

HOLMES: Efficient Distribution Testing for Secure Collaborative Learning

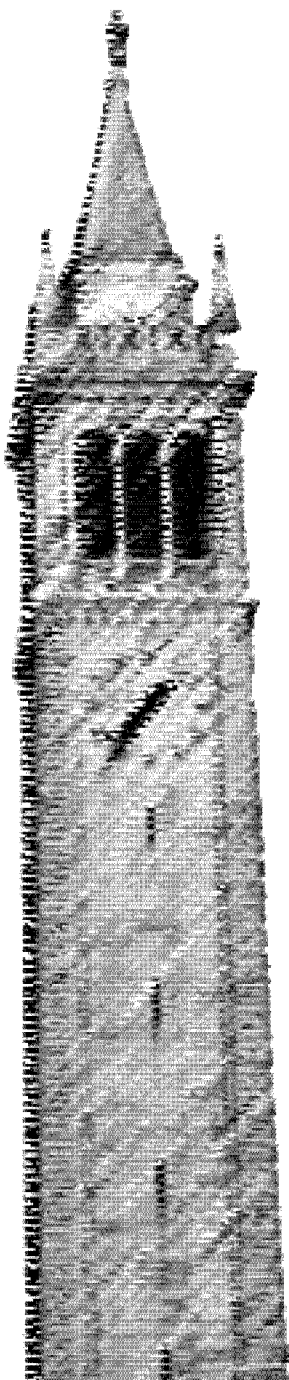
*Ian Chang
Katerina Sotiraki
Weikeng Chen
Murat Kantarcioglu
Raluca Ada Popa*

Electrical Engineering and Computer Sciences
University of California, Berkeley

Technical Report No. UCB/EECS-2023-171

<http://www2.eecs.berkeley.edu/Pubs/TechRpts/2023/EECS-2023-171.html>

May 12, 2023



Copyright © 2023, by the author(s).
All rights reserved.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission.

Acknowledgement

I would like to thank my research advisor, Professor Raluca Ada Popa, for her guidance on this work and support in my academic career, as well as my research mentors, Weikeng Chen and Katerina Sotiraki, for guiding me throughout the entirety of research continuing from my undergraduate years. This work would not have been possible without the efforts of Weikeng Chen, Katerina Sotiraki, Murat Kantarcioglu, and Raluca Ada Popa who are equal partners in the authorship of this work. At last, I give my many thanks to the Sky Computing Lab, all of whom have helped me flourish, given me great feedback to improve as much as a researcher as possible. Finally I would like to thank my family for their continued support and love.

HOLMES: Efficient Distribution Testing for Secure Collaborative Learning

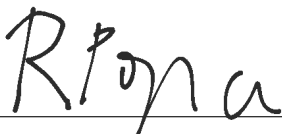
Ian Chang

Research Project

Submitted to the Department of Electrical Engineering and Computer Sciences,
University of California at Berkeley, in partial satisfaction of the requirements for
the degree of **Master of Science, Plan II**.

Approval for the Report and Comprehensive Examination:

Committee



Raluca Ada Popa
Research Advisor

May 12, 2023

(Date)

★ ★ ★ ★ ★ ★ ★



Natacha Crooks
Second Reader

May 12, 2023

(Date)

Abstract

Using *secure multiparty computation (MPC)*, organizations which own sensitive data (e.g., in healthcare, finance or law enforcement) can train machine learning models over their joint dataset without revealing their data to each other. At the same time, secure computation restricts operations on the joint dataset, which impedes computation to assess its quality. Without such an assessment, deploying a jointly trained model is potentially illegal. Regulations, such as the European Union’s General Data Protection Regulation (GDPR), require organizations to be legally responsible for the errors, bias, or discrimination caused by their machine learning models. Hence, testing data quality emerges as an *indispensable* step in secure collaborative learning. However, performing distribution testing is prohibitively expensive using current techniques, as shown in our experiments.

We present HOLMES, a protocol for performing distribution testing *efficiently*. In our experiments, compared with three non-trivial baselines, HOLMES achieves a speedup of more than $10\times$ for classical distribution tests and up to $10^4\times$ for multidimensional tests. The core of HOLMES is a hybrid protocol that integrates MPC with zero-knowledge proofs and a new ZK-friendly and naturally oblivious sketching algorithm for multidimensional tests, both with *significantly lower* computational complexity and concrete execution costs.

Acknowledgement

I would like to thank my research advisor, Professor Raluca Ada Popa, for her guidance on this work and support in my academic career, as well as my research mentors, Weikeng Chen and Katerina Sotiraki, for guiding me throughout the entirety of research continuing from my undergraduate years. This work would not have been possible without the efforts of Weikeng Chen, Katerina Sotiraki, Murat Kantarcioglu, and Raluca Ada Popa who are equal partners in the authorship of this work. At last, I give my many thanks to the Sky Computing Lab, all of whom have helped me flourish, given me great feedback to improve as much as a researcher as possible. Finally I would like to thank my family for their continued support and love.

HOLMES: Efficient Distribution Testing for Secure Collaborative Learning

Ian Chang
UC Berkeley

Katerina Sotiraki
UC Berkeley

Weikeng Chen
UC Berkeley & DZK Labs

Murat Kantarcioglu
University of Texas at Dallas & UC Berkeley

Raluca Ada Popa
UC Berkeley

Abstract

Using *secure multiparty computation (MPC)*, organizations which own sensitive data (e.g., in healthcare, finance or law enforcement) can train machine learning models over their joint dataset without revealing their data to each other. At the same time, secure computation restricts operations on the joint dataset, which impedes computation to assess its quality. Without such an assessment, deploying a jointly trained model is potentially illegal. Regulations, such as the European Union’s General Data Protection Regulation (GDPR), require organizations to be legally responsible for the errors, bias, or discrimination caused by their machine learning models. Hence, testing data quality emerges as an *indispensable* step in secure collaborative learning. However, performing distribution testing is prohibitively expensive using current techniques, as shown in our experiments.

We present HOLMES, a protocol for performing distribution testing *efficiently*. In our experiments, compared with three non-trivial baselines, HOLMES achieves a speedup of more than $10\times$ for classical distribution tests and up to $10^4\times$ for multidimensional tests. The core of HOLMES is a hybrid protocol that integrates MPC with zero-knowledge proofs and a new ZK-friendly and naturally oblivious sketching algorithm for multidimensional tests, both with *significantly lower* computational complexity and concrete execution costs.

1 Introduction

The MIT Technology Review article, “AI is sending people to jail—and getting it wrong” [1] is a reminder that machine learning models are determining people’s fate. The article explains how COMPAS¹, a system that rates people’s risk of future crime and decides if one should be held in jail before trial, has been shown to disproportionately target low-income populations and minorities [1]. Other examples where AI has infiltrated our everyday life include models to detect traditionally unrecognized anxiety and depression from speech [2, 3]

¹Correctional Offender Management Profiling for Alternative Sanctions.

or diagnose attention deficiency [4]. In these cases, accuracy is crucial for people to receive the right treatment.

However, if the training data is inaccurate, skewed, or affected by systemic biases, without any special attention to this issue, the trained model will also be biased [5]. There are many approaches (e.g., [6–9]) to guarantee fairness, starting from detecting and removing imbalances from the training data. For example, if a dataset has a large number of negative records (e.g., low credit scores) toward a certain group, one can reduce the imbalance by subsampling.

The situation in *secure collaborative learning*, where the model is trained in a way that none of the parties has access to the whole dataset, is vastly different. If the quality of the joint dataset is good, we expect the model trained on the joint dataset to be superior to models trained on individual datasets [10–14]. However, due to privacy considerations (e.g., due to GDPR [15]) organizations cannot know if the data is indeed of high quality. GDPR also requires organizations using such models to *prevent* errors, bias, and discrimination and *take liability* of the model [16]. Hence, organizations face the following conundrum. How can an organization take the (unknown) risk for another organization’s *untested* data?

Organizations will use collaborative secure computation only if there are *data quality* guarantees about the joint dataset. Hence, the ability to check data quality in secure computation is as important as data confidentiality and integrity. The first step toward this direction is to perform *distribution testing* over the joint dataset to examine:

- **one-dimensional properties**, such as the histogram of values of a specific attribute (e.g., income), or as basic as checking the range of each entry (e.g., age should be ≤ 200).
- **multidimensional properties**, such as whether the distribution of several attributes (e.g., age, gender, and income) *fits* into a desired distribution (e.g., represents a balanced demographic composition).

Distribution testing is a prominent method in statistical analysis. For instance, in clinical trials, the NIH Collaboratory [17] recommends comparing the distribution of different

datasets to detect data discrepancies. However, this useful tool is missing in prior works of secure collaborative learning (e.g., [18–21]), often left as an open problem. This is likely due to its *extremely high* overhead in MPC.

We present HOLMES, a protocol that performs such distribution testing efficiently, often at only a small fraction of the cost of secure collaborative learning. In our experiments, compared to three non-trivial baselines that we describe below, HOLMES achieves a speedup of more than $10\times$ for classical one-dimensional tests and up to $10^4\times$ for multidimensional tests. The core of HOLMES is a hybrid protocol that integrates MPC with zero-knowledge proofs and a new ZK-friendly, naturally oblivious sketching algorithm for multidimensional tests. HOLMES is already open-sourced in GitHub (anonymously): <https://github.com/holmes-inputcheck/>

1.1 Utilising IZK

Intuitively, to guarantee privacy no data-dependent computation should be performed outside the MPC protocol used in secure collaborative learning. So, bypassing the inefficiency of MPC seems impossible. Our insight is that the computation in distribution testing can be divided into parts that involve a single party’s dataset. Hence, distribution testing is mostly a *verification* task, instead of a direct computation. This brings us to the non-trivial first baseline described below.

First Baseline: Each party provides some auxiliary information based on their individual dataset, called *witness*, for the distribution tests. All parties verify each party’s computation using the witnesses and proceed to compute the distribution tests in MPC. Verifying a computation can be significantly faster than directly computing. For example, verifying that a value x is in the range $[a, b]$ typically involves computing the bit decomposition of $(x - a)$ and $(b - x)$. Hence, it becomes easier when each party provides this information as a witness.

However, this solution is still costly to implement in secure collaborative systems, when we require malicious security with a dishonest majority. SPDZ [22] and SCALE-MAMBA [23] are the fastest well-known MPC protocols for generic arithmetic computations with a dishonest and malicious majority, and are commonly used in secure collaborative learning (e.g., [19]). Since efficient distribution testing algorithms are based primarily on linear arithmetic operations, SPDZ-type protocols have lower computational overhead than other general MPC approaches, such as garbled circuits.

However, the cost of performing distribution tests still grows rapidly in this setting. Assume that t parties, each with the same amount of data, want to check the data quality of each individual dataset using distribution tests. We specifically use the standardized measure "wall-clock time" to describe the computational overhead, which refers to the amount of time **each** party takes until the MPC protocol finishes.

- In the best case, all t parties agree on the same set of distribution tests for all individual datasets. If the computation of

the tests on a single dataset requires C multiplications, the computation for t datasets has cost $C \cdot t$ multiplications, and the online phase of SCALE-MAMBA has computational overhead $O(C \cdot t)$ field elements². Moreover, the offline computational overhead of SCALE-MAMBA is proportional to the product of number of parties and the size of the online computation, i.e., $O(C \cdot t^2)$, since for each multiplication we need to produce t pieces of information, one for each party. Thus, running the distribution tests on t datasets in this t -party MPC leads to a wall-clock time of $O(C \cdot t^2)$.

- In the worst case, each party provides a different set of distribution tests for each individual dataset. Assuming that each set of tests has computational overhead $O(C)$ in SCALE-MAMBA, the online computational overhead using a t -party protocol is $O(C \cdot t^2)$; this is because the online cost for all sets of tests on each individual dataset is $O(C \cdot t)$. Taking into account the offline phase, running the distribution tests on all t datasets in this t -party MPC leads to a wall-clock time of $O(C \cdot t^3)$.

We present a solution that reduces the wall-clock time to $O(C \cdot t)$ in all cases.

HOLMES: efficiency via IZK. Zero-knowledge proofs are protocols that allow verifying a computation on a single party’s dataset, without revealing any other information. Thus, they are more suited for efficient distribution testing. Specifically, using interactive zero-knowledge (IZK) proofs, which are 2-party protocols involving a prover and a verifier, each party proves to $(t - 1)$ parties the distribution tests for its individual dataset, and verifies the tests for the other $(t - 1)$ datasets. In contrast to the first baseline, now the pairwise IZK protocols can be run concurrently, i.e., while \mathcal{P}_1 runs an IZK with \mathcal{P}_2 , the parties \mathcal{P}_3 and \mathcal{P}_4 can also run an IZK etc. Hence, if the computational overhead in IZK for the distribution tests on an individual dataset is $O(C)$, the wall-clock time for the datasets in IZK is $O(C \cdot t)$: each party proves $t - 1$, potentially different, statements and verifies $t - 1$ statements.

IZK reduces the computational overhead as performing distribution tests on an individual dataset, in essence, is not a computation that requires input from t parties; thus, it does not actually need t -party MPC. This new computational model leaves two challenges:

1. How can we ensure that the data in the MPC for collaborative learning is the same data used in IZK?
2. Are IZK-based solutions concretely efficient?

We answer the first question in §3.1 by proposing a lightweight consistency check. Additionally, our evaluation in §4 confirms that the answer to the second question is yes.

Remark: choice of IZK protocol. A generic approach for IZK is to use pairwise secure two-party computation (2PC) protocols. Namely, each pair of parties performs a 2PC protocol to show that a set of tests on its individual dataset passes

²Additions do not require communication in SCALE-MAMBA.

and to verify a (potentially different) set of tests on the other party’s dataset. However, there are also specialized IZK protocols that focus on the functionality of proving statements instead of a general two-party computation. In §4.4.1, we provide a comparison between HOLMES and other baselines. In HOLMES, we use QuickSilver [24] as the underlying IZK protocol. We show that the baseline of pairwise generic 2PC protocols, e.g., using SCALE-MAMBA, is up to $46\times$ slower. Additionally, we investigate baselines based on non-interactive ZK (NIZK) protocols, e.g., Spartan [25], which may asymptotically reduce the communication of each party to $O(C)$, as in this case each party only needs to produce a single proof for the tests on its individual dataset and broadcast it to the other $t - 1$ parties. However, NIZK protocols typically suffer from large proving times. Indeed, we find that producing an NIZK proof (i.e., using Spartan) is slower than performing $t - 1$ IZK protocols (i.e., using QuickSilver) for various dataset sizes and up to $t = 10$ parties. Since we are interested in the most efficient solution, including the computational overhead of each party, we use QuickSilver. We leave as an open problem to find an NIZK that is concretely more efficient in our setting.

1.2 Testing multidimensional properties

We now turn our focus to the task of *efficiently* testing multidimensional properties of distributions. For example, in a financial dataset, instead of focusing on a single attribute (i.e., one dimension) such as “debt”, we want to understand the properties of the joint distribution of several attributes (i.e., multiple dimensions), for example gender, age, race, income level, and debt. A property we might want to ensure is that the histogram of these attributes is not far from a balanced distribution, in which different genders, ages, races, and income levels are fairly represented.

In plaintext systems, this can be done by directly computing the histogram; namely, putting the data into buckets, where each bucket represents a different combination of values for the attributes gender, age, race, income level, and debt. Then, the distribution test is done by performing a Pearson’s χ^2 -test. In secure computation, the bucketing must be *oblivious*, in that it cannot reveal to which bucket a sample belongs. This either requires a linear scan of the buckets for each sample, as shown in Fig. 1, or oblivious RAM (ORAM) [26–33] in MPC, which theoretically has better complexity but in practice is concretely expensive.

The curse of oblivious bucketing. As shown in Fig. 1, in linear scan besides the real update (indicated by a black line), there are many dummy updates (indicated by red lines). In particular, the total number of updates is equal to the number of buckets. If the i -th attribute takes D_i distinct values, then the number of buckets is $\prod_i D_i = D_1 \cdot D_2 \cdot \dots \cdot D_d$, which quickly becomes prohibitively large. Therefore, the cost of oblivious bucketing for N data points is $O(N \cdot \prod_i D_i)$. Our experiments

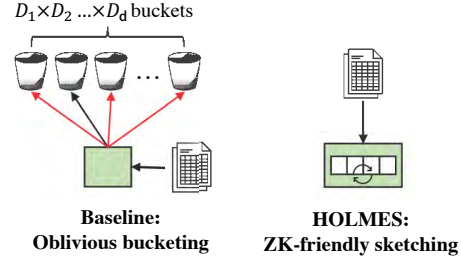


Figure 1: Methods for multidimensional distribution testing.

in §4 show that such oblivious bucketing for 5 attributes with values ranging from 1 to 10 takes 10^5 seconds.

Streaming and sketching to the rescue. We avoid oblivious bucketing by utilizing two concepts from algorithm design: *streaming* and *sketching*.

- **Streaming:** an algorithm that takes as input a sequence and only needs access to *limited* memory.
- **Sketching:** an algorithm that (approximately) performs the computation, using a *compressed* representation of the data.

We present a *sketching* algorithm that, given pseudorandomness, compresses the histogram of a dataset while preserving the necessary information for applying Pearson’s χ^2 -test. We also show how to compute this compressed representation in a *streaming* fashion, i.e., by accessing sequentially each entry in the dataset. This algorithm applies a *random linear projection* that approximately preserves the ℓ_2 -norm, according to the *Johnson-Lindenstrauss lemma* [34].

A challenge that arises in our algorithm is how to efficiently obtain pseudorandomness in IZK, as running classical pseudorandom functions, such as SHA-256, is impractical. Instead, we strive to find a *tailored* way to obtain pseudorandomness for our algorithm. We call this construction “ZK-friendly sketching”.

Finding pseudorandomness for random projection. Our random projection requires r pseudorandom maps with an one-bit output $b \in \{-1, 1\}$. Although we can use any pseudorandom function and extract one bit from it through bit decomposition, we find that many ZK-friendly hash functions [35, 36] are still costly. Instead, we discover that Legendre PRF, which has been studied recently [37–40] and is conjectured to be a universal one-way hash function (UOWHF) [41, 42] with a one-bit extractor, is a natural fit, and is extremely cheap—it only requires 8 input or multiplication gates in IZK.

With these techniques, HOLMES’s multidimensional distribution testing only requires $O(r \cdot N)$ computation, where r is the output size of the random linear projection. In our experiments, we show that this solution is up to $10^4\times$ faster compared to linear scan.

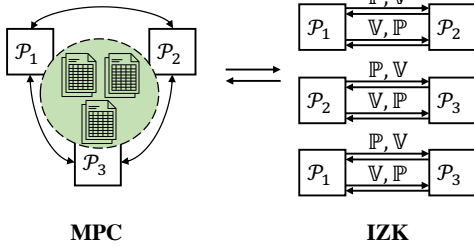


Figure 2: System model of HOLMES.

1.3 Summary of contributions

In summary, HOLMES’s contributions are:

- a new hybrid protocol that integrates MPC, IZK, and a lightweight consistency check for distribution testing, which has lower complexity and is concretely much more efficient (at least $10\times$ in our experiments) compared to three non-trivial baselines;
- a new efficient multidimensional distribution testing procedure via ZK-friendly sketching, which has lower complexity and is concretely much more efficient (up to $10^4\times$ in our experiment) compared to the naive multidimensional distribution testing.

Using HOLMES, we create a library of well-known statistical tests, such as z -test, t -test, F -test, and χ^2 -test. We also perform extensive experimental evaluation of HOLMES, including evaluation using real-world datasets from bank marketing [43, 44], healthcare [45, 46], and online auctions [47].

2 Protocol Overview

We assume that t parties want to participate in secure collaborative learning based on a t -party MPC protocol (e.g., SCALE-MAMBA). In HOLMES, the parties perform two types of computation: MPC and IZK, as shown in Fig. 2. During the distribution testing, parties participate in a t -party MPC, which is also used for the collaborative learning, and run IZK in a *pairwise* and *bidirectional* manner, where each party communicates with every other party, and both parties take turns as a prover and verifier to perform distribution tests on each other’s individual dataset. We require that MPC and IZK allow parties to load their datasets before learning the distribution tests, which prevents adaptively changing their input. We formalize this property in §3.1.

Tests. A distribution test is a predicate over an individual or joint (i.e., from multiple parties) dataset. Examples include well-known statistical tests, such as mean equality z -test (when the variance of the dataset is known) and t -test (when the variance is unknown), variance equality F -test, and Pearson’s χ^2 -test. These tests check a property between two populations, or between a population and a public distribution.

Each population can be an individual or a joint dataset and can contain data points with multiple attributes (e.g., age, gender, income, etc.). HOLMES implements these tests and various distribution test gadgets.

Workflow. HOLMES runs as a subroutine in the early stages of secure collaborative learning when all parties have loaded their datasets in MPC. HOLMES is invoked to perform distribution tests before MPC starts to run the actual training algorithm, as follows.

1. **Input loading in IZK:** Each party loads its dataset in IZK. This prevents the party from changing the input adaptively after seeing the revealed distribution tests.
2. **Revealing the distribution tests:** The parties reveal the distribution tests they want to perform.
3. **Consistency check:** Parties perform a consistency check (e.g., as in §3.1) to verify that the input loaded in IZK and MPC is the same. Parties reject if the check fails.
4. **IZK verification:** Each pair of parties acts as the prover and the verifier in IZK. The prover proves the correct calculation of some specified statistics about their dataset. The verifier verifies the proof and rejects if IZK fails.
5. **MPC finishing touches:** For distribution tests over joint datasets (e.g., z , t , F -tests), parties decide whether the data passes the test in MPC. Here, MPC only performs a small computation over *statistics* verified in IZK; thus we call it the “finishing touches”.

If all parties accept the distribution tests in HOLMES, the secure collaborative learning continues.

2.1 Protocol Components

HOLMES combines three underlying protocols:

1. an MPC protocol Π_{mpc} , which performs computations involving two or more parties;
2. an IZK protocol Π_{izk} , which verifies computation performed by a single party;
3. a consistency check (CC) protocol Π_{cc} (§3.1), which ensures that a party loads the same inputs in the MPC and IZK.

The underlying MPC and IZK protocols have to additionally satisfy the following properties.

Definition 1. An MPC protocol Π enables input-loading if in the beginning of the protocol each party \mathcal{P}_i with input \mathbf{x}_i , without access to the function \mathcal{F} , computes and sends $\text{load}_{\text{mpc}}(j, \mathbf{x}_i)$ to party \mathcal{P}_j . Then, the parties receive \mathcal{F} and proceed to the protocol with inputs $\text{load}_{\text{mpc}}(j, \mathbf{x}_i)$, without further access to the inputs \mathbf{x}_i .

Definition 2. An IZK protocol Π for the language $\mathcal{L} = \{x : \exists w \text{ s.t. } (w, x) \in \mathcal{R}\}$ enables input-loading if in the beginning of the protocol the prover sends $l := \text{load}_{\text{izk}}(w)$, where load_{izk} is a cryptographic commitment and during the protocol the

prover proves that $\varkappa \in \mathcal{L}$ and l is a commitment to w , i.e., runs an IZK for the language $\mathcal{L}' = \{(\varkappa, l) : \exists w \text{ s.t. } (w, \varkappa) \in \mathcal{R} \text{ and } l = \text{Com}(w)\}$.

We remark that the majority of MPC and IZK protocols are input-loading. For MPC, load_{mpc} typically corresponds to secret sharing (e.g. [23, 48, 49]).

The protocol of HOLMES is described in Fig. 3. We define the security of HOLMES in the real/ideal-world paradigm using the standard definition for (standalone) malicious security [50]. In malicious security, up to $t - 1$ of the parties can collude (statically) and arbitrarily deviate from the protocol. The ideal functionality $\mathcal{F}_{\text{HOLMES}}$ (Fig. 3) takes as input the list of distribution tests and the datasets; it outputs whether the datasets pass the tests or fail. Finally, in the full version of the paper [51] we provide the security proof of the following theorem.

Theorem 1. *If Π_{mpc} is a secure MPC protocol that enables input-loading, Π_{izk} is an IZK protocol that enables input-loading, and Π_{cc} is a consistency check protocol (see full version [51] for the formal definition), then HOLMES securely computes $\mathcal{F}_{\text{HOLMES}}$.*

Choice of protocols. HOLMES can be instantiated with a variety of MPC and IZK protocols, as long as they satisfy Definition 1 and Definition 2. In §3.1, we present a lightweight consistency check for MPC and IZK with additional, yet typical, properties, and in §4, we discuss how different protocol choices affect the efficiency of HOLMES.

Leakage from distribution tests. Note that any distribution testing leaks one-bit information – whether the test passed or failed. Hence, it is important that a party does not participate in distribution tests that may leak sensitive information. We leave as an open direction quantifying the privacy loss due to distribution tests, as discussed in §6.

Input-loading as a safeguard. Informally, the input-loading phase impedes parties from trying to tune their input data to pass the distribution tests. Even if parties abort after learning the distribution tests (e.g., pretending that the network is down) and request to redo the distribution testing, input-loading enforces that parties in the second execution must use the initial data.

3 Our Design

3.1 Consistency Check

HOLMES ensures that data used in MPC and IZK is the same via a *consistency check*, which is a protocol between a prover and $t - 1$ verifiers. We define this functionality formally in the full version of the paper [51]. We now describe our lightweight protocol for specific types of MPC and IZK, which we use in our implementation of HOLMES. The proof of Lemma 1 appears in the full version of the paper [51].

Lemma 1. *The construction below is a consistency check.*

Construction. Let Π_{mpc} be a secure MPC protocol over field \mathbb{F}_p that enables input-loading through additive secret sharing; namely, for an input $\mathbf{x} = (x_1, \dots, x_N) \in \mathbb{F}_p^N$, $(\text{load}_{\text{mpc}}(j, \mathbf{x}))_{j \in [t]} = (\text{share}_j(x_1), \dots, \text{share}_j(x_N))_{j \in [t]}$.

Let Π_{izk} be an IZK protocol over field \mathbb{F}_p that enables input-loading through homomorphic commitments; namely for an input $\mathbf{x} = (x_1, \dots, x_N) \in \mathbb{F}_p^N$, $(\text{load}_{\text{izk}}(j, \mathbf{x}))_{j \in [t]} = (\text{Com}_{\text{ck}_j}(x_1), \dots, \text{Com}_{\text{ck}_j}(x_N))_{j \in [t]}$. The protocol Π_{cc} works as follows.

1. Prover \mathbb{P} samples a random number $r \leftarrow \mathbb{F}_p$ and sends $(\text{load}_{\text{mpc}}(j, r), \text{load}_{\text{izk}}(j, r)) := (\text{share}_j(r), \text{Com}_{\text{ck}_j}(r))$ to each verifier \mathbb{V}_j .
2. Parties run a coin toss protocol [52] to obtain a random challenge $\beta \leftarrow \mathbb{F}_p$.
3. Parties run an MPC protocol to compute $\rho := r + \sum_{k=1}^N x_k \cdot \beta^k$. Note that each \mathbb{V}_j knows $(\text{share}_j(x_1), \dots, \text{share}_j(x_N), \text{share}_j(r))$.
4. Prover \mathbb{P} runs an IZK proof with each verifier \mathbb{V}_j for proving that $\text{Com}_{\text{ck}_j}(\rho) - \text{Com}_{\text{ck}_j}(r) - \sum_{k=1}^N \text{Com}_{\text{ck}_j}(x_k) \cdot \beta^k = \text{Com}_{\text{ck}_j}(0)$.
5. Verifiers accept if all IZK proofs are valid. Otherwise, they reject.

Cost analysis. Computing ρ and $\text{Com}_{\text{ck}_j}(\rho) - \text{Com}_{\text{ck}_j}(r) - \sum_{k=1}^N \text{Com}_{\text{ck}_j}(x_k) \cdot \beta^k$ requires $O(N)$ local operations since the MPC uses additive secret sharing and the commitment scheme is homomorphic. Hence, the main cost is due to the coin toss protocol and the IZK for proving that a value is a commitment to 0. Both of these functionalities are very efficient in state-of-the-art systems.

3.2 One-Dimensional Distribution Tests

We now provide more details about how IZK and MPC work together in each distribution test. We describe one-dimensional tests in this section, and present multidimensional tests in §3.3. We split the distribution tests into two steps: (1) compute statistical properties of an individual dataset, which is done in IZK, as described in §3.2.1; and (2) perform distribution tests using the computed statistical properties. If the distribution test is over an individual dataset, the second step is performed in IZK, but if it is over a joint dataset, the second step is done in MPC, as described in §3.2.2. We discuss the cost of these distribution tests in the full version of the paper [51].

3.2.1 IZK: Verifying basic statistics

HOLMES verifies basic statistical properties in IZK over an individual dataset: range, histogram, mean, variance, and trimmed mean.

Range. To prove that $a \leq x \leq b$ where $x \in \mathbb{F}_p$, the prover \mathbb{P} shows that $x - a \geq 0$ and $b - x \geq 0$. We describe in detail

HOLMES protocol

For t parties $\mathcal{P}_1, \dots, \mathcal{P}_t$ with inputs $(\mathbf{x}_1, \dots, \mathbf{x}_t)$ respectively, the protocol proceeds as follows:

1. **MPC input-loading:** Each party \mathcal{P}_i sends $\ell_{\text{mpc}}[i, j] := \text{load}_{\text{mpc}}(j, \mathbf{x}_i)$ to party \mathcal{P}_j for $j \in [t] \setminus \{i\}$.
2. **IZK input-loading:** Each party \mathcal{P}_i initiates $t - 1$ IZK protocols as a prover \mathbb{P} with $w = \mathbf{x}$ with parties $\mathbb{V} = \mathcal{P}_j$ for $j \in [t] \setminus \{i\}$. \mathcal{P}_i sends $\ell_{\text{izk}}[i, j] := \text{load}_{\text{izk}}(j, \mathbf{x}_i)$ to party \mathcal{P}_j for $j \in [t] \setminus \{i\}$, where $\text{load}_{\text{izk}}(j, \cdot)$ is the load_{izk} function of the IZK with $\mathbb{P} = \mathcal{P}_i$ and $\mathbb{V} = \mathcal{P}_j$.
3. **Distribution tests:** Each party \mathcal{P}_i sends a set of distribution tests tests_i to each party \mathcal{P}_j for $j \in [t] \setminus \{i\}$.
4. **Consistency check:** Each party \mathcal{P}_i performs Π_{cc} with the other $t - 1$ parties to show that $\ell_{\text{mpc}}[i, j] = \text{load}_{\text{mpc}}(j, \mathbf{x}_i)$ and $\ell_{\text{izk}}[i, j] = \text{load}_{\text{izk}}(j, \mathbf{x}_i)$.
5. **IZK protocol:** Each party \mathcal{P}_i performs Π_{izk} with \mathcal{P}_j for each test in tests_j as described in §3.2 and §3.3.
6. **MPC protocol:** Parties perform Π_{mpc} for each test in tests_i for $i \in [t]$ as described in §3.2 and §3.3. For each test, the parties either output pass or fail.

Ideal functionality $\mathcal{F}_{\text{HOLMES}}$

For t parties $\mathcal{P}_1, \dots, \mathcal{P}_t$ and a simulator Sim , $\mathcal{F}_{\text{HOLMES}}$ proceeds as follows: Upon receiving a message $(\mathcal{P}_i, \mathbf{x}_i, \text{tests}_i)$ from each of the t parties (or from Sim if that party is corrupted),

1. **Abort:** $\mathcal{F}_{\text{HOLMES}}$ awaits a message deliver or abort from Sim to decide whether the computation should move forward. $\mathcal{F}_{\text{HOLMES}}$ proceeds to the next step if the message is deliver. Otherwise, $\mathcal{F}_{\text{HOLMES}}$ sends abort to each \mathcal{P}_i and Sim , and halts.
2. **Distribution tests:** $\mathcal{F}_{\text{HOLMES}}$ runs the distribution tests specified in $(\text{tests}_1, \dots, \text{tests}_t)$ with respect to the inputs $(\mathbf{x}_1, \dots, \mathbf{x}_t)$. For each distribution test, $\mathcal{F}_{\text{HOLMES}}$ sends either pass or fail to each \mathcal{P}_i and Sim , and halts.

Figure 3: HOLMES protocol and ideal functionality $\mathcal{F}_{\text{HOLMES}}$.

the range check for a single element $x \in \mathbb{F}_p$ in the full version of the paper [51]. The gadget $\text{range}(\langle S \rangle, \text{attr}, [a, b])$ performs a range check on each data point of the population S with respect to attribute attr .

Histogram. The histogram of a population S for an attribute attr with values in \mathbb{F}_p counts the data points in a set of non-overlapping buckets. Each bucket might correspond to a single value or a range $[a, b]$. For instance, for the attribute “marital status”, we have single-value buckets (e.g., single, married, etc.), whereas for the attribute “age” we might be interested in range-buckets (e.g., 0-10, 11-20, etc.). In our setting, all values are elements in \mathbb{F}_p , so a single value a can be described by the range $[a, a]$. Hence, we focus on the case of range-buckets.

The prover \mathbb{P} first computes a one-hot encoding (OHE) $\vec{\sigma} = (\sigma_k)_{k=1}^D = (0, 0, \dots, 1, \dots, 0, 0)$ for each entry such that if σ_k is 1, the entry belongs to the k -th bucket. In the full version of the paper [51], we describe how the prover \mathbb{P} proves that $\vec{\sigma}$ is a valid one-hot encoding and how to perform the histogram check. We also discuss the extension to the multidimensional case, where each data point has d attributes (i.e., is in \mathbb{F}_p^d) and each bucket corresponds to d ranges $([a_i, b_i])_{i=1}^d$, one for each attribute. The gadget $\text{count_histogram}(\langle S \rangle, (\text{attr}_1, \dots, \text{attr}_d), (\mathbf{b}_1, \dots, \mathbf{b}_D)) \rightarrow \text{count}$ performs a histogram check on attributes $\text{attr}_1, \dots, \text{attr}_d$ for the population S with respect to buckets $\mathbf{b}_1, \dots, \mathbf{b}_D$.

Mean and variance. Mean and variance are essential in many tests, such as z -tests and t -tests. To prove that \bar{x} is

the mean of values $(x_j)_{j=1}^N$, the prover \mathbb{P} shows that $N \cdot \bar{x} \approx \sum_{j=1}^N x_j$. In practice, we want to keep a few decimal places for \bar{x} (e.g., $\bar{x} = 12.34$ with two decimal places). This is done by defining $\bar{x}' = 1234$, a fixed-point representation of \bar{x} , and asking \mathbb{P} to show that $N \cdot \bar{x}' \leq 100 \cdot \sum_{j=1}^N x_j < N \cdot (\bar{x}' + 1)$ using the range algorithm. To prove the correct calculation of the variance s^2 , \mathbb{P} first proves the calculation of the mean \bar{x} and of the mean of the square of each value, \bar{y} . The variance can be verified by checking that $s^2 \approx \frac{N}{N-1} (\bar{y} - \bar{x}^2)$.³ We provide the two gadgets $\text{mean}(\langle S \rangle, \text{attr}) \rightarrow \bar{x}$ and $\text{variance}(\langle S \rangle, \text{attr}) \rightarrow s^2$.

Trimmed mean. Trimmed mean is similar to mean, but it only considers entries with values within a certain range $[0, \theta]$. This statistic is useful as it can remove extreme values before computing the mean. This check combines range checks and a mean check as shown in the full version of the paper. HOLMES implements trimmed mean in the gadget $\text{trimmedMean}(\langle S \rangle, \text{attr}, \theta) \rightarrow \tilde{x}$.

3.2.2 MPC: Finishing touches

Distribution tests use the basic statistics of §3.2.1. If the tests involve an individual dataset, they are computed in IZK. If they involve dataset from multiple parties, the final computation is done in MPC. Since basic statistics are verified in IZK,

³The term $N/(N-1)$ corrects the bias of the variance because \bar{x} is computed from the data [53].

only the “finishing touches” (i.e., a small computation) are performed in MPC. We discuss how to perform well-known statistical tests: z -test, t -test, F -test, and Pearson’s χ^2 -test. Below we denote by yellow the basic statistics that depend on private datasets and have been verified in IZK.

z -test. This distribution test checks whether the means of two populations S_1 and S_2 (of size N_1, N_2) for the attribute attr are equal, assuming known variances. The parties provide the means \bar{x}_1 and \bar{x}_2 , and the test passes if:

$$(\bar{x}_1 - \bar{x}_2) / \sqrt{\sigma_1^2/N_1 + \sigma_2^2/N_2} \stackrel{?}{\leq} T_{\alpha, N_1, N_2},$$

where T_{α, N_1, N_2} is the critical value determined by the significance level α , N_1 , and N_2 and is computed outside of MPC.

t -test. This distribution test checks whether the means of two populations S_1 and S_2 for the attribute attr are equal, when the variances are not known. Parties provide the means \bar{x}_1 and \bar{x}_2 , and the variances s_1^2 and s_2^2 . The test passes if:

$$(\bar{x}_1 - \bar{x}_2) / \sqrt{s_1^2/N_1 + s_2^2/N_2} \stackrel{?}{\leq} T_{\alpha, \text{df}}$$

where $T_{\alpha, \text{df}}$ is the critical value determined by the significance level α and the degrees of freedom, which are defined as follows.

$$\text{df} = \frac{(s_1^2/N_1 + s_2^2/N_2)^2}{(s_1^2/N_1)^2/(N_1 - 1) + (s_2^2/N_2)^2/(N_2 - 1)},$$

The value $T_{\alpha, \text{df}}$ is computed in MPC using a lookup table for df . In our implementation, the lookup table for df ranges from 1 to 100. When $\text{df} > 100$, $T_{\alpha, \text{df}}$ is approximated by 1.645.

F -test. This distribution test checks whether the variances of populations S_1 and S_2 for the attribute attr are equal. Parties provide the variance s_1^2 and s_2^2 , and the test passes if:

$$s_1^2 / s_2^2 \stackrel{?}{\leq} T_{\alpha, N_1, N_2}$$

where T_{α, N_1, N_2} is determined by the significance level α , N_1 , and N_2 and is computed outside of MPC.

(One-dimensional) Pearson’s χ^2 -test. This distribution test checks whether the attribute attr of population S (which can be a joint dataset) follows a public distribution. Parties provide the histogram count of their joint dataset over the attribute attr which has D buckets, and the test passes if:

$$\sum_{j=1}^D (\text{count}[j] - Np_j)^2 / (Np_j) \stackrel{?}{\leq} T_{\alpha, D},$$

where N is the number of entries and p_j is the probability mass for the j -th bucket of the public distribution $\vec{p} = (p_1, \dots, p_D)$. The critical value $T_{\alpha, D}$ is determined by the significance level α and D and is computed outside of MPC.

3.2.3 Subsampling with malicious security

HOLMES allows distribution tests to be performed on a *random* subset of the dataset, which is decided after the data has been loaded to IZK and MPC. The random subset is chosen using a pseudorandom function, with a seed that comes from a coin toss protocol among the t parties [52]. Though this might sacrifice accuracy, it boosts efficiency and allows more tests to be performed with a given computational budget. We leave as an open direction identifying applications where the subsampling is beneficial, as discussed in §6.

3.3 Multidimensional Distribution Tests

We now discuss the setting where we want to test the distribution over multiple attributes (i.e., dimensions). Particularly, we want to test if the distribution of a dataset is close to a public distribution (e.g., a balanced distribution where different groups are represented appropriately) using Pearson’s χ^2 -test. Note that in this case, the number of buckets is $\prod_{i=1}^d D_i$ where D_i is the number of buckets of the i -th attribute.

Baseline: multidimensional bucketing. We can naturally extend the one-dimensional Pearson’s χ^2 -test by creating multidimensional buckets. In particular, given the histogram count over the attributes $(\text{attr}_1, \dots, \text{attr}_d)$, the test checks if:

$$\sum_{j=1}^{D_1 \cdot D_2 \cdot \dots \cdot D_d} (\text{count}[j] - Np_j)^2 / (Np_j) \stackrel{?}{\leq} T_{\alpha, D},$$

where N is size of population S , D_j is the number of distinct buckets of the j -th attribute attr_j , $D = \prod_{i=1}^d D_i$, and p_j is the probability mass for the j -th bucket of the public distribution $\vec{p} = (p_1, \dots, p_D)$. The critical value $T_{\alpha, D}$ is determined outside of MPC by the significance level α and the number of buckets D . We illustrate this test in Fig. 4.

Cost analysis of the baseline. This baseline becomes impractical when the number of buckets D is high.

- **IZK cost:** A multidimensional histogram with D buckets is computed obliviously. This requires an arithmetic circuit of size $O(N \cdot (D + d\ell))$, where each bucket contains ranges of size at most 2^ℓ , and becomes impractical when D is large.
- **MPC cost:** The MPC performs the final computation for Pearson’s χ^2 -test, which involves the histogram count of length D , which has been verified in IZK. The cost in MPC is $O(D \cdot \text{cost}_{\div})$ operations and one comparison, where cost_{\div} is the cost of division in MPC.

The linear growth with respect to D is discouraging. In our experiments in §4.4.3, performing distribution testing over four attributes—age, jobs, marital status, and education—results in $D = 37,500$ and takes 10^5 seconds to compute.

New test: unnormalized χ^2 -test. HOLMES uses another test for goodness-of-fit, called unnormalized χ^2 -test inspired by the work of Arias-Castro, Pelletier, and Saligrama [54]. This test has a more complicated critical value, but it requires no divisions. The test checks if:

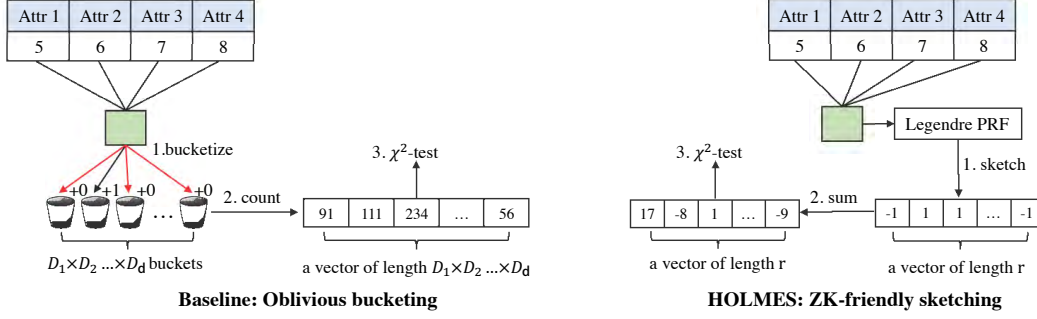


Figure 4: Two methods of multidimensional distribution testing.

$$\sum_{j=1}^{D_1 \cdot D_2 \cdot \dots \cdot D_d} (\text{count}[j] - N p_j)^2 \stackrel{?}{\leq} T_{\alpha, N, p_1, p_2, \dots, p_D},$$

where the critical value $T_{\alpha, N, p_1, p_2, \dots, p_D}$ is computed from a variant of the generalized χ^2 distribution with parameters (p_1, \dots, p_D) [55–57]. We provide the details of the statistical test in the full version of the paper [51]. Since the parameters are public, $T_{\alpha, N, p_1, p_2, \dots, p_D}$ is computed outside MPC.

Workflow in HOLMES. Naively computing the unnormalized χ^2 -test has the same IZK cost as the baseline, and hence it is still prohibitively expensive for large D . We use sketching and pseudorandom functions (PRFs) to reduce this cost. The overall workflow for computing the unnormalized χ^2 -test is also illustrated in Fig. 4. We show why our approach approximates well the unnormalized χ^2 -test in §3.3.1.

- **PRF keys:** After the input-loading in IZK and MPC, parties run a coin toss protocol to sample r keys for the Legendre PRF, denoted by k_1, k_2, \dots, k_r . These PRF keys will be used to produce a pseudorandom matrix consisting of values in $\{-1, 1\}$ for the random linear projection. We expand upon our sketching using random linear projections in §3.3.1. Details about the Legendre PRF can be found in the full version of the paper [51].
 - **IZK sketching:** Each party \mathcal{P}_i whose dataset is involved in the computation proves the following in IZK:
 1. computation of the linear index o (i.e., the index of value 1 in the one-hot encoding) of each data point for the attributes $(\text{attr}_1, \dots, \text{attr}_d)$;
 2. computation of $\text{PRF}_{k_j}(o)$ for each $j \in [r]$ and the linear index o of each data point. For the k -th data point, this result in a r -vector \vec{u}_k consisting of elements in $\{-1, 1\}$;
 3. computation of $\vec{\text{sum}}_i$ as $\sum_{k=1}^{N_i} \vec{u}_k$, where N_i is the dataset size of \mathcal{P}_i .
 - **MPC finishing touches:** Similarly to the other tests, the computation in MPC is lightweight:
 1. compute $\vec{\text{sum}} = \sum_{i=1}^t \vec{\text{sum}}_i$, where t is the number of parties whose dataset is involved in the test;
 2. check if $\sum_{v=1}^r (\text{sum}[v] - q_v)^2 \stackrel{?}{\leq} r \cdot T_{\alpha, N, p_1, p_2, \dots, p_D}$.
- The vector \vec{q} is computed outside of MPC as follows: Let $\mathbf{R} \in \mathbb{F}_p^{r \times D}$ be a matrix such that $\mathbf{R}[i][j] = \text{PRF}_{k_i}(j)$, then

$$\vec{q} = N \cdot \mathbf{R} \cdot \vec{p}.$$

Cost analysis of HOLMES’s approach. The IZK requires an arithmetic circuit of size $O(N \cdot r \cdot \text{cost}_{\text{PRF}})$ for a joint dataset of size N , where cost_{PRF} is the cost of a PRF evaluation in IZK. The cost in MPC is $O(t + r)$ operations in \mathbb{F}_p . The cost of computing \vec{q} is $O(D \cdot r \cdot \text{cost}_{\text{PRF}})$ local operations. However, this computation is public, so its cost is negligible. In our experiments §4.4.2, HOLMES’s approach can be 10^4 times faster than the baseline.

3.3.1 Why sketching works?

We now explain why HOLMES’s approach approximates well the unnormalized χ^2 -test.

Approximating the unnormalized χ^2 -test. Observe that the unnormalized χ^2 -test compares the Euclidean distance of $\vec{\text{count}}$ and $N \vec{p}$, $\text{dist}(\vec{\text{count}}, N \vec{p})$, to the value $T_{\alpha, N, p_1, p_2, \dots, p_D}$. From well-known results in statistics [34, 58] which we review in the full version of the paper [51], the Euclidean distance of two D -dimensional vectors \vec{x} and \vec{y} , $\text{dist}(\vec{x}, \vec{y})$, can be approximated by $\text{dist}(\mathbf{R} \cdot \vec{x}, \mathbf{R} \cdot \vec{y})/r$ where $\mathbf{R} \cdot \vec{x}$ and \mathbf{R} is a $r \times D$ matrix with entries (pseudo-)randomly chosen from $\{-1, 1\}$ and r is sufficiently large, but independent of D . Hence, for suitable r

$$\text{dist}(\vec{\text{count}}, N \vec{p}) \approx \text{dist}(\mathbf{R} \cdot \vec{\text{count}}, \vec{q})/r.$$

Proving that $\vec{\text{sum}} = \mathbf{R} \cdot \vec{\text{count}}$. The histogram $\vec{\text{count}}$ stores the number of elements in each bucket. Hence, if $\vec{\sigma}_k$ is the one-hot encoding of the k -th value in the dataset, then

$$\vec{\text{count}} = \sum_{k=1}^N \vec{\sigma}_k.$$

Using this equality, it follows that

$$\mathbf{R} \cdot \vec{\text{count}} = \mathbf{R} \cdot \sum_{k=1}^N \vec{\sigma}_k = \sum_{k=1}^N \mathbf{R} \cdot \vec{\sigma}_k.$$

Since $\vec{\sigma}_k$ is a one-hot encoding, it contains a single element of value 1. Let o_k be the index of the element of value 1; we call o_k the linear index of $\vec{\sigma}_k$. Then,

$$\mathbf{R} \cdot \vec{\sigma}_k = \mathbf{R}[o_k],$$

where $\mathbf{R}[o_k]$ is the o_k -th column of the matrix \mathbf{R} . In our sketching, the matrix \mathbf{R} is pseudorandom, and each element $\mathbf{R}[i][j]$ is equal to $\text{PRF}_{k_i}(j)$, i.e., the evaluation on input j of a pseudorandom function with key k_i . Hence,

$$\mathbf{R}[o_k] = (\text{PRF}_{k_1}(o_k), \dots, \text{PRF}_{k_r}(o_k)),$$

which is by definition equal to \vec{u}_k . Overall, since $\vec{\text{sum}} = \sum_{k=1}^N \vec{u}_k$, it holds that

$$\vec{\text{sum}} = \mathbf{R} \cdot \vec{\text{count}}.$$

Choice of parameter r . In our implementation, we choose $r = 200$, which empirically results in 1.1 approximation factor. We discuss this choice in the full version of the paper [51].

3.3.2 Choosing an IZK-friendly PRF

Why use a PRF? The naive solution for our sketching is to sample \mathbf{R} with random elements, without the special structure related to PRF evaluations. In this case, for the k -th data point, party \mathcal{P}_i provides the computation $\mathbf{R}[o_k]$ obviously, i.e., performs a lookup on the D columns of the *random matrix* \mathbf{R} . Using a linear scan, this requires an arithmetic circuit of size $O(N \cdot r \cdot D)$ in IZK. In our IZK sketching, the matrix \mathbf{R} is produced from pseudorandom functions and this allows computing $\mathbf{R}[o_k]$ directly, without a linear scan. Hence, for a dataset to size N , \mathcal{P}_i produces a proof that $\vec{\text{sum}}_i$ was computed correctly with $N \cdot r$ PRF evaluations. When the PRF evaluation cost in IZK is small, our solution is more efficient.

Our choice: Legendre PRF. A concern with PRF evaluations in IZK is that the cost for common PRFs (e.g., SHA-256) is prohibitive. Thus, new ZK-friendly PRFs have been developed, e.g., Rescue [36] and Poseidon [35]. In our ZK-friendly sketching, we identify Legendre PRF [37–40] as the most suitable choice. Recall that in the sketching algorithm, the output of each PRF evaluation is an element in $\{-1, 1\}$. Legendre PRF, whose output is the Legendre symbol of a value modulo a prime, already has this property. In contrast, for other PRFs we need to extract these bits from a longer output, which incurs extra cost. We provide details about the Legendre PRF and compare its cost with other PRFs in the full version of the paper.

4 Implementation and Evaluation

In this section, we present and discuss the evaluation results of HOLMES, which answer the following questions:

- How well do HOLMES’s distribution tests hold up against corruptions to both simulated and real-world data? (§4.3)
- How does HOLMES compare to the baselines, as well as alternative efficient system implementations? (§4.4)
- What is the overhead of HOLMES on real-world datasets? What contributes to this overhead? (§4.4.3)

4.1 Setup

We run our experiments on AWS c5.9xlarge instances, each with 36 cores and 72 GB memory. Each party has its own c5.9xlarge instance. We limit each instance’s bandwidth to 2 Gbps and add a round-trip latency of 20 ms. We standardize data inputs across all protocols as field elements in \mathbb{F}_p where $p = 2^{62} - 2^{16} + 1$. Text labels of an attribute are mapped and converted to field elements in \mathbb{F}_p . A d -dimensional input is formed as a vector of d field elements in \mathbb{F}_p , where the k -th vector entry represents the range bucket that the data point falls into for the k -th attribute. Decimals arising from divisions are stored in fixed-point representation (§3.2.1), where we multiply the operand by 10^2 to achieve a precision up to two decimal places by default. This fixed-point accuracy can be easily changed anytime by the parties for their use case.

General parameters for all setups include the statistical security parameter $\lambda = 30$, computational security parameter $\kappa = 128$, the input size, and the fixed-point accuracy.

HOLMES and the baselines are implemented using state-of-the-art cryptographic libraries, as follows.

HOLMES. We use QuickSilver [24] due to the lower prover overhead for IZK. The version of QuickSilver we use has integrated the latest techniques in Silver [59]. The number of concurrent threads run in a single prover-verifier protocol is defaulted to 32 to maximize multithreading. The input form that we use is the `IntFp` datatype, which represents a field element in \mathbb{F}_p .

We use SCALE-MAMBA [23] and MP-SPDZ [22, 60] for MPC, where the Low Gear protocol in MP-SPDZ is used for the offline phase of MPC and SCALE-MAMBA is used for the online phase. They are the state-of-the-art MPC protocols for arithmetic computations over large prime fields with a dishonest and malicious majority. For MP-SPDZ and SCALE-MAMBA, we use the Full Threshold Linear Secret Sharing Scheme with prime set to $p = 2^{62} - 2^{16} + 1$. Furthermore, we compile our circuits with the `sint` bit length limited to 32 bits, statistical security parameter $\lambda = 30$, and prime modulus size limited to 64 bits.

We compare HOLMES with three baselines—generic MPC, pairwise generic 2PC, NIZK/SNARK—to quantify the efficiency advantages. We now describe the setup of these systems.

Baseline 1: Generic MPC. The baseline runs HOLMES’s MPC setup in entirety, using SCALE-MAMBA and MP-SPDZ. The data are only loaded once, and since there is no other auxiliary protocol there is no need for the consistency check.

Baseline 2: Pairwise 2PC. Each pair of parties runs a 2PC with the same setup as the generic MPC baseline using SCALE-MAMBA and MP-SPDZ. We instantiate a single party to host and execute 2PC protocols to all other parties on concurrent threads. Each individual 2PC protocol is run on a separate network port.

Baseline 3: NIZK & SNARK. We use a state-of-the-art NIZK and SNARK system, Spartan [25], with low communication overhead and small verification time. To utilize data-parallel circuit uniformity in Spartan, we copy the subcircuits of the distribution tests for each data entry. Furthermore, we parallelize verifications of NIZK/SNARK proofs on concurrent threads. We adapt the 62-bit field with modulus $p = 2^{62} - 2^{16} + 1$.

4.2 Artifacts

We released HOLMES in an open-sourced anonymous repository in GitHub. The implementation consists of three parts:

- **Compute engine:** The original QuickSilver is not compatible with many efficient MPC protocols because it works on a special prime field.⁴ We perform an extensive search for a prime with low Hamming weight that is compatible with such MPC preprocessing protocols, and we settle to $p = 2^{62} - 2^{16} + 1$. We contribute a fork of EMP-ZK, called EMP-ZK-HOLMES⁵, which includes a highly tuned, specialized implementation for modular reduction and learning-parity-with-noise (LPN) map for this prime.
- **Distribution tests:** We implement distribution tests for range, histogram, mean, variance, trimmed mean, z -test, t -test, F -test, and χ^2 -test, including both oblivious bucketing and our ZK-friendly sketching. The codebase also includes integration tests, unit tests, and individual benchmarks.⁶
- **Examples and benchmarks:** We assemble distribution tests for the baselines and three real-world datasets (described in §4.4.3) and benchmark their performance. We also include accuracy evaluations for HOLMES’ statistical tests against corruptions to simulated and real-world data. Finally, we provide a QuickSilver-to-SCALE-MAMBA source-to-source compiler and an online-only SCALE-MAMBA for ease of benchmarking with the baselines.

4.3 Accuracy evaluation

We show how HOLMES performs in face of specific corruption scenarios. Our goal is to better understand the fraction of data that needs to be corrupted for a distribution test to fail. Our statistical testing suite for the accuracy evaluation includes the z -test, t -test, F -test, naive normalized χ^2 -test, naive unnormalized χ^2 -test, and HOLMES’s ZK-friendly sketching χ^2 -test. We evaluate the accuracy on both simulated and real datasets.

General setup. We start with a dataset of positive integer values. We randomly divide the dataset into two equal parts; one

⁴QuickSilver is restricted to a Mersenne prime $p = 2^{61} - 1$. However, this prime is not compatible with MPC preprocessing protocols based on ring learning-with-error (LWE) and would force the MPC to choose preprocessing protocols based on oblivious transfer, which are slower.

⁵<https://github.com/holmes-inputcheck/emp-zk>

⁶<https://github.com/holmes-inputcheck/holmes-library>

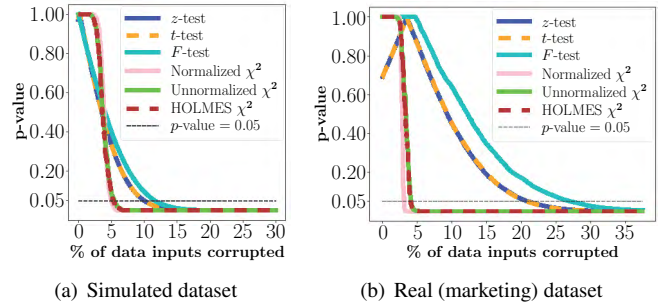


Figure 5: Accuracy of distribution tests after corrupting the dataset

part—which we call the “corrupted dataset”—is the dataset controlled by the adversary, and the other is an honest dataset. Our corruption model is as follows: at each iteration, we randomly select a data point from the corrupted dataset; we then modify, depending on the test, its value by 1, assuming that it remains within the acceptable bounds. For t and z -tests, we add 1 to the value, for F and χ^2 -test we choose to either increase or decrease by 1. Note that this is a minimal amount of corruption per iteration. For instance, for the χ^2 -tests, this corresponds to moving the data point to the next or previous bucket. The two populations in each test refer to the corrupted and the honest dataset⁷.

We compute the test statistic (i.e., the value before the final comparison with T_α) via the formulas of §3.2.2 and §3.3. We then calculate the p -value using inverse cumulative probability functions⁸. Typically, a test fails if the p -value is less than 0.05, which corresponds to a significance level $\alpha = 0.95$.

Simulated dataset. We simulate a dataset of randomly sampled values from $\mathcal{N}(\mu, \sigma^2)$, the normal distribution with mean μ and variance σ^2 . We sample 40000 entries with $\mu = 17$ and $\sigma = 10$ to prevent field integer overflow and underflow while ensuring a diverse set of values in the dataset. For our χ^2 -tests, we establish a histogram with 35 buckets ($[0, 1], [1, 2], \dots, [35, \infty)$) to avoid excessive outliers at the boundary buckets (i.e. $[0, 1]$ and $[35, \infty)$), and ensure there are no empty buckets. For the other tests, we use the parameters $\mu = 17$ and $\sigma = 10$ and we calculate \bar{x}_2 and s_2^2 using the honest dataset. We graph the statistical p -values after corruptions in Fig. 5(a).

Real dataset. We use a bank marketing dataset [43, 44], which we describe further and time in our evaluation of HOLMES on real-world datasets in §4.4.3. We perform the z -test, t -test, and F -test over the telemarketing call duration

⁷In χ^2 -test, where there is a single population, we use the entire dataset to compute the public distribution \vec{p} and we ignore the honest dataset.

⁸For the unnormalized χ^2 -test with or without ZK-friendly sketching, we compute a lookup table for the inverse cumulative probability functions as outlined in the full version of the paper [51]. For the lookup table, we use 5000 random samples.

attribute. We perform the χ^2 -tests (normalized, unnormalized with and without our sketching) over the attributes age, job, educational level, and marital status. We graph the statistical p -values after corruptions in Fig. 5(b).

Results. For the simulated dataset, the z -test and the t -test, both of which test the mean between two populations, fail at around 10% corruptions and they follow the same trend as a function of the fraction of corruptions. The normalized, unnormalized, and HOLMES χ^2 -tests follow the same trend, which confirms the accuracy of HOLMES χ^2 -test. The χ^2 -tests fail when approximately 5% of the dataset is corrupted, much faster than all other statistical tests. Finally, the F -test is the most robust to corruptions and fails at about 13% corruptions.

For the marketing dataset, the z -test and the t -test have the same trend. However, now there might be an initial increase in p -values, depending on the random split of the dataset into honest and corrupted, since we use the honest dataset to compute the underlying properties. For instance, if the honest dataset has initially larger mean than the corrupted, increasing the values of data points initially leads to an increase of the p -value. Hence, the p -values drop at an offsetted percentage of corrupted data entries for the z -test, t -test, and F -test. Similarly to the simulated dataset, the χ^2 -tests drop off at the same rate and lead to test failures at approximately the same point.

4.4 Evaluation discussion

We evaluate the overheads for our distribution tests. Our results are as follows: the histogram in HOLMES is about 5–11 \times more efficient than generic MPC; the mean and variance are 2–3 \times more efficient; and the trimmed mean is about 5–10 \times more efficient.

Next, in §4.4.1 we compare the overhead of range checks and ZK-friendly sketching with alternative applicable systems; these are the two most expensive gadgets supported in HOLMES. Range check is the main source of overhead for many HOLMES gadgets (§3.2.1, etc.), while ZK-friendly sketching is the bottleneck for HOLMES’s χ^2 -test. Later on, in §4.4.2 we depict the drastic overhead reduction of HOLMES’s χ^2 -test over the naive χ^2 -test. Finally, in §4.4.3 we show that HOLMES performs efficiently in practical settings compared to our generic MPC baseline by modeling distribution test workflows on real-world datasets.

4.4.1 Comparison of HOLMES with the Baselines

We evaluate range checks and ZK-friendly sketching for number of parties $t = 2, 6, 10$ and input size per party $N_{\text{ind}} = 100\text{k}, 200\text{k}, 500\text{k}$.

Range and HOLMES’s χ^2 -test Setup. For the range check, we vary the sizes of the range $[a, b]$ (i.e., $b - a$) as 2^ℓ for $\ell \in \{8, 12, 16, 20, 24\}$. We run the range check algorithm with these inputs and parameters. Our results are listed in Tab. 1. In the ZK-friendly sketching, for simplicity, we assume that

all attributes take the same number of distinct values. We consider the number of attributes $d \in \{2, 4\}$, and vary the buckets per attribute as $D_0 = \dots = D_d \in \{5, 10, 50\}$. We perform the IZK check described in §3.3, i.e, for each data point, we compute its linear index (detailed in the full version of the paper [51]) and feed it into $r = 200$ Legendre PRFs with polynomial degree 3, and quadratic nonresidue $7 \in \mathbb{F}_p$; the $r = 200$ unique keys are generated from a random oracle based on SHA-256. We run the Legendre PRF algorithm with these inputs and parameters. Our results are listed in Tab. 3.

Baseline 1: Generic MPC . The overhead grows quadratically in the number of parties and linearly in the input size; hence, generic MPC is the slowest baseline in our comparison. The baseline is 10–256 \times and 35–198 \times slower than QuickSilver for the range check and the ZK-friendly sketching, respectively.

Baseline 2: Pairwise 2PC . Pairwise 2PC provides higher throughput than generic MPC due to the parallelization of the offline phases and online phases. It also has lower latency than generic MPC, since 2PC reduces the creation of authenticated shares from the entire t -party combined dataset of size $t \cdot N_{\text{ind}}$ to the two-party combined dataset of size $2 \cdot N_{\text{ind}}$. MP-SPDZ still needs to preprocess $t - 1$ different inputs of size $2 \cdot N_{\text{ind}}$, and as a result, the preprocessing phase contributes to most of the overhead. In sum, the overhead grows linearly to the number of parties and linearly to the input size.

Pairwise 2PC is faster than generic MPC, but slower than HOLMES. Pairwise 2PC is 4–32 \times slower for the range check and 13–36 \times slower for the ZK-friendly sketching than QuickSilver. For 10 parties, 2PC is 18 \times and 13 \times faster than generic MPC for the range check and the ZK-friendly sketching, respectively. For 6 parties, these numbers become 10 \times and 8 \times , depicting a speedup factor of around $O(t)$ over generic MPC.

Baseline 3: NIZK & SNARK . Spartan_{NIZK} is the second fastest system behind HOLMES. NIZK and SNARK systems scale well for a large number of parties, since each party only generates a single proof for its dataset, and parties can concurrently verify other parties’ proofs on multiple cores. Thus, this approach has overhead that remains relatively constant to the number of parties (up to the number of threads in our machine). However, the overhead to prove dense circuits with lots of constraints is large relative to arithmetic-based IZKs, and still grows linearly to the input size; for instance, in the two-party case, Spartan_{NIZK} is 2.4–16 \times slower for the range check and 4–45 \times slower for the ZK-friendly sketching than QuickSilver. Hence, for extremely large and dense circuits (e.g. ZK-friendly sketching), we extrapolate our small input size experiments to larger input sizes.

For a small number of parties, Spartan_{NIZK} is quite inefficient and has speeds comparable to pairwise 2PC. However, for larger numbers of parties, e.g., 10, it is approximately 4–5 \times and 3–4 \times faster for the range check and the ZK-friendly sketching, respectively, than pairwise 2PC. Spartan_{NIZK} is

Table 1: Overhead of range checking on different protocols with varying input sizes. For a range $[a, b]$, the range size is $b - a$. $\text{Spartan}_{\text{SNARK}}$ is significantly slower than $\text{Spartan}_{\text{NIZK}}$, so we omit benchmarks for input sizes 200k and 500k.

N_{ind}	Range Size	Number of parties = 2					Number of parties = 6					Number of parties = 10				
		2^8	2^{12}	2^{16}	2^{20}	2^{24}	2^8	2^{12}	2^{16}	2^{20}	2^{24}	2^8	2^{12}	2^{16}	2^{20}	2^{24}
100k	QuickSilver	3.4s	3.7s	4.2s	4.7s	4.9s	17.0s	18.6s	21.0s	23.4s	24.5s	30.6s	33.4s	37.8s	42.1s	44.1s
	Paired 2PC	35.0s	51.9s	73.2s	85.9s	103.4s	88.5s	133.1s	173.3s	216.4s	267.9s	145.3s	217.0s	288.7s	359.3s	432.3s
	MPC	35.0s	51.9s	73.2s	85.9s	103.4s	894.7s	1309s	1739s	2134s	2568s	2566s	3795s	5024s	6241s	7603s
	$\text{Spartan}_{\text{NIZK}}$	30.4s	48.5s	59.9s	77.1s	108.0s	30.8s	49.0s	60.3s	77.6s	108.5s	31.5s	49.2s	60.6s	77.9s	109.8s
	$\text{Spartan}_{\text{SNARK}}$	275.5s	526.5s	526.9s	528.0s	1070s	276.7s	528.7s	529.2s	530.2s	1073s	278.4s	530.9s	531.4s	532.4s	1076s
200k	QuickSilver	4.7s	5.6s	6.2s	7.2s	7.8s	23.7s	27.8s	30.9s	36.1s	38.9s	42.6s	50.1s	55.6s	65.1s	70.1s
	Paired 2PC	71.0s	104.5s	139.4s	173.7s	207.9s	178.9s	265.9s	345.2s	439.4s	523.3s	293.4s	433.3s	578.9s	716.8s	864.4s
	MPC	71.0s	104.5s	139.4s	173.7s	207.9s	1757s	2682s	3423s	4278s	5303s	5092s	7576s	10142s	12691s	15064s
	$\text{Spartan}_{\text{NIZK}}$	66.2s	108.8s	130.9s	153.6s	208.2s	66.1s	108.6s	131.0s	154.2s	208.5s	66.4s	107.5s	132.1s	152.5s	210s
	QuickSilver	9.3s	10.3s	12.9s	14.8s	16.5s	46.7s	51.3s	64.4s	74.3s	82.4s	84.1s	92.4s	115.9s	133.7s	148.3s
500k	Paired 2PC	177.3s	263.7s	361.1s	441.1s	530.9s	443.3s	663.7s	872.9s	1081.6s	1294.4s	729.9s	1092s	1446s	1797s	2142s
	MPC	177.3s	263.7s	361.1s	441.1s	530.9s	4370s	6543s	8799s	10707s	13405s	12792s	19411s	25299s	31333s	38063s
	$\text{Spartan}_{\text{NIZK}}$	155.4s	244.3s	299.7s	356.5s	483.0s	155.5s	244.7s	300.1s	356.8s	483.3s	158.1s	242.6s	305.8s	357.3s	487.0s

1–3× slower for the range check and 4–5× slower for the ZK-friendly sketching than QuickSilver.

$\text{Spartan}_{\text{SNARK}}$ has succinct proof size and verification time but with massive prover overhead. Even for the most complicated benchmark, the ZK-friendly sketching with $N_{\text{ind}} = 100k$, verification time is ~ 1 –2s. However, the computational overhead for the same circuit in $\text{Spartan}_{\text{SNARK}}$ is around 9× larger than $\text{Spartan}_{\text{NIZK}}$ and 2× slower than generic MPC.

4.4.2 Cost of multidimensional tests

We now measure the efficiency of HOLMES’s multidimensional χ^2 -test. Namely, we test HOLMES’s ZK-friendly sketching against the oblivious bucketing approach, which we call the "naive χ^2 -test". For both HOLMES χ^2 -test and naive χ^2 -test, we run the benchmarks entirely in our choice of IZK: QuickSilver. We vary the number of attributes as $d \in \{2, 3, 4, 5\}$, and the buckets per attributes as $D_0 = \dots = D_d \in \{5, 10, 15, 20, 25\}$. Since the naive approach is extremely expensive, we were able to run the experiment only on a small scale, so extrapolate our small scale experiments to a larger scale. In the full version of the paper [51], we show that HOLMES and the naive χ^2 -test have drastically different growth patterns. The overhead of the naive χ^2 -test grows exponentially in the number of buckets per attribute. The overhead of HOLMES multidimensional χ^2 -test is no longer dominated by the number of multidimensional buckets as in the naive χ^2 -tests; instead, the new overhead of our sketching approach, $O(N \cdot r \cdot \text{cost}_{\text{PRF}})$, is now dominated by the input size and the number of PRF keys. Based on our experiments, we find that HOLMES’s sketching approach greatly improves the efficiency of multidimensional tests when the number of attributes and distinct values per attribute are large. For instance, when $d = 4$ and $D_i = 25$, we observe an efficiency increase of around $10^4 \times$.

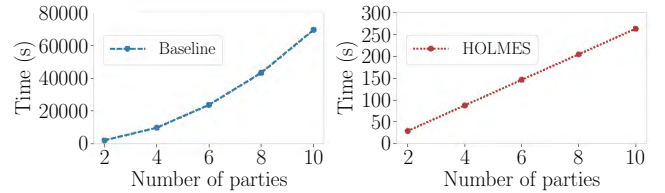


Figure 6: Marketing dataset overhead in generic MPC and HOLMES.

Table 2: Cost breakdown for the marketing dataset (two-party).

Number of entries (41188×2)	82376
Number of attributes	21
Total time	29.19 s
Average time per entry	0.35 ms
Loading the data to IZK	0.27 s
Range tests for all attributes	5.40 s
Histogram and χ^2 test on <i>age</i>	0.76 s
Multidimensional χ^2 test on <i>age, job, marital status, education</i>	22.3 s
Mean, variance, and <i>t</i> test on <i>call duration</i>	0.12 s
Consistency check	< 0.01 s

4.4.3 Evaluation on real-world dataset

We apply HOLMES and the generic MPC baseline to a real-world bank marketing dataset and study the overhead. We include two additional real-world examples of dataset testing workflows and study their overheads in the full version of the paper [51].

In our experiment, we vary the number of parties from 2 to 10 to see how HOLMES and the generic MPC baseline scale with more parties. In secure collaborative learning with more parties we expect to have access to more data, so we assume that each party provides the same amount of data N_{ind} ; when there are t parties, there are $t \cdot N_{\text{ind}}$ data. For example, for our marketing data, we assume that each party provides $N_{\text{ind}} = 41188$ entries of data. When there are 10 parties, the entire distribution testing would be over $N_{\text{ind}} \cdot t = 411880$

entries.

Marketing dataset workflow. The dataset [43, 44] consists of telemarketing records for financial products. It includes client profile and call records. We choose a distribution test workflow that fits the common use case of untrusted banks who wish to jointly train a model to predict the success of the campaign. Before training they want to ensure that the dataset has a balanced number of customers from different backgrounds. Therefore, banks may consider the following tests:

- Pearson’s (histogram) χ^2 -test over age grouped into the buckets 10–19, 20–29, \dots , 90–99 to ensure that the dataset distribution is similar to the national census age distribution,
- Pearson’s (multidimensional) χ^2 -test over age, job, educational level, and marital status to check if the dataset is balanced across customers with different backgrounds,
- t -test over the telemarketing call duration to check whether their telemarketing records are similar enough to train a model together, and
- range checks for all attributes.

Results. We present our results in Fig. 6(a) and Fig. 6(b). HOLMES’s approach outperforms the generic MPC baseline by 77–264 \times . The overhead gap widens as the number of parties increases. As expected, HOLMES’s overhead grows linearly to t , while the baseline overhead grows quadratically to t . We also present the cost breakdown in Tab. 2. We see that the range and multidimensional tests contribute to a large portion of the overhead compared to all other tests. The consistency check between IZK and MPC has a small overhead.

5 Related Work

We summarize related works and explain their connection to HOLMES.

Secure multiparty computation frameworks. A rich body of works propose MPC protocols [49, 61, 62] for malicious adversaries and dishonest majority, with SPDZ [63–65] and authenticated garbling [66–69] being the state-of-the-art. HOLMES uses SPDZ since it is more suitable for arithmetic computation that is used for secure collaborative learning.

Zero-knowledge proofs. Zero-knowledge proofs [70] enable a party to prove a statement without leaking any information. Constructing practical ZK has gained much attention, especially since succinct non-interactive proofs [25, 71–74] have been used in blockchains. New protocols for interactive zero-knowledge proofs based on silent OT [24, 75–80] are currently being studied for their efficiency. Although not currently ready for implementation, a subset of these protocols known as line-point zero-knowledge (LPZK) [81, 82] promise greater flexibility with primes and smaller prover and verification overhead than QuickSilver. HOLMES in the future can be extended to these newer protocols.

Statistics and range checks. There have been works [83, 84] whose goal is to perform statistical tests privately using MPC. In contrast to our protocol, they mostly focus on the two-party case and consider different threat models. Also, range checks [72, 85, 86] are frequently used to limit the effect of misreported values in secure computation. As an example, Prio [85] (or Prio+ [86]) is a system that aggregates statistics over multiple clients who wish to preserve the confidentiality of their individual data but relies on the existence of non-colluding semi-honest servers. HOLMES offers security guarantees even with a dishonest majority.

Secure collaborative computation systems. Multiple works build systems for data analytics and machine learning against malicious adversaries [18, 19, 87–100], but they do not address the issue of corrupted input datasets or group fairness, which is often left as an open question. We envision an integration of HOLMES to secure collaborative computation systems as an efficient method for distribution testing.

6 Conclusion

We first discuss some challenges that HOLMES does not solve and we identify several exciting directions for future work. Finally, we conclude with a summary of our contributions.

Identifying necessary tests. HOLMES enables parties to perform distribution tests tailored to their use case. It does not, however, decide what the necessary tests are. The parties have to specify tests depending on their application. We are not aware of a systematic approach that identifies the necessary tests for a specific application, such as measuring the data quality or identifying bias in clinical trials. Since this is relevant even without any privacy considerations, this question is orthogonal to the goal of HOLMES: privately computing distribution tests. A compelling future direction is to combine HOLMES’s rich class of distribution tests with a systematic approach to identify necessary tests in practical applications.

Privacy leakage from distribution tests. Any distribution test leaks one-bit information – whether the test passed or failed – which leads to potential attacks. For instance, assume that an organization wants to check that the mean value of another organization’s dataset is a close to a specific value, e.g., in a medical study we might want to prove that the mean efficacy of a drug is 0.9. A malicious organization, who is not supposed to know this mean can recover it by requesting multiple distributions tests. For example, in the medical study example, the adversary can ask whether the mean is 0.01, 0.02, etc., until the distribution test succeeds.

Potential mitigations for this problem include having a curated list of allowed distribution tests (e.g., proving that the mean of the age of a population is below 150 does not leak any sensitive information), or enforce a rate limit on the tests. An interesting direction for future research is to devise attacks

Table 3: Overhead of the ZK-friendly sketching with varying parameters. For simplicity, we assume that we have the same number of buckets for each attribute. The dimension setup refers to [number of attributes, buckets of each attribute]. Spartan_{SNARK} is significantly slower than Spartan_{NIZK}, so we omit benchmarks for input sizes 200k and 500k.

N_{ind}	Dimension setup	Number of parties = 2			Number of parties = 6			Number of parties = 10		
		[1, 10]	[4, 10]	[4, 50]	[1, 10]	[4, 10]	[4, 50]	[1, 10]	[4, 10]	[4, 50]
100k	QuickSilver	57.3s	57.9s	58.5s	286.8s	289.9s	292.6s	516.3s	521.9s	526.7s
	Paired 2PC	2064s	2074s	2074s	4340s	4385s	4354s	6927s	6955s	6966s
	MPC	2064s	2074s	2074s	36333s	36330s	36522s	91387s	90982s	90648s
	Spartan _{NIZK}	2602s	2610s	2613s	2521s	2559s	2579s	2667s	2651s	2689s
	Spartan _{SNARK}	20752s	20985s	20212s	20754s	20988s	20214s	20757s	20990s	20217s
200k	QuickSilver	114.1s	111.8s	111.9s	570.6s	558.8s	559.6s	1027s	1006s	1007s
	Paired 2PC	4086s	4098s	4106s	8734s	8737s	8708s	13915s	13788s	13935s
	MPC	4086s	4098s	4106s	72600s	72934s	72863s	182917s	184209s	183890s
	Spartan _{NIZK}	5065s	5087s	5090s	5152s	5146s	5127s	5135s	5140s	5153s
	Spartan _{SNARK}	280.5s	277.8s	276.9s	1402s	1389s	1385s	2524s	2500s	2492s
500k	QuickSilver	10136s	10163s	10201s	21914s	21794s	21773s	34820s	34941s	34840s
	Paired 2PC	10136s	10163s	10201s	181400s	182747s	181889s	456853s	456366s	463618s
	MPC	10002s	9989s	10152s	10501s	10388s	10441s	10549s	10417s	10521s
	Spartan _{NIZK}									
	Spartan _{SNARK}									

that exploit this leakage and identify mitigations in specific applications.

Improving HOLMES’s efficiency. Even though HOLMES outperforms the baselines based on state-of-the-art systems in our benchmarks, there are specific cases that other systems have a small efficiency advantage. For example, in the range check test with $\ell = 8$ and $N_{\text{ind}} = 100k$, at 11 parties or more we expect Spartan_{NIZK} (31.5s) to be faster than QuickSilver (34s) since Spartan’s overhead remains relatively constant to a growing number of parties. We leave as an open question how to build a system that is more efficient in all settings.

New protocols for interactive zero-knowledge proofs based on silent OT [24, 75–80] are currently being studied for their efficiency. Although not ready for implementation, a subset of these protocols known as line-point zero-knowledge (LPZK) [81, 82] promises greater flexibility with primes and smaller prover overhead than QuickSilver. HOLMES in the future can be extended to these newer protocols.

Additionally, HOLMES supports distribution tests performed on a random subset of the dataset to boost efficiency. When the datasets are sufficiently large, intuitively subsampling should not affect accuracy. However, we are not aware of specific applications where this feature can be tested.

Adversarial machine learning. HOLMES is a useful tool for identifying bias in datasets used in machine learning training without compromising their privacy. Even though we have experimented with the accuracy of HOLMES in specific adversarial scenarios, our protocol does not offer any formal guarantees. In the realm of adversarial machine learning, data poisoning shows that it is possible to corrupt a machine learning model by using datasets practically indistinguishable from the honest ones. As a mitigation, robust statistics [101–103] focuses on statistics that are resilient to any corrupted input distribution. A fascinating future direction is to augment HOLMES with robust statistics that not only detect bias, but can reduce its effect in the final machine learning application.

In conclusion, we present HOLMES, a protocol for performing distribution testing in secure collaborative learning efficiently. The core of HOLMES consists of two contributions:

- a new hybrid protocol that integrates MPC, IZK, and a lightweight consistency check for distribution testing, which is concretely more efficient than non-trivial baselines, and
- a novel, efficient multidimensional distribution testing procedure that utilizes sketching and pseudorandom functions to avoid the severe penalty of oblivious computation.

These two tools significantly improve the performance of distribution testing. Efficient support for distribution testing can be seen as the first step towards detecting different types of incorrect (or even malicious) inputs for secure computation in general, which is an essential for practical secure collaborative learning. HOLMES is open-sourced in GitHub: <https://github.com/holmes-inputcheck/>.

Acknowledgements

We thank the anonymous reviewers and our shepherd for their helpful feedback. This work is supported by NSF CISE Expeditions Award CCF-1730628, NSF CAREER 1943347, and gifts from the Alibaba, Amazon Web Services, Ant Group, Astronomer, Ericsson, Facebook, Futurewei, Google, IBM, Intel, Lacework, Microsoft, Nexla, Nvidia, Samsung, Scotiabank, Splunk, and VMware.

References

- [1] Karen Hao. *AI is sending people to jail—and getting it wrong*. <https://www.technologyreview.com/2019/01/21/137783/algorithms-criminal-justice-ai/>.

- [2] Ganes Kesari. *AI Can Now Detect Depression From Your Voice, And It's Twice As Accurate As Human Practitioners*. <https://bit.ly/32HdcUQ>.
- [3] *Ellipsis Health*. <https://www.ellipsishealth.com/>.
- [4] *Qbtech*. <https://www.qbtech.com/>.
- [5] McKinsey Global Institute. *Tackling bias in artificial intelligence (and in humans)*. <https://mck.co/3Ge65B8>.
- [6] Zhe Yu, Joymallya Chakraborty, and Tim Menzies. "FairBalance: Improving Machine Learning Fairness on Multiple Sensitive Attributes With Data Balancing". In: *Arxiv:2107.08310*.
- [7] Angelina Wang, Arvind Narayanan, and Olga Russakovsky. "REVISE: A tool for measuring and mitigating bias in visual datasets". In: *ECCV '20*.
- [8] Faisal Kamiran and Toon Calders. "Data preprocessing techniques for classification without discrimination". In: *Knowledge and Information Systems* 33.1 (2012), pp. 1–33.
- [9] Thomas Davidson, Debasmita Bhattacharya, and Ingmar Weber. "Racial Bias in Hate Speech and Abusive Language Detection Datasets". In: *Workshop on Abusive Language Online '19*.
- [10] Nature Communications Editorial. "Data sharing and the future of science". In: *Nature Communications '18*.
- [11] Heather A. Piwowar and Todd J. Vision. "Data reuse and the open data citation advantage". In: *PeerJ '13*.
- [12] Milton Packer. "Data sharing in medical research". In: *British Medical Journal '18*.
- [13] Liina Kamm, Dan Bogdanov, Sven Laur, and Jaak Vilo. "A new way to protect privacy in large-scale genome-wide association studies". In: *Bioinformatics '13*.
- [14] Emmanuel A Abbe, Amir E Khandani, and Andrew W Lo. "Privacy-preserving methods for sharing financial risk exposures". In: *American Economic Review '12*.
- [15] *General Data Protection Regulation*. <https://gdpr-info.eu/>.
- [16] *Rights related to automated decision making including profiling*. <https://bit.ly/3ALUJ6m>.
- [17] *Assessing Data Quality for Healthcare Systems Data Used in Clinical Research*. https://dcricollab.dcri.duke.edu/sites/NIHKKR/KR/Assessing-data-quality_V1%200.pdf.
- [18] Rishabh Poddar, Sukrit Kalra, Avishay Yanai, Ryan Deng, Raluca Ada Popa, and Joseph M. Hellerstein. "Senate: A Maliciously-Secure MPC Platform for Collaborative Analytics". In: *SEC '20*.
- [19] Wenting Zheng, Raluca Ada Popa, Joseph E. Gonzalez, and Ion Stoica. "Helen: Maliciously Secure Cooperative Learning for Linear Models". In: *S&P '19*.
- [20] Sameer Wagh, Shruti Tople, Fabrice Benhamouda, Eyal Kushilevitz, Prateek Mittal, and Tal Rabin. "Falcon: Honest-Majority Maliciously Secure Framework for Private Deep Learning". In: *PETS '21*.
- [21] Sameer Wagh, Divya Gupta, and Nishanth Chandran. "SecureNN: 3-Party Secure Computation for Neural Network Training". In: *PETS '19*.
- [22] *Multi-Protocol SPDZ*. <https://github.com/data61/MP-SPDZ>.
- [23] *SCALE and MAMBA*. <https://github.com/KULeuven-COSIC/SCALE-MAMBA>.
- [24] Kang Yang, Pratik Sarkar, Chenkai Weng, and Xiao Wang. "QuickSilver: Efficient and Affordable Zero-Knowledge Proofs for Circuits and Polynomials over Any Field". In: *CCS '21*.
- [25] Srinath Setty. "Spartan: Efficient and General-Purpose zkSNARKs Without Trusted Setup". In: *CRYPTO '20*.
- [26] Oded Goldreich. "Towards a theory of software protection". In: *CRYPTO '86*.
- [27] Oded Goldreich. "Towards a theory of software protection and simulation by oblivious RAMs". In: *STOC '87*.
- [28] Oded Goldreich and Rafail Ostrovsky. "Software protection and simulation on oblivious RAMs". In: *JACM '96*.
- [29] Benny Pinkas and Tzachy Reinman. "Oblivious RAM revisited". In: *CRYPTO '10*.
- [30] Elaine Shi, T-H Hubert Chan, Emil Stefanov, and Mingfei Li. "Oblivious RAM with $O(\log N^3)$ worst-case cost". In: *ASIACRYPT '11*.
- [31] Emil Stefanov, Elaine Shi, and Dawn Song. "Towards practical oblivious RAM". In: *NDSS '12*.
- [32] Emil Stefanov, Marten Van Dijk, Elaine Shi, Christopher Fletcher, Ling Ren, Xiangyao Yu, and Srinivas Devadas. "Path ORAM: An extremely simple oblivious RAM protocol". In: *CCS '13*.
- [33] Craig Gentry, Kenny A Goldman, Shai Halevi, Charanjit Julta, Mariana Raykova, and Daniel Wichs. "Optimizing ORAM and using it efficiently for secure computation". In: *PETS '13*.
- [34] William B Johnson and Joram Lindenstrauss. "Extensions of Lipschitz mappings into a Hilbert space". In: *Contemporary mathematics '84*.

- [35] Lorenzo Grassi, Dmitry Khovratovich, Christian Rechberger, Arnab Roy, and Markus Schafneggger. “POSEIDON: A New Hash Function for Zero-Knowledge Proof System”. In: *SEC '21*.
- [36] Abdelrahman Aly, Tomer Ashur, Eli Ben-Sasson, Siemen Dhooghe, and Alan Szepieniec. “Design of Symmetric-Key Primitives for Advanced Cryptographic Protocols”. In: *FSE '20*.
- [37] Ivan Bjerre Damgård. “On the randomness of Legendre and Jacobi sequences”. In: *CRYPTO '88*.
- [38] Lorenzo Grassi, Christian Rechberger, Dragos Rotaru, Peter Scholl, and Nigel P. Smart. “MPC-friendly symmetric key primitives”. In: *CCS '16*.
- [39] Dmitry Khovratovich. “Key recovery attacks on the Legendre PRFs within the birthday bound”. In: *IACR ePrint 2019/862*.
- [40] Alexander May and Floyd Zveydinger. “Legendre PRF (Multiple) Key Attacks and the Power of Preprocessing”. In: *IACR ePrint 2021/645*.
- [41] Mark N. Wegman and J. Lawrence Carter. “New hash functions and their use in authentication and set equality”. In: *JCSS '81*.
- [42] Moni Naor and Moti Yung. “Universal one-way hash functions and their cryptographic applications”. In: *STOC '89*.
- [43] Sérgio Moro, Paulo Cortez, and Paulo Rita. “A data-driven approach to predict the success of bank telemarketing”. In: *Decision Support Systems '14*.
- [44] *Bank Marketing Data Set*. <https://archive.ics.uci.edu/ml/datasets/Bank+Marketing>.
- [45] Beata Strack, Jonathan P. DeShazo, Chris Gennings, Juan L. Olmo, Sebastian Ventura, Krzysztof J. Cios, and John N. Clore. “Impact of HbA1c measurement on hospital readmission rates: Analysis of 70,000 clinical database patient records”. In: *BioMed Research International '14*.
- [46] *Diabetes 130-US hospitals for years 1999-2008 Data Set*. <https://archive.ics.uci.edu/ml/datasets/Diabetes+130-US+hospitals+for+years+1999-2008>.
- [47] *Real Time Advertiser's Auction*. <https://www.kaggle.com/saurav9786/real-time-advertisers-auction>.
- [48] Ran Canetti, Yehuda Lindell, Rafail Ostrovsky, and Amit Sahai. “Universally composable two-party and multi-party secure computation”. In: *STOC '02*.
- [49] Oded Goldreich, Silvio Micali, and Avi Wigderson. “How to Play any Mental Game or A Completeness Theorem for Protocols with Honest Majority”. In: *STOC '87*.
- [50] Yehuda Lindell. “How to simulate it: A tutorial on the simulation proof technique”. In: *Tutorials on the Foundations of Cryptography: Dedicated to Oded Goldreich*. 2017.
- [51] Ian Chang, Katerina Sotiraki, Weikeng Chen, Murat Kantarcioglu, and Raluca Ada Popa. *HOLMES: Efficient Distribution Testing for Secure Collaborative Learning*. Cryptology ePrint Archive, Paper 2021/1517. <https://eprint.iacr.org/2021/1517>. 2021.
- [52] Manuel Blum. “Coin flipping by telephone a protocol for solving impossible problems”. In: *ACM SIGACT News '83*.
- [53] Carl Friedrich Gauss. “Theoria combinationis observationum erroribus minimis obnoxiae”. In: *Carl Friedrich Gauss Werke*. 1823.
- [54] Ery Arias-Castro, Bruno Pelletier, and Venkatesh Saligrama. “Remember the curse of dimensionality: the case of goodness-of-fit testing in arbitrary dimension”. In: *Journal of Nonparametric Statistics '18*.
- [55] Robert B Davies. “Algorithm AS 155: The distribution of a linear combination of χ^2 random variables”. In: *Applied Statistics* (1980), pp. 323–333.
- [56] J Sheil and I O’Muircheartaigh. “Algorithm AS 106: The distribution of non-negative quadratic forms in normal variables”. In: *Journal of the Royal Statistical Society. Series C (Applied Statistics)* 26.1 (1977), pp. 92–98.
- [57] Jean-Pierre Imhof. “Computing the distribution of quadratic forms in normal variables”. In: *Biometrika* 48.3/4 (1961), pp. 419–426.
- [58] Dimitris Achlioptas. “Database-Friendly Random Projections: Johnson-Lindenstrauss with Binary Coins”. In: *Journal of Computer and System Sciences '03*.
- [59] Geoffroy Couteau, Peter Rindal, and Srinivasan Raghuraman. “Silver: Silent VOLE and Oblivious Transfer from Hardness of Decoding Structured LDPC Codes”. In: *CRYPTO '21*.
- [60] Marcel Keller. “MP-SPDZ: A Versatile Framework for Multi-Party Computation”. In: *CCS '20*.
- [61] Michael Ben-Or, Shafi Goldwasser, and Avi Wigderson. “Completeness Theorems for Non-Cryptographic Fault-Tolerant Distributed Computation”. In: *STOC '88*.
- [62] Andrew Chi-Chih Yao. “Protocols for Secure Computations”. In: *FOCS '82*.
- [63] Ivan Damgård, Valerio Pastro, Nigel P. Smart, and Sarah Zakarias. “Multiparty Computation from Somewhat Homomorphic Encryption”. In: *CRYPTO '12*.

- [64] Ivan Damgård, Marcel Keller, Enrique Larraia, Valerio Pastro, Peter Scholl, and Nigel P. Smart. “Practical Covertly Secure MPC for Dishonest Majority - Or: Breaking the SPDZ Limits”. In: *ESORICS '13*.
- [65] Marcel Keller, Valerio Pastro, and Dragos Rotaru. “Overdrive: Making SPDZ Great Again”. In: *EUROCRYPT '18*.
- [66] Xiao Wang, Samuel Ranellucci, and Jonathan Katz. “Global-Scale Secure Multiparty Computation”. In: *CCS '17*.
- [67] Carmit Hazay, Peter Scholl, and Eduardo Soria-Vazquez. “Low Cost Constant Round MPC Combining BMR and Oblivious Transfer”. In: *ASIACRYPT '17*.
- [68] Kang Yang, Xiao Wang, and Jiang Zhang. “More Efficient MPC from Improved Triple Generation and Authenticated Garbling”. In: *CCS '20*.
- [69] Kang Yang, Chenkai Weng, Xiao Lan, Jiang Zhang, and Xiao Wang. “Ferret: Fast Extension for Correlated OT with small communication”. In: *CCS '20*.
- [70] Shafi Goldwasser, Silvio Micali, and Charles Rackoff. “The Knowledge Complexity of Interactive Proof-Systems”. In: *STOC '85*.
- [71] Jonathan Bootle, Andrea Cerulli, Pyrros Chaidos, Jens Groth, and Christophe Petit. “Efficient Zero-Knowledge Arguments for Arithmetic Circuits in the Discrete Log Setting”. In: *EUROCRYPT '16*.
- [72] Benedikt Bünz, Jonathan Bootle, Dan Boneh, Andrew Poelstra, Pieter Wuille, and Gregory Maxwell. “Bulletproofs: Short Proofs for Confidential Transactions and More”. In: *S&P '18*.
- [73] Tiancheng Xie, Jiaheng Zhang, Yupeng Zhang, Charalampos Papamanthou, and Dawn Song. “Libra: Succinct Zero-Knowledge Proofs with Optimal Prover Computation”. In: *CRYPTO '19*.
- [74] Jiaheng Zhang, Tiancheng Xie, Yupeng Zhang, and Dawn Song. “Transparent Polynomial Delegation and Its Applications to Zero Knowledge Proof”. In: *S&P '20*.
- [75] Chenkai Weng, Kang Yang, Xiang Xie, Jonathan Katz, and Xiao Wang. “Mystique: Efficient Conversions for Zero-Knowledge Proofs with Applications to Machine Learning”. In: *SEC '21*.
- [76] Chenkai Weng, Kang Yang, Xiao Wang, and Jonathan Katz. “Wolverine: Fast, Scalable, and Communication-Efficient Zero-Knowledge Proofs for Boolean and Arithmetic Circuits”. In: *S&P '21*.
- [77] Samuel Dittmer, Yuval Ishai, and Rafail Ostrovsky. “Line-Point Zero Knowledge and Its Applications”. In: *ITC '21*.
- [78] Carsten Baum, Alex J. Malozemoff, Marc B. Rosen, and Peter Scholl. “Mac’n’Cheese: Zero-Knowledge Proofs for Boolean and Arithmetic Circuits with Nested Disjunctions”. In: *CRYPTO '21*.
- [79] Elette Boyle, Geoffroy Couteau, Niv Gilboa, and Yuval Ishai. “Compressing Vector OLE”. In: *CCS '18*.
- [80] Elette Boyle, Geoffroy Couteau, Niv Gilboa, Yuval Ishai, Lisa Kohl, and Peter Scholl. “Efficient Pseudorandom Correlation Generators: Silent OT Extension and More”. In: *CRYPTO '19*.
- [81] Samuel Dittmer, Yuval Ishai, and Rafail Ostrovsky. “Line-Point Zero Knowledge and Its Applications”. In: *ITC '21*.
- [82] Samuel Dittmer, Yuval Ishai, Steve Lu, and Rafail Ostrovsky. *Improving Line-Point Zero Knowledge: Two Multiplications for the Price of One*. Cryptology ePrint Archive, Paper 2022/552. <https://eprint.iacr.org/2022/552>. 2022. URL: <https://eprint.iacr.org/2022/552>.
- [83] Alexandr Andoni, Tal Malkin, and Negev Shekel Nosatzki. “Two party distribution testing: Communication and security”. In: *arXiv preprint arXiv:1811.04065* (2018).
- [84] Varun Narayanan, Manoj Mishra, and Vinod M Prabhakaran. “Private two-terminal hypothesis testing”. In: *ISIT '20*.
- [85] Henry Corrigan-Gibbs and Dan Boneh. “Prio: Private, robust, and scalable computation of aggregate statistics”. In: *NSDI '17*.
- [86] Surya Addanki, Kevin Garbe, Eli Jaffe, Rafail Ostrovsky, and Antigoni Polychroniadou. “Prio+: Privacy Preserving Aggregate Statistics via Boolean Shares”. In: *IACR ePrint 2021/576*.
- [87] Payman Mohassel and Yupeng Zhang. “SecureML: A system for scalable privacy-preserving machine learning”. In: *S&P '17*.
- [88] Melissa Chase, Ran Gilad-Bachrach, Kim Laine, Kristin E. Lauter, and Peter Rindal. “Private Collaborative Neural Network Learning”. In: *IACR ePrint 2017/762*.
- [89] Irene Giacomelli, Somesh Jha, Marc Joye, C David Page, and Kyonghwan Yoon. “Privacy-preserving ridge regression with only linearly-homomorphic encryption”. In: *ACNS '18*.
- [90] Jian Liu, Mika Juuti, Yao Lu, and Nadarajah Asokan. “Oblivious neural network predictions via MiniONN transformations”. In: *CCS '17*.
- [91] Chiraag Juvekar, Vinod Vaikuntanathan, and Anantha Chandrakasan. “GAZELLE: A low latency framework for secure neural network inference”. In: *SEC '18*.

- [92] M Sadegh Riazi, Mohammad Samragh, Hao Chen, Kim Laine, Kristin Lauter, and Farinaz Koushanfar. “XONN: XNOR-based oblivious deep neural network inference”. In: *SEC '19*.
- [93] Valerie Chen, Valerio Pastro, and Mariana Raykova. “Secure Computation for Machine Learning With SPDZ”. In: *NeurIPS '18*.
- [94] Anselme Tueno, Florian Kerschbaum, and Stefan Katzenbeisser. “Private Evaluation of Decision Trees using Sublinear Cost”. In: *PETS '19*.
- [95] Valeria Nikolaenko, Udi Weinsberg, Stratis Ioannidis, Marc Joye, Dan Boneh, and Nina Taft. “Privacy-preserving ridge regression on hundreds of millions of records”. In: *S&P '13*.
- [96] Adrià Gascón, Phillipp Schoppmann, Borja Balle, Mariana Raykova, Jack Doerner, Samee Zahur, and David Evans. “Privacy-Preserving Distributed Linear Regression on High-Dimensional Data.” In: *PETS '17*.
- [97] Hao Chen, Miran Kim, Ilya Razenshteyn, Dragos Rotaru, Yongsoo Song, and Sameer Wagh. “Maliciously secure matrix multiplication with applications to private deep learning”. In: *ASIACRYPT '20*.
- [98] Wenting Zheng, Ryan Deng, Weikeng Chen, Raluca Ada Popa, Aurojit Panda, and Ion Stoica. “Cerebro: A Platform for Multi-Party Cryptographic Collaborative Learning”. In: *SEC '21*.
- [99] Christopher A. Choquette-Choo, Natalie Dullerud, Adam Dziedzic, Yunxiang Zhang, Somesh Jha, Nicolas Papernot, and Xiao Wang. “CaPC Learning: Confidential and Private Collaborative Learning”. In: *Arxiv:2102.05188*. 2021.
- [100] Mark Abspoel, Daniel Escudero, and Nikolaj Volgushev. “Secure training of decision trees with continuous attributes”. In: *PETS '21*.
- [101] Frank R Hampel, Elvezio M Ronchetti, Peter J Rousseeuw, and Werner A Stahel. *Robust statistics: The approach based on influence functions*. Vol. 196. John Wiley & Sons, 2011.
- [102] Peter J Huber. *Robust statistics*. Vol. 523. John Wiley & Sons, 2004.
- [103] Ricardo A Maronna, R Douglas Martin, Victor J Yohai, and Matías Salibián-Barrera. *Robust statistics: Theory and methods (with R)*. John Wiley & Sons, 2019.