceanStore: An architecture for global-scale persistent storage. *ASPLOS*, 2000.
- [43] A. B. Lewko and B. Waters. Decentralizing Attribute-Based Encryption. *EUROCRYPT 2011*.
- [44] MQTT. [www.mqtt.org/](http://www.mqtt.org/)
- [45] D. Naor, and M. Naor, and J. B. L. Lotspiech. Revocation and tracing schemes for stateless receivers. *CRYPTO 2001*.
- [46] E. Onica, P. Felber, H. Mercier, and E. Rivière. Confidentiality-preserving publish/subscribe: A survey. *ACM Comput. Surv.*, 2016.
- [47] R. Padilha and F. Pedone. Belisarius: BFT Storage with confidentiality. *NCA 2011*.
- [48] G. S. Ramachandran et al. Trinity: A Byzantine fault-tolerant distributed publish-subscribe system with immutable blockchain-based persistence. *ICBC*.
- [49] M. K. Reiter and K. Birman. How to securely replicate services. *ACM TOPLAS*, vol. 16 issue 3, pp. 986–1009, ACM, 1994.
- [50] V. Shoup and R. Gennaro. Securing threshold cryptosystems against chosen ciphertext attack. *EUROCRYPT '98*.
- [51] M. Srivatsa and L. Liu. Securing publish-subscribe overlay services with EventGuard. *CCS 2005*.
- [52] M. A. Tariq, B. Koldehofe, and K. Rothermel. Securing Broker-Less Publish/Subscribe Systems Using Identity-Based Encryption. *IEEE Trans. on Parallel and Distributed Systems* 25, 2 (2014), 518–528.
- [53] Tendermint core. <https://github.com/tendermint/tendermint>
- [54] S. Vinoski. 2006. Advanced message queuing protocol. *IEEE Internet Computing* 10, 6 (2006).
- [55] G. Wood. Ethereum: A secure decentralised generalised transaction ledger. Byzantium Version.
- [56] A. Wun and H.-A. Jacobsen. A Policy Management Framework for Content-Based Publish/Subscribe Middleware. *Middleware 2007*.
- [57] J. Yin et al. Separating agreement from execution for Byzantine fault tolerant services. *SOSP 2003*.
- [58] K. Zhang, V. Muthusamy, H.-A. Jacobsen. Total order in content-based publish/subscribe systems. *ICDCS 2012*.
- [59] Y. Zhao and D.C. Sturman. Dynamic Access Control in a Content-based Publish/Subscribe System with Delivery Guarantees. *ICDCS 2006*.
- [60] N. Zupan, K. Zhang, and H.-A. Jacobsen. Hyperpubsub: a decentralized, permissioned, publish/subscribe service using blockchains: demo. *Middleware Posters and Demos 2017*.

APPENDIX A  
NOTATION

To facilitate reference to the notation used in the paper, we present the following summary. (The table does not include the notation for the cryptographic schemes which are self-contained.)

|               |   |
|---------------|---|
| $m$           | operation, transaction, or message  |
| $n$           | the number of brokers (servers, replicas)   |
| $f$           | the upper bound of the number of faulty brokers   |
| $t$           | the threshold for threshold encryption and PRF; set as $f + 1$ in this paper              |
| cid           | client identity   |
| ts            | timestamp (an increasing sequence number in this paper)                                   |
| tid           | unique tag (unique session ID); set as cid  ts  |
| $\mathcal{A}$ | adversary   |
| $N$           | the number of potential subscribers; the number of recipients in the broadcast encryption |
| $l$           | security parameter  |
| $L$           | label for threshold encryption; a vector in $\{0, 1\}^{**}$ (the space of all vectors)    |
| op            | operation type (e.g., pub, sub)   |
| hr            | header for a publication  |
| $p$           | publication   |
| ac            | access control rules  |
| tp            | topic   |
| sn            | sequence number assigned to each operation by BFT servers                                 |

Table II  
NOTATION

APPENDIX B

OUR THRESHOLD ENCRYPTION SCHEME

Our VIL threshold encryption scheme, TDH2-VIL, extends the TDH2 threshold encryption by Shoup and Gennaro [50].

TDH2-VIL is the same as the TDH2 scheme, except in handling labels. TDH2-VIL supports an *arbitrary* number of vectors, each of which can be of *arbitrary* length.

We describe the new VIL threshold encryption scheme, TDH2-VIL in Fig. 8. TDH2 is secure against CCA attacks, under the Decisional Diffie-Hellman (DDH) assumption in the random oracle model (ROM). The proof of TDH2 can be extended to show that TDH2-VIL is also secure against CCA attacks, under the DDH assumption in the ROM.

The threshold encryption is chosen ciphertext attack (CCA) secure against an adversary that controls up to  $t-1$  servers. We also require consistency of decryptions, i.e., no adversary that controls up to  $t-1$  servers can produce a ciphertext and two  $t$ -size sets of valid decryption shares (i.e., where Vrf returns  $b = 1$  for each share) such that they yield different plaintexts.

APPENDIX C  
MODULES IN CHIOS

Using VIL threshold encryption accomplishes general access control. Moreover, even when subscriptions and access control rules change frequently, publication matching and access control can be fully executed at the brokers. The drawback is that threshold encryption is relatively expensive for cases where each replica has limited computational power.

Let  $G$  be a group of prime order  $q$  with generator  $g$ . Let  $\mathbb{Z}_q$  be the additive group, integers modulo a prime  $q$ . Let  $m \in \{0, 1\}^l$ , where  $l$  is a security parameter. Choose the following four hash functions:  $H_1: G \rightarrow \{0, 1\}^l$ ,  $H_2: \{0, 1\}^l \times \{0, 1\}^l \times G^4 \rightarrow \mathbb{Z}_q$ ,  $H_3: G^3 \rightarrow \mathbb{Z}_q$ ,  $H_4: \{0, 1\}^* \rightarrow \{0, 1\}^l$ . Define a new hash function  $H_5: \{0, 1\}^{**} \rightarrow \{0, 1\}^l$  so that  $H_5(L_1, \dots, L_z) = H_4(H_4(1, |L_1|, L_1), \dots, H_4(z, |L_z|, L_z), z)$ , where  $z$  is a positive integer.

• **TGen**( $1^{ql}$ ): Choose random points  $a_0, \dots, a_{t-1} \xleftarrow{\$} \mathbb{Z}_q$  and define a polynomial  $F(X) = \sum_{j=0}^{t-1} a_j X^j \in \mathbb{Z}_q[X]$ . For  $i \in [1..n]$ , set  $x_i = F(i) \in \mathbb{Z}_q$  and  $h_i = g^{x_i}$ . Set  $x = F(0)$  and  $h = h_0 = g^x$ . Choose  $\bar{g} \xleftarrow{\$} G$ . Set  $pk = (G, g, h, \bar{g})$ ,  $vk = (pk, h_1, \dots, h_n)$ , and  $sk_i = (pk, x_i)$  for  $i \in [1..n]$ . Return  $(pk, vk, sk)$ .

• **TEnc**( $pk, m, L$ ): Choose at random  $r, s \xleftarrow{\$} \mathbb{Z}_q$  and compute  $c_1 = H_1(h^r) \oplus m$ ,  $u = g^r$ ,  $w = g^s$ ,  $\bar{u} = \bar{g}^r$ ,  $\bar{w} = \bar{g}^s$ ,  $L' = H_5(L)$ ,  $e = H_2(c, L', u, w, \bar{u}, \bar{w})$ ,  $f = s + re$ . Return  $c = (c_1, u, \bar{u}, e, f)$ .

• **ShareDec**( $sk_i, c, L$ ): Given a label  $L$  and a ciphertext  $c = (c_1, u, \bar{u}, e, f)$ , server  $i \in [1..n]$  verifies whether  $c$  is well-formed, i.e.,  $e = H_2(c_1, L', u, w, \bar{u}, \bar{w})$ , where  $L' = H_5(L)$ ,  $w = g^f / u^e$ ,  $\bar{w} = \bar{g}^f / \bar{u}^e$ . If  $c$  is not well-formed, the server returns  $(i, \perp)$ . If  $c$  is well-formed, it chooses at random  $s_i \xleftarrow{\$} \mathbb{Z}_q$ , computes  $u_i = u^{x_i}$ ,  $\hat{u}_i = u^{s_i}$ ,  $\hat{h}_i = g^{s_i}$ ,  $e_i = H_3(u_i, \hat{u}_i, \hat{h}_i)$ ,  $f_i = s_i + x_i e_i$ , and returns  $\tau = (i, u_i, e_i, f_i)$ .

• **Vrf**( $vk, c, L, \tau$ ): Check whether  $c$  is well-formed, and  $\tau = (i, u_i, e_i, f_i)$  is well-formed, i.e.,  $e_i = H_3(u_i, \hat{u}_i, \hat{h}_i)$ , where  $\hat{u}_i = u^{f_i} / u_i^{e_i}$ ,  $\hat{h}_i = g^{f_i} / h_i^{e_i}$ . If both  $c$  and  $\tau$  are well-formed, return  $b = 1$ . If either is not well-formed, return  $b = 0$ .

• **Comb**( $vk, c, L, \{\tau_j\}$ ): Given  $c$  and  $t$  valid decryption shares, if  $c$  is well-formed, run the Lagrange Interpolation in the Exponent on  $\{\tau_j\}$  to get  $c' = H_1(h^r)$  for the  $r$  used to compute  $c$ . Return  $m = c' \oplus c_1$ .

Figure 8. **TDH2-VIL** = (TGen, TEnc, ShareDec, Vrf, Comb) is a  $(t, n)$  labeled threshold encryption scheme for a label vector  $L \in \{0, 1\}^{**}$ . We include all the algorithms for completeness.

**Module 3: “Hybrid” Encryption** When subscriptions and access control rules do not change frequently, or simply are static, it is beneficial to use “hybrid encryption.” Clients can use labeled threshold encryption to encrypt a random session key and store the ciphertext in the brokers. According to subscriptions and access control rules, the decryption shares of the ciphertext of the random session key will be sent to authorized subscribers. Later, the publications will simply be encrypted using the session key. However, a fresh session key needs to be established when subscriptions, access control rules, or even client attributes change. If any of the above situations occur frequently, Module 3 would not work well.

**Module 4: Combining Threshold Encryption with Broadcast Encryption** Module 4 combines threshold encryption with broadcast encryption [26], [45], [13] to overcome the drawback of Module 3. In Module 4, one just needs to set

up broadcast encryption keys *once* (just as session keys) and uses the keys to transmit multiple messages with different access control policies. In our full paper, we provide a generic construction and an efficient instantiation, both of which are technically novel and of independent interests.

#### APPENDIX D CORRECTNESS PROOF

*Theorem 1:* If  $f \leq \lfloor (n-1)/3 \rfloor$ , Chios satisfies all reliability and security goals in Sec. III-A.

**Proof:** Agreement 1, Total Order 1, and Liveness 1 follow from Agreement, Total Order, and Liveness of the underlying BFT protocol, respectively, as operations related to these properties are conventional BFT (write) operations. We are left to show the correctness of the rest of the properties.

We begin with some lemmas. Due to the decryption consistency property of the underlying threshold encryption scheme, there is only one valid plaintext  $p$  for any threshold encryption ciphertext  $c$ . More precisely, we have the following lemma:

*Lemma 2:* For a publication in ciphertext  $c$ , given any  $f+1$  valid decryption shares of  $c$  from any set of  $f+1$  servers, there exists a unique publication  $p$  in plaintext (which may be a distinguished symbol denoting decryption failure).

For simplicity, we let  $c$  denote a publication in threshold encryption ciphertext and  $p$  denote the corresponding publication. We also need the following lemma:

*Lemma 3:* If a subscriber receives  $f+1$  matching and valid decryption shares of the form  $(\tau_i, ps, \tau_i)$ , the subscriber will eventually deliver the underlying publication  $p$ .

*Proof of Lemma 3:* If a subscriber receives  $f+1$  matching and valid decryption shares of the form  $(\tau_i, ps, \tau_i)$ , it will first run Comb to obtain the underlying publication  $p$ . According to our protocol, we just need to show *all* publications with smaller sequence numbers than  $ps$  are either delivered or skipped. To prove this, we first observe that due to the Total Order property of the underlying BFT protocol, all replicas have delivered publications with smaller sequences numbers than  $ps$ . For interested and authorized subscribers, replicas will send their decryption shares (whose sequence numbers are smaller than  $ps$ ), and these publications will deliver before  $ps$  deliver. For interested but unauthorized subscribers, replicas will send empty messages, and these publications are skipped. The lemma thus follows naturally.

**Agreement 2.** If a correct subscriber delivers a publication  $p$  matching a subscription  $T$ , the corresponding publication in ciphertext  $c$  was delivered by at least  $f+1$  replicas. Among the  $f+1$  replicas, *at least* one of them is correct and delivers  $c$ , as there are at most  $f$  faulty replicas. By the Agreement property of the BFT protocol, every correct replica delivers  $c$ . Therefore, every correct replica will broadcast their decryption shares for the publication to each subscriber who has the same subscription  $T$  and has access to  $p$ . Accordingly, interested and authorized subscribers will receive  $f+1$  valid decryption shares from possibly different sets of  $f+1$  replicas. We can now directly apply Lemma 2 and Lemma 3 to conclude that all interested and authorized subscribers deliver the publication  $p$ .

**Total Order 2 (Publication total order).** If a correct subscriber has delivered  $p_1, p_2, \dots, p_s$  for a subscription  $T$ , according to Agreement 1, all correct subscribers will deliver these publications. We are left to show subscribers will deliver these publications in the same order. In Chios, subscribers deliver authorized publications according to the per-topic total order assigned by servers which are totally ordered and unauthorized publications are skipped by all correct subscribers. Total Order 2 follows easily.

**Liveness 2 (Publication liveness).** We now prove Liveness 2. We begin with the first part of Liveness 2. If a publisher is correct and submits a publication  $p$  matching a subscription  $T$ , then servers will check the access control rules and send decryption shares to authorized subscribers with the subscription  $T$  and will send unauthorized but interested subscribers an empty message (using a distinguished symbol  $\perp$ ). Therefore, authorized and interested subscribers will receive at least  $f+1$  matching and valid decryption shares and according to Lemma 3, they will deliver the publication  $p$ .

Analogously, we can prove the second part of Liveness 2. If a subscriber issues a subscription  $T$ , servers will send all matching publications and all authorized and interested subscribers will receive all publications matching  $T$ .

**No Creation, No Duplication, and Confidentiality.** No Creation follows from channel authentication between publishers and servers and between servers and subscribers. More specifically, if a subscriber delivers a publication, then the corresponding publication in ciphertext  $c$  was delivered by at least  $f+1$  replicas. Among the  $f+1$  replicas, *at least* one of them is correct and delivers  $c$ , as there are at most  $f$  faulty replicas. This implies some publisher has published the publication.

No Duplication easily follows from the fact that each subscriber maintains a log of received publications and deliver each of them in sequence number order and by once.

Confidentiality and access control follow from the chosen-ciphertext security of our vector-label-input threshold encryption.

This completes the proof of the theorem. ■

#### APPENDIX E MODULE 4: A GENERAL CONSTRUCTION AND AN EFFICIENT INSTANTIATION

##### A. Building Blocks

**Threshold PRF.** We review threshold PRF (e.g., [16]), where a public key is associated with the system and a PRF key is shared among all the servers. A  $(t, n)$  threshold PRF scheme for a function  $F$  consists of the following algorithms. A probabilistic key generation algorithm  $F_{\text{Gen}}$  takes as input a security parameter  $l$ , the number  $n$  of total servers, and threshold parameter  $t$ , and outputs  $(pk, vk, sk)$ , where  $pk$  is the public key,  $vk$  is the verification key, and  $sk = (sk_1, \dots, sk_n)$  is a list of private keys. A PRF share evaluation algorithm  $F_{\text{Eva}}$  takes a public key  $pk$ , a PRF input  $m$ , and a private key  $sk_i$ , and outputs a PRF share  $y_i$ . A deterministic share verification

algorithm  $Vrf$  takes as input the verification key  $vk$ , a PRF input  $m$ , and a PRF share  $y_i$ , and outputs  $b \in \{0, 1\}$ . A deterministic combining algorithm  $FCom$  takes as input the verification key  $vk$ , a PRF input  $m$ , and a set of  $t$  valid PRF shares, and outputs a PRF value  $y$ . We require the threshold PRF value to be unpredictable against an adversary that controls up to  $t - 1$  servers. We also require the threshold PRF to be robust in the sense the combined PRF value for  $m$  is equal to  $F(m)$ .

**Broadcast encryption.** Broadcast encryption [26] enables a broadcaster to transmit encrypted content over a broadcast channel so that only users qualified can decrypt the content. Broadcast encryption can flexibly deal with a dynamically changing group of qualified users. Broadcast encryption requires a trusted party (may it be the broadcaster) to set up the system keys and user keys.

A broadcast encryption system consists of the following three algorithms. A probabilistic key generation algorithm  $BGen$  takes as input the total number of potential receivers  $N$ , and outputs  $N$  private keys  $\{d_i\}_{i=1}^N$ . A probabilistic encryption algorithm  $BEnc$  takes as input a subset  $S \subseteq [1..N]$ , private keys  $\{d_i\}_{i=1}^N$ , and a message  $m$ , and outputs  $(S, hdr, k, C_k(m))$ , where  $hdr$  is a header and  $C$  is a symmetric encryption algorithm. The broadcast encryption ciphertext consists of  $(S, hdr, C_k(m))$ . A deterministic decryption algorithm takes as input a subset  $S \subseteq [1..N]$ , a user id  $i \in [1..N]$ , the private key  $d_i$  for user  $i$ , and a header  $hdr$ , and outputs  $k$  used to decrypt the symmetric ciphertext.

A broadcast encryption is said to be  $(t, n)$  collusion-resistant if for any adversary who learns the secret keys of at most  $t - 1$  revoked users, the broadcasts do not leak any information to the adversary. In broadcast encryption, if a broadcaster is faulty, different recipients may obtain different plaintexts.

### B. A General Construction

The Module 4 process can be divided into a *key distribution phase* and a *content distribution phase*. The key distribution process needs only to be done once, which contrasts with Module 3. We first describe Module 4 for generic broadcast encryption and then our tailored and “optimized” instantiation using the NNL broadcast encryption that is more efficient than the trivial instantiation.

The key distribution phase is depicted in Fig. 9. A publisher generates broadcast encryption keys for  $N$  subscribers. Each publisher uses vector-label-input threshold encryption to encrypt all keys for  $N$  subscribers, with the label  $L$  specifying access rules. The ciphertexts are stored at the brokers which will, *at some point*, send decryption shares to subscribers so that each subscriber can combine a broadcast encryption key. The brokers fully control when broadcast encryption keys are distributed: either when the subscribers register themselves, express their interests, pay their fees, or when they meet the access control rules.

Brokers need to update publishers and let them know the indices of keys to be revoked, and publishers broadcast new ciphertexts accordingly. Brokers should distribute the keys in a

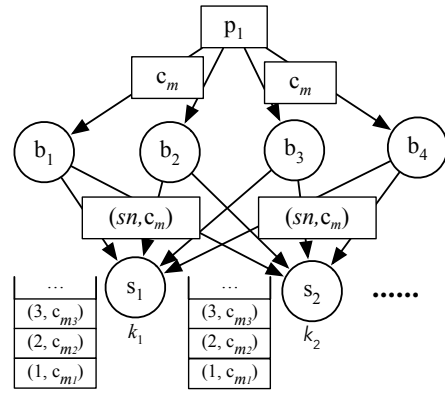


Figure 9. **Key distribution.** In the key distribution process, a publisher sends brokers all the  $N$  broadcast encryption keys encrypted using threshold encryption.

random order (using threshold PRF) to subscribers to maintain the decoupling property. To see why a random order is needed, let us consider the situation when subscribers unsubscribe their interests. In this situation, if the keys are distributed in some known order, their keys need to be known by the broadcaster (i.e., the publisher) when subscribers unsubscribe their interests. This allows the publisher to know the identities of the subscribers and what publications they have received. But in Chios, we require brokers to distribute the keys in a random order, indices do not correspond to real identities of subscribers and the publishers and subscribers remain decoupled and anonymous.

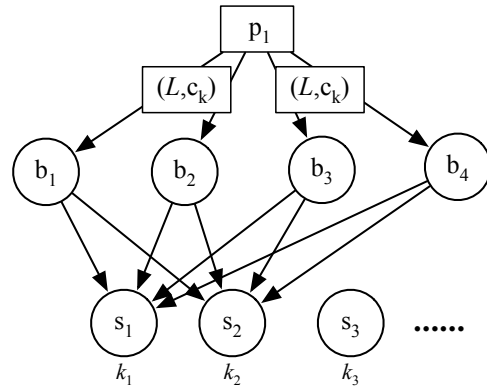


Figure 10. **Content distribution.** In the content distribution process, the publications are encrypted using broadcast encryption.

As shown in Fig. 10, each publisher then uses broadcast encryption to encrypt the publications. The ciphertexts will be sent to brokers who forward them to authorized subscribers.

Module 4 leads to several major benefits: being computationally more efficient and supporting dynamic membership without interaction. There are two drawbacks for the module. First, one would need to use different broadcast encryption for different settings, as there is no one-size-fits-all broadcast encryption. Second, most broadcast encryption cannot handle malicious publishers, because broadcast encryption does not have decryption consistency. If a broadcaster is faulty, different recipients may obtain different results.

### C. Optimized Module 4 Using NNL

In this subsection, we begin by describing the NNL broadcast encryption [45] and then introduce our optimized Module 4 protocol using NNL.

**NNL broadcast encryption.** The complete subtree broadcast encryption by Naor, Naor, and Lotspiech [45] (NNL) that Chios implements is efficient when the number of revoked recipients is small. Let the number of all potential recipients be  $N$ . Imagine a full-binary tree with  $N$  leaves corresponding to the  $N$  recipients. Let  $\mathcal{N}$  denote all these leaves. Let  $v_i$  be a node in the tree. Let  $S_i$  be the *subtree* set of all leaves in the subtree of  $v_i$ . As a result, there are  $2N - 1$  nodes and  $2N - 1$  complete subtrees. In NNL, the broadcaster assigns a random key  $k_i$  to every node  $v_i$  in the tree. Each recipient (leaf) stores all  $\log N + 1$  keys along the path to the root. For instance, if  $N = 2^{32}$ , each recipient stores 33 keys.

In NNL, for a given set  $\mathcal{R}$  of revoked receivers, let  $u_1, \dots, u_r$  be the leaves corresponding to the elements of  $\mathcal{R}$ . Let  $ST(\mathcal{R})$  be the directed Steiner Tree induced by the set  $\mathcal{R}$  of vertices and the root, i.e., the minimal subtree of the full binary tree connecting all leaves in  $\mathcal{R}$ . Let  $S_{i_1}, \dots, S_{i_m}$  be all the subtrees whose roots  $v_1, \dots, v_m$  are adjacent to nodes of outdegree 1 in  $ST(\mathcal{R})$ , but they are not in  $ST(\mathcal{R})$ .

Let  $E$  be a blockcipher and  $F$  be a symmetric encryption secure against IND-CPA attacks. To broadcast  $M$  to  $\mathcal{N} \setminus \mathcal{R}$ , the broadcaster sends  $(i_1, \dots, i_m, E_{k_{i_1}}(K), \dots, E_{k_{i_m}}(K), F_K(M))$ . Each non-revoked user can first decrypt the ciphertext corresponding to its subtree to obtain  $K$  and then  $M$ .

In NNL, the length of the ciphertext is  $r \log N / r$ , the number of keys stored at a recipient is  $\log N$ , and the number of decryptions for a receiver is  $\mathcal{O}(1)$ .

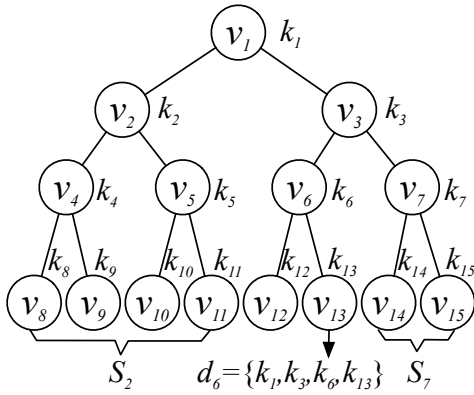


Figure 11. NNL broadcast encryption key setup.

We describe a pictorial example in Fig. 11. Suppose there are  $N = 8$  potential recipients. We generate a full-binary tree with 8 leaves corresponding to the recipients. Let  $v_i$  ( $i \in [1..15]$ ) be a node in the tree. The set of leaves  $\mathcal{N} = [v_8..v_{15}]$  corresponds to the set of recipients. For example, we have  $S_2 = [v_8..v_{11}]$  and  $S_7 = \{v_{14}, v_{15}\}$ . For  $i \in [1..15]$ , the broadcaster assigns a random key  $k_i$  to every node  $v_i$  in the tree. The 1-st recipient  $v_8$  has the key

Let  $n$  and  $N$  be the number of servers and the upper bound of potential subscribers, respectively. We set up a  $(f + 1, n)$  vector-label-input threshold encryption (TGen, TEnc, ShareDec, Vrf, Comb) so that a public key  $pk$  and verification keys  $vk$  are associated with the system, while a secret key is shared among all servers, with server  $i$  having a key  $sk_i$  for  $i \in [1..n]$ . Analogously, we also set up a threshold PRF for  $F: [1..N] \rightarrow [1..N]$ . Initially, we set  $\Omega = [1..N]$ .

- A client runs the NNL key generation algorithm to generate  $2N - 1$  keys  $[k_1..k_{2N-1}]$  corresponding to each node in the NNL tree. The client specifies an ac policy in  $L$ , computes  $[c_1..c_{2N-1}] = [\text{TEnc}(pk, k_1, L).. \text{TEnc}(pk, k_{2N-1}, L)]$ , and sends to the servers  $(L, [c_1..c_{2N-1}])$ .
- Each time a subscriber  $j$  needs to obtain an NNL broadcast encryption key, the servers
  - run the threshold PRF to select a random index  $\lambda \xleftarrow{\$} \Omega$ . The value  $\lambda$  corresponds to the leaf of index  $v_{(N-1)+\lambda}$  in the NNL tree.  $\Omega = \Omega / \{\lambda\}$ .
  - find the all ciphertexts corresponding to the path from root  $v_1$  to  $v_{(N-1)+\lambda}$ , i.e.,  $[c_{j_1}..c_{j_{\log N+1}}]$ .
  - send  $j$  the decryption shares for  $[c_{j_1}..c_{j_{\log N+1}}]$ .
- The subscriber  $j$  combines the decryption shares to recover the keys from root  $v_1$  to  $v_{(N-1)+\lambda}$  in the NNL tree.

Figure 12. Key distribution phase for the confidentiality module using NNL.

$d_1 = \{k_1, k_2, k_4, k_8\}$ , and the 6-th recipient  $v_{13}$  has the key  $d_6 = \{k_1, k_3, k_6, k_{13}\}$ . Suppose  $\mathcal{R} = \{v_{12}, v_{13}\}$ . To encrypt a message  $m$  to  $\mathcal{N} \setminus \mathcal{R}$ , the broadcaster will only need to send the encrypted message to the set  $S_2$  and  $S_7$ . More specifically, the broadcaster randomly selects a key  $K$  and broadcasts a ciphertext  $(2, 7, E_{k_2}(K), E_{k_7}(K), F_K(m))$ . Recipients in the set  $S_2$  and the set  $S_7$  can use their respective keys to first obtain  $K$  and then  $m$ . For instance,  $v_8, v_9, v_{10}, v_{11}$  can use  $k_2$  to get  $K$ , while  $v_{14}$  and  $v_{15}$  can use  $k_7$  to get  $K$ .

In Chios, we extend the NNL key distribution performed by a trusted dealer to achieve *distributed* key distribution. We also extend NNL to support an unbounded number of recipients.

**Optimized confidentiality module 4 using NNL.** If we trivially follow the general approach to instantiate the NNL broadcast encryption, the number of keys stored at a recipient is  $\log N$  and the total keys stored at the brokers would be  $N \log N$ . Chios supports an optimized NNL-based Module 4 reducing the number of keys stored from  $N \log N$  to  $2N - 1$ .

We describe the algorithm in Fig. 12. In the approach, instead of sending all NNL decryption keys (encrypted using threshold encryption) for each potential subscribers, a publisher only needs to send all encrypted keys associated to nodes in the NNL tree. Meanwhile, the brokers fully control how the keys are distributed. The brokers can distribute the keys according to the tree structure agreed beforehand, and each subscriber can combine the individual keys from the encrypted

keys. Doing so helps reduce the number of keys transmitted from  $N \log N$  to  $2N - 1$ . This, however, breaks the decoupling feature of pub/sub systems, as keys are distributed to subscribers according to the NNL tree. To overcome this issue, we require brokers to distribute the keys randomly to subscribers. Our approach can easily allow subscribers to unsubscribe their interests without compromising the anonymity of them. Note that in this case, the publisher needs to know the identity of subscriber. To solve the problem, we require that the indices do not correspond to the user identities and that when updating subscriptions, brokers just need to update publishers the indices of keys to be revoked.

For NNL, the number of recipients need not be bounded. When the initial tree does not suffice, publishers will generate a larger one doubling the original one so that the original one becomes its subtree.

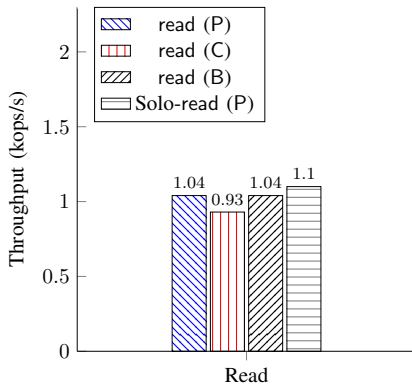


Figure 13. Throughput for read operations in the LAN using 600 clients.

## APPENDIX F SECURELY GENERATING KEYS FOR THRESHOLD CRYPTOGRAPHY

Chios uses threshold cryptography extensively, including (vector-input-length) threshold encryption and threshold PRF. We describe the following three options for generating keys for these schemes.

- The first option is to use a trusted dealer to set up the keys for these schemes. The trusted dealer is only needed in a setup phase, after which the system has no single point of failure. This is the usual assumption that is used in almost all secure, BFT-based distributed systems using threshold cryptography (e.g., [57], [10]).
- The second approach is to use distributed key generation protocols [37] to generate the keys in a distributed manner.
- Adida et al. [1] described an approach that requires trustees (shareholders) to meet in person to securely generate keys for regular public-key encryption schemes. We extend the idea to handle key generation for threshold cryptography. Specifically, a number of trustees are selected to meet physically. For the meeting, a brand-new laptop and a number of brand-new USB sticks will be brought. The hard disk drives of the laptop need to be removed and the wireless network card

needs to be disabled. Then, the laptop will boot up using a standard Linux live-CD, and the regular key generation code as in the first option (inspected by the external experts) will be loaded on the machines through a USB stick. Keys will be then generated in the laptop and secret keys will be placed on individual folders. During the meeting, each shareholder can use a USB stick to get his or her key only. All shareholders should monitor the process. The laptop and the Linux live-CD will be then destroyed after all shareholders obtain their secret keys.

## APPENDIX G ADDITIONAL IMPLEMENTATION DETAILS

We use JPyPe from the client-side to directly invoke the Java library and send operations to replicas. We disable the message handling and reply functionalities in BFT-SMaRt and use Python to handle the replies. Specifically, after the replicas reach a consensus on an operation, the operation and the sequence number are relayed to the Python library. We implement a shim, using Java and Python, as illustrated in Fig. 16. At the Java side, the shim includes a sender thread collecting totally ordered (committed) batches of client operations. The thread simply delivers the ordered batches from the Java library to the Python library at the same node using a socket. At the Python side, the shim includes a receiver thread collecting batches of ordered operations from the socket. The client operations from the batches are then relayed to the Chios modules. The ordered client operations are then processed according to their orders and the operation types.

Chios is flexible and modular, currently supporting four modules designed to meet different goals. For each module, Chios supports three tunable parameters: BatchSize, Block-Size, and the number of threads used at various stages. First, Chios uses a *batch-process, block-store* approach, where operations are batched, ordered, processed, and the results are stored in the database in blocks. BatchSize and BlockSize can be tuned according to the underlying hardware, workloads, and applications. Second, Chios uses multithreading extensively at various stages and allows adjusting the number of threads used for each stage *independently* according to the underlying hardware and workload. Besides, for the threshold encryption module, one can choose to decrypt the threshold encryption ciphertexts when they are written to the system or when they are read. This flexibility allows Chios to provide further trade-offs between read and write performance.

We divide operations into operations changing the system state (i.e., reg, advertise, pub, sub) and operations that do not change system state (i.e., read, notify). For the first case, the replies, client IDs, client operations, and the sequence numbers are delivered to the BlockStore module. When the number of outputs reaches the BlockSize, the operations are stored in the database as a block, the replies are sent to the reply handler thread pool. The replies are then sent to the clients in parallel. For the second case, the pub/sub handler threads obtain the

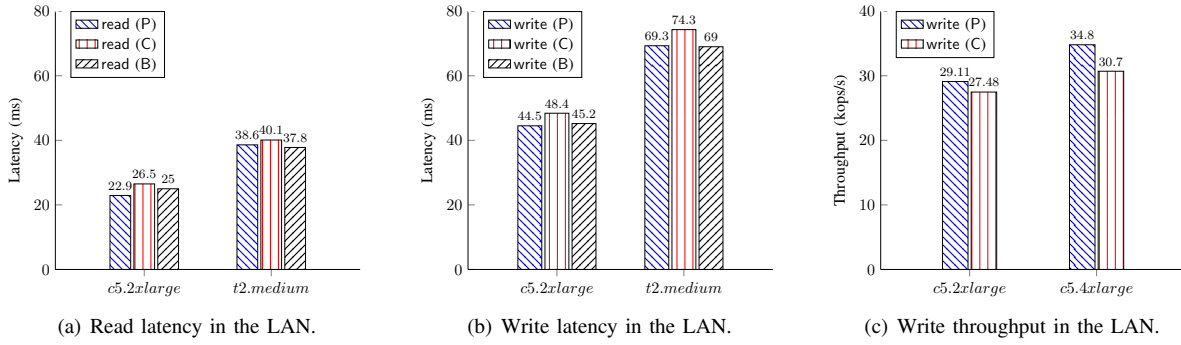
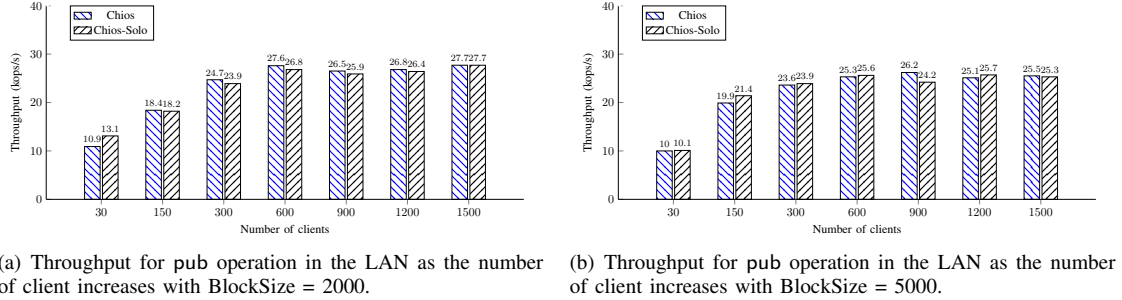


Figure 14. Performance of Chios using different hardware.



(a) Throughput for pub operation in the LAN as the number of client increases with BlockSize = 2000. (b) Throughput for pub operation in the LAN as the number of client increases with BlockSize = 5000.

Figure 15. Throughput of Chios and Chios-Solo for pub operations.

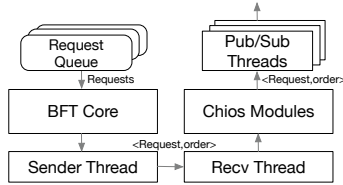


Figure 16. Chios library shim.

data from the database and generate replies and replies are sent to the clients via the reply handler threads.

## APPENDIX H ADDITIONAL EVALUATION

**Read throughput.** We show the peak read throughput of different Chios modules and Chios-Solo for the LAN setting in Fig. 13. We show Module P and Module B are only 5.5% slower than Chios-Solo. Module C is 15.5% slower than Chios-Solo, as Module C requires running decryption operations for threshold encryption. As Chios can be configured to pre-compute decryption shares (Sec. VI), Module C with pre-computation (not shown) is only 6.3% slower than Chios-Solo.

**Performance using different hardware.** As illustrated in Fig. 14, we assess the performance using two different EC2 VM types (with different CPUs and number of threads supported). For all these experiments, the clients use the *c5.2xlarge* type to provide consistent client workloads. When evaluating throughput, we use 200 total number of clients with 5000 BlockSize. As observed from the figures, the latency for the *c5.2xlarge* type is consistently lower than that of the *t2.medium*. For throughput, we compared the

*c5.2xlarge* type with *c5.4xlarge* type. The throughput for the *c5.2xlarge* type is also consistently lower than that of the *c5.4xlarge*. This is, however, expected, as Chios uses multithreading extensively in various stages of the protocol and multithreading offers no help to reduce latency.

**Comparison with unrepliated version.** We compare the throughput of Chios-Solo with Chios Module P in the LAN setting. We vary the number of clients and assess the throughput for pub (P) using BlockSize 2000 (Fig. 15(a)) and 5000 (Fig. 15(b)). We find the throughput for both Chios-Solo and Chios with both sizes is similar.