

Taking proof-based verified computation a few steps closer to practicality

Srinath Setty, Victor Vu, Nikhil Panpalia, Benjamin Braun, Andrew J. Blumberg, and Michael Walfish
The University of Texas at Austin

Abstract. We describe GINGER, a built system for unconditional, general-purpose, and nearly practical verification of outsourced computation. GINGER is based on PEPPER, which uses the PCP theorem and cryptographic techniques to implement an *efficient argument* system (a kind of interactive protocol). GINGER slashes the query size and costs via theoretical refinements that are of independent interest; broadens the computational model to include (primitive) floating-point fractions, inequality comparisons, logical operations, and conditional control flow; and includes a parallel GPU-based implementation that dramatically reduces latency.

1 Introduction

We are motivated by *outsourced computing*: cloud computing (in which clients outsource computations to remote computers), peer-to-peer computing (in which peers outsource storage and computation to each other), volunteer computing (in which projects outsource computations to volunteers’ desktops), etc.

Our goal is to build a system that lets a client outsource computation verifiably. The client should be able to send a description of a computation and the input to a server, and receive back the result together with some auxiliary information that lets the client *verify* that the result is correct. For this to be sensible, the verification must be faster than executing the computation locally.

Ideally, we would like such a system to be *unconditional*, *general-purpose*, and *practical*. That is, we don’t want to make assumptions about the server (trusted hardware, independent failures of replicas, etc.), we want a setup that works for a broad range of computations, and we want the system to be usable by real people for real computations in the near future.

In principle, the first two properties above have been achievable for almost thirty years, using powerful tools from complexity theory and cryptography. Interactive proofs (IPs) and probabilistically checkable proofs (PCPs) show how one entity (usually called the *verifier*) can be convinced by another (usually called the *prover*) of a given mathematical assertion—without the verifier having to fully inspect a proof [5, 6, 19, 32]. In our context, the mathematical assertion is that a given computation was carried out correctly; though the proof is as long as the computation, the theory implies—surprisingly—that the verifier need only inspect the proof in a small number of randomly-chosen locations or query the prover a relatively small number of times.

The rub has been the third property: practicality. These protocols have required expensive encoding of compu-

tations, monstrously large proofs, high error bounds, prohibitive overhead for the prover, and intricate constructions that make the asymptotically efficient schemes challenging to implement correctly.

However, a line of recent work indicates that approaches based on IPs and PCPs are closer to practicality than previously thought [21, 44, 45, 49]. More generally, there has been a groundswell of work that aims for potentially practical verifiable outsourced computation, using theoretical tools [11, 12, 20, 24, 25].

Nonetheless, these works have notable limitations. Only a handful [21, 44, 45, 49] have produced working implementations, all of which impose high costs on the verifier and prover. Moreover, their model of computation is *arithmetic circuits* over finite fields, which represent non-integers awkwardly, control flow inefficiently, and comparisons and logical operations only by degenerating to verbose *Boolean* circuits. Arithmetic circuits are well-suited to integer computations and numerical straight line computations (e.g., multiplying matrices and computing second moments), but the intersection of these two domains leaves few realistic applications.

This paper describes a built system, called GINGER, that addresses these problems, thereby taking general-purpose proof-based verified computation several steps closer to practicality. GINGER is an *efficient argument* system [37, 38]: an interactive proof system that assumes the prover to be computationally bounded. Its starting point is the PEPPER protocol [45] (which is summarized in Section 2). GINGER’s contributions are as follows.

(1) GINGER *demonstrates the strength of linear commitment* (§3). This paper proves that PEPPER’s commitment primitive [45], which generalizes the commitment primitive of Ishai et al. [35], is surprisingly powerful: it not only commits an untrusted entity to a function and extracts evaluations of that function (as previously shown) but also ensures that the function is linear. (The primitive embeds a *strong linearity test*.) This result sharply reduces the required number of queries (from 500 to 3) and a key error bound, and hence overhead.

(2) GINGER *supports a general-purpose programming model* (§4). Although the model does not handle looping concisely, it includes primitive floating-point quantities, inequality comparisons, logical expressions, and conditional control flow. Moreover, we have a compiler (derived from Fairplay [39]) that transforms computations expressed in a general-purpose language to an executable verifier and prover. The core technical challenge is representing computations as additions and multiplications over a finite field (as required by the verification proto-

col); for instance, “not equal” and “if/else” do not obviously map to this formalism, inequalities are problematic because finite fields are not ordered, and fractions compound the difficulties. GINGER overcomes these challenges with techniques that, while not deep, require care and detail.¹ These techniques should apply to other protocols that use arithmetic constraints or circuits.

(3) GINGER exploits parallelism to slash latency (§5). The prover can be distributed across machines, and some of its functions are implemented in graphics hardware (GPUs). Moreover, GINGER’s verifier can use a GPU for its cryptographic operations. Allowing the verifier to have a GPU models the present (many computers have GPUs) and a plausible future in which specialized hardware for cryptographic operations is common.²

We have implemented and evaluated GINGER (§6). Compared to PEPPER [45], its base, GINGER lowers network costs by 1–2 orders of magnitude (to hundreds of KB or less in our experiments). The verifier’s costs drop by multiples and possibly orders of magnitude, depending on the cost of encryption; if we model encryption as free, the verifier can gain from outsourcing when batch-verifying as few as 20 computations (down from 3900 in PEPPER). The prover’s CPU costs drop by 10–15%, which is not much, but our parallel implementation reduces latency with near-linear speedup. Computing with rational numbers in GINGER is roughly three times more expensive than computing with integers, and arithmetic constraints permit far smaller representations than a naive use of Boolean or arithmetic circuits.

Despite all of the above, GINGER is not quite ready for the big leagues. However, PEPPER and GINGER have made argument systems far more practical (in some cases improving costs by 23 orders of magnitude over a naive implementation). We are thus optimistic about ultimately achieving true practicality.

2 Problem statement and background

Problem statement. A computer V , known as the *verifier*, has a computation Ψ and some desired input x that it wants a computer P , known as the *prover*, to perform. P returns y , the purported output of the computation, and then V and P conduct an efficient interaction. This interaction should be cheaper for V than locally computing $\Psi(x)$. Furthermore, if P returned the correct answer, it should be able to convince V of that fact; otherwise, V should be able to reject the answer as incorrect, with high probability. (The converse will not hold: rejection does not imply that P returned incorrect output, only that

it misbehaved somehow.) Our goal is that this guarantee be *unconditional*: it should hold regardless of whether P obeys the protocol (given standard cryptographic assumptions about P ’s computational power). If P deviates from the protocol at any point (computing incorrectly, proving incorrectly, etc.), we call it *malicious*.

2.1 Tools

In principle, we can meet our goal using PCPs. The PCP theorem [5, 6] says that if a set of constraints is satisfiable (see below), there exists a *probabilistically checkable* proof (a PCP) and a verification procedure that accepts the proof after querying it in only a small number of locations. On the other hand, if the constraints cannot be satisfied, then the verification procedure rejects *any* purported proof, with probability at least $1 - \epsilon$.

To apply the theorem, we represent the computation as a set of quadratic constraints over a finite field. A *quadratic constraint* is an equation of degree 2 that uses additions and multiplications (e.g., $A \cdot Z_1 + Z_2 - Z_3 \cdot Z_4 = 0$). A set of constraints is *satisfiable* if the variables can be set to make all of the equations hold simultaneously; such an assignment is called a *satisfying assignment*. In our context, a set of constraints \mathcal{C} will have a designated input variable X and output variable Y (this generalizes to multiple inputs and outputs), and $\mathcal{C}(X = x, Y = y)$ denotes \mathcal{C} with variable X bound to x and Y bound to y .

We say that a set of constraints \mathcal{C} is *equivalent* to a desired computation Ψ if: for all x, y , $\mathcal{C}(X = x, Y = y)$ is satisfiable if and only if $y = \Psi(x)$. As a simple example, increment-by-1 is equivalent to the constraint set $\{Y = Z + 1, Z = X\}$. (For convenience, we will sometimes refer to a given input x and purported output y implicitly in statements such as, “If constraints \mathcal{C} are satisfiable, then Ψ executed correctly”.) To verify a computation $y = \Psi(x)$, one could in principle apply the PCP theorem to the constraints $\mathcal{C}(X = x, Y = y)$.

Unfortunately, PCPs are too large to be transferred. However, if we assume a computational bound on the prover P , then *efficient arguments* apply [37, 38]: V issues its PCP queries to P (so V need not receive the entire PCP). For this to work, P must commit to the PCP *before* seeing V ’s queries, thereby simulating a fixed proof whose contents are independent of the queries. V thus extracts a cryptographic commitment to the PCP (e.g., with a collision-resistant hash tree [40]) and verifies that P ’s query responses are consistent with the commitment.

This approach can be taken a step further: not even P has to materialize the entire PCP. As Ishai et al. [35] observe, in some PCP constructions, which they call *linear PCPs*, the PCP itself is a linear function: the verifier submits queries to the function, and the function’s outputs serve as the PCP responses. Ishai et al. thus design a *linear commitment primitive* in which P can commit to

¹We elide some of these details for space; they are documented in a longer version of this paper [46].

²One may wonder why, if the verifier has this hardware, it needs to outsource. GPUs are amenable only to certain computations (which include the cryptographic underpinnings of GINGER).

a linear function (the PCP) and V can submit function inputs (the PCP queries) to P , getting back outputs (the PCP responses) as if P itself were a fixed function.

PEPPER [45] refines and implements the outline above. In the rest of the section, we summarize the linear PCPs that PEPPER incorporates, give an overview of PEPPER, and provide formal definitions. Additional details are in Appendix A.1.

2.2 Linear PCPs, applied to verifying computations

Imagine that V has a desired computation Ψ and desired input x , and somehow obtains purported output y . To use PCP machinery to check whether $y = \Psi(x)$, V compiles Ψ into equivalent constraints \mathcal{C} , and then asks whether $\mathcal{C}(X = x, Y = y)$ is satisfiable, by consulting an *oracle* π : a fixed function (that depends on \mathcal{C}, x, y) that V can query. A *correct* oracle π is the proof (or PCP); V should accept a correct oracle and reject an incorrect one.

A correct oracle π has three properties. First, π is a *linear function*, meaning that $\pi(a) + \pi(b) = \pi(a + b)$ for all a, b in the domain of π . A linear function $\pi: \mathbb{F}^n \rightarrow \mathbb{F}$ is determined by a vector w ; i.e., $\pi(a) = \langle a, w \rangle$ for all $a \in \mathbb{F}^n$. Here, \mathbb{F} is a finite field, and $\langle a, b \rangle$ denotes the inner (dot) product of two vectors a and b . The parameter n is the size of w ; in general, n is quadratic in the number of variables in \mathcal{C} [5], but we can sometimes tailor the encoding of w to make n smaller [45].

Second, one set of the entries in w must be a redundant encoding of the other entries. Third, w encodes the actual satisfying assignment to $\mathcal{C}(X = x, Y = y)$.

A surprising aspect of PCPs is that each of these properties can be tested by making a small number of queries to π ; if π is constructed incorrectly, the probability that the tests pass is upper-bounded by $\epsilon > 0$. A key test for us—we return to it in Section 3—is the *linearity test* [16]: V randomly selects q_1 and q_2 from \mathbb{F}^n and checks if $\pi(q_1) + \pi(q_2) = \pi(q_1 + q_2)$. The other two PCP tests are the *quadratic correction test* and the *circuit test*.

The completeness and soundness properties of linear PCPs are defined in Section 2.4. A detailed explanation of why the mechanics above satisfy those properties is outside our scope but can be found in [5, 13, 35, 45].

2.3 Our base: PEPPER

We now walk through the three phases of PEPPER [45], which is depicted in Figure 1. The approach is to compose a *linear PCP* and a *linear commitment primitive* that forces the prover to act like an oracle.

Specify and compute. V transforms its desired computation, Ψ , into a set of equivalent constraints, \mathcal{C} . V sends Ψ (or \mathcal{C}) to P , or P may come with them installed.

To gain from outsourcing, V must amortize the costs of compiling Ψ to \mathcal{C} and generating queries. Thus, V verifies computations in batches [45] (although they need not be

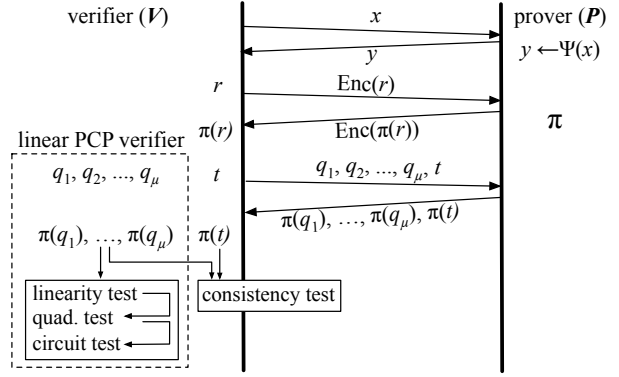


Figure 1—The PEPPER protocol [45], which is GINGER’s base. Though not depicted, many of the protocol steps happen in parallel, to facilitate batching.

executed in a batch). A *batch* (of size β) refers to a set of computations in which Ψ is the same but the inputs are different; a member of the batch is called an *instance*. In the protocol, V has inputs x_1, \dots, x_β that it sends to P (not necessarily all at once), which returns y_1, \dots, y_β ; for each instance i , y_i is supposed to equal $\Psi(x_i)$.

For each instance i , an honest P stores a proof vector w_i that encodes a satisfying assignment to $\mathcal{C}(X = x_i, Y = y_i)$; w_i is constructed as described in Section 2.2. Being a vector, w_i can also be regarded as a linear function π_i —or an oracle of the kind described above.

Extract commitment. V cannot inspect $\{\pi_i\}$ directly (they are functions; written out, they would have an entry for each value in a huge domain). Instead, V extracts a *commitment* to each π_i . To do so, V randomly generates a *commitment vector* $r \in \mathbb{F}^n$. V then homomorphically encrypts each entry of r under a public key pk to get a vector $\text{Enc}(pk, r) = (\text{Enc}(pk, r_1), \text{Enc}(pk, r_2), \dots, \text{Enc}(pk, r_n))$. We emphasize that $\text{Enc}(\cdot)$ need not be fully homomorphic encryption [27] (which remains unfeasibly expensive); PEPPER uses ElGamal [23, 45].

V sends $(\text{Enc}(pk, r), pk)$ to P . If P is honest, then π_i is linear, so P can use the homomorphism of $\text{Enc}(\cdot)$ to compute $\text{Enc}(pk, \pi_i(r))$ from $\text{Enc}(pk, r)$, without learning r . P replies with $(\text{Enc}(pk, \pi_1(r)), \dots, \text{Enc}(pk, \pi_\beta(r)))$, which is P ’s commitment to $\{\pi_i\}$. V then decrypts to get $(\pi_1(r), \dots, \pi_\beta(r))$.

Verify. V now generates PCP queries $q_1, \dots, q_\mu \in \mathbb{F}^n$, as described in Section 2.2. V sends these queries to P , along with a *consistency query* $t = r + \sum_{j=1}^{\mu} \alpha_j \cdot q_j$, where each α_j is randomly chosen from \mathbb{F} (here, \cdot represents scalar multiplication).

For ease of exposition, we focus on a single proof π_i ; however, the following steps happen β times in parallel, using the same queries for each of the β instances. If P is honest, it returns $(\pi_i(q_1), \dots, \pi_i(q_\mu), \pi_i(t))$. V checks that $\pi_i(t) = \pi_i(r) + \sum_{j=1}^{\mu} \alpha_j \cdot \pi_i(q_j)$; this is known as

the *consistency test*. If P is honest, then this test passes, by the linearity of π . Conversely, if this test passes then, *regardless* of P 's honesty, V can treat P 's responses as the output of an oracle (this is shown in previous work [35, 45]). Thus, V can use $\{\pi_i(q_1), \dots, \pi_i(q_\mu)\}$ in the PCP tests described in Section 2.2.

2.4 PCPs and arguments defined more formally

The definitions of PCPs [5, 6] and argument systems [19, 32] below are borrowed from [35, 45].

A *PCP protocol* with soundness error ϵ includes a probabilistic polynomial time verifier V that has access to a constraint set \mathcal{C} . V makes a constant number of queries to an oracle π . This process has the following properties:

- **PCP Completeness.** If \mathcal{C} is satisfiable, then there exists a linear function π such that, after V queries π , $\Pr\{V \text{ accepts } \mathcal{C} \text{ as satisfiable}\} = 1$, where the probability is over V 's random choices.
- **PCP Soundness.** If \mathcal{C} is not satisfiable, then $\Pr\{V \text{ accepts } \mathcal{C} \text{ as satisfiable}\} < \epsilon$ for *all* purported proof functions $\tilde{\pi}$.

An argument (P, V) with soundness error ϵ comprises P and V , two probabilistic polynomial time (PPT) entities that take a set of constraints \mathcal{C} as input and provide:

- **Argument Completeness.** If \mathcal{C} is satisfiable and P has access to a satisfying assignment z , then the interaction of $V(\mathcal{C})$ and $P(\mathcal{C}, z)$ makes $V(\mathcal{C})$ accept \mathcal{C} 's satisfiability, regardless of V 's random choices.
- **Argument Soundness.** If \mathcal{C} is not satisfiable, then for every malicious PPT P^* , the probability over V 's random choices that the interaction of $V(\mathcal{C})$ and $P^*(\mathcal{C})$ makes $V(\mathcal{C})$ accept \mathcal{C} as satisfiable is less than ϵ .

3 Protocol refinements in GINGER

In principle, PEPPER solves the problem of verified computation. The reality is less attractive: PEPPER's computational burden is high, its network costs are absurd, and its applicability is limited (to straight line numerical computations). Our system, GINGER, mitigates these issues: it lowers costs through protocol refinements (presented in this section), and it applies to a much wider class of computations (as we discuss in Section 4).

GINGER's refinements eliminate many queries, by relying on a new analysis of the base commitment primitive. To motivate this analysis, we note that there is something seemingly redundant in the base machinery (see Figure 1): why does the linear PCP require a linearity test (§2.2) if an honest prover depends on the linearity of its function π to pass the linear commitment protocol's consistency test (§2.3)? Can we remove one of these tests, or combine them somehow? The reason that

PEPPER appears to need both tests is that their guarantees are (so far) subtly different:

- *Consistency test* (§2.3): First, an honest prover is guaranteed to pass this test. Second, if the prover—even a cheating one—passes this test, then it is very likely bound to *some* function (as shown in [35, 45]).
- *Linearity test* (§2.2): This test is needed in case the prover cheats—it establishes that π is linear (as required by the rest of the PCP protocol). More accurately, if π is far from being linear, the test is somewhat likely to catch that case.

Yet, it seems unsatisfying that both tests are required when composing linear commitment and the linear PCP: can a prover really pass the consistency test systematically with a function that the linearity test would reject? In fact, our intuitive dissatisfaction is well-founded: this paper proves that the commitment primitive (which includes the consistency test) is far stronger than the linearity test. Put simply, even a cheating prover is very likely bound to a function that is linear, or almost so.

Practically, this result saves query generation and response costs. For one thing, we can eliminate linearity tests from the protocol. More significantly, we eliminate *amplification*: PEPPER needed to repeat the protocol to turn the linearity test's guarantee of "somewhat likely" into "very likely". In contrast, our result already gives a guarantee of "very likely", so no repetition is required.

More broadly, this result means that the commitment primitive is considerably more powerful than was realized—it efficiently commits an untrusted entity to a linear function and extracts evaluations of that function—and may apply elsewhere.

Details. The protocol refinements are rooted in a strengthened soundness analysis. *Soundness error* (for example, ϵ in Section 2.4) refers to the probability that a protocol or test succeeds when the condition that it is verifying or testing is actually false. The ideal is to have a small upper-bound on soundness error.

The soundness of the PCP protocol in Section 2.2 and Appendix A.1 is connected to the soundness of linearity testing [16]. Specifically, the base analysis proves that if the prover returns $y \neq \Psi(x)$, then the prover survives all tests (linearity, quadratic correction, circuit) with probability less than $7/9$ (requiring ρ runs to make $(7/9)^\rho$ small). The $7/9$ derives from a result [8] that if the proof oracle is not "somewhat close" to linear, then the linearity test passes with probability $< 7/9$ (though fascinating, this result is inconveniently weak in our context).

Our analysis, detailed in Appendix A.2, establishes that the commitment protocol binds the prover to a function that is extremely close to linear (otherwise, the prover could break the semantic security of the homo-

| | PEPPER [45] | GINGER |
|--|--|--|
| PCP encoding size (n) | $s^2 + s$, in general | $s^2 + s$, in general |
| V's per-instance CPU costs | | |
| Issue commit queries | $(e + 2c) \cdot n/\beta$ | $(e + 2c) \cdot n/\beta$ |
| Process commit responses | d | d |
| Issue PCP queries | $\rho \cdot (\chi \cdot f + \ell' \cdot f + 5c) \cdot n/\beta$ | $(\chi \cdot f + \ell \cdot f + 2c) \cdot n/\beta$ |
| Process PCP responses | $\rho \cdot (2\ell' + x + y) \cdot f$ | $(2\ell + x + y) \cdot f$ |
| P's per-instance CPU costs | | |
| Issue commit responses | $h \cdot n$ | $h \cdot n$ |
| Issue PCP responses | $(\rho \cdot \ell' + 1) \cdot f \cdot n$ | $(\ell + 1) \cdot f \cdot n$ |
| Network cost (per instance) | $((\rho \cdot \ell' + 1) \cdot p + \xi) \cdot n/\beta$ | $((\ell + 1) \cdot p + \xi) \cdot n/\beta$ |
| PCP soundness error | $(7/9)^\rho = 2.3 \cdot 10^{-8}$ | $\kappa = 2.6 \cdot 10^{-6}$ |
| Overall soundness error | $2.4 \cdot 10^{-8}$ | $4.5 \cdot 10^{-6}$ |

$|x|, |y|$: # of elements in input, output
 n : # of components in linear function π (§2.2)
 s : # of variables in constraint set (§2.1)
 χ : # of constraints in constraint set (§2.1)
 $\ell = 3$: # of high-order PCP queries in GINGER (§A.2, §A.3)
 $\ell' = 7$: # of high-order PCP queries in PEPPER (§A.1)
 $\rho = 70$: # of PCP reps. in base scheme (§A.1)
 β : batch size (# of instances) (§2.3)
 e : cost of encrypting an element in \mathbb{F}
 d : cost of decrypting an encrypted element
 f : cost of multiplying in \mathbb{F}
 h : cost of ciphertext add plus multiply
 c : cost to generate 192-bit pseudorandom #
 $|p|$: length of an element in \mathbb{F}
 $|\xi|$: length of an encrypted element in \mathbb{F}

Figure 2—High-order costs and error in GINGER, compared to its base (PEPPER [45]), for a computation represented as χ constraints over s variables (§2.1). The soundness error depends on field size (Appendix A.2); the table assumes $|\mathbb{F}| = 2^{128}$. Many of the cryptographic costs enter through the commitment protocol (see Section 2.3 or Figure 12); Section 6 quantifies the parameters. The “PCP” row include the consistency query and check. The network costs slightly underestimate by not including query responses.

morphic encryption used by GINGER and PEPPER). This results in the PCP soundness error improving from $7/9$ to κ , where $\kappa \approx 4\sqrt[6]{1/|\mathbb{F}|}$; this analysis does not depend on linearity tests, so they can be dropped.

The soundness error is somewhat low by cryptographic standards, but in practice, a failure rate (when the prover is malicious) of 1 in 200,000 is reasonable.

A further optimization. GINGER reuses some queries across the quadratic correction and circuit tests; this refinement is detailed and justified in Appendix A.3.

Savings. Most significantly, V can take advantage of the lower soundness error to run $\rho = 1$ instead of $\rho = 70$ repetitions of the PCP protocol. Also, per repetition, V 's work to generate pseudorandom queries decreases by $3/5$ ($2/5$ coming from the elimination of linearity tests and $1/5$ from reusing queries). These gains are depicted in Figure 2, most notably in the reduction from $\rho \cdot \ell' \approx 500$ to $\ell = 3$ total PCP queries.

The total savings for the verifier depend on the relative cost of pseudorandom number generation (encapsulated by c) and encryption (encapsulated by e). These savings show up in β^* , the minimum batch size (§2.3) at which V gains from outsourcing. As shown in Section 6.1, the reduction in β^* can be several orders of magnitude (when e is small). Finally, taking $|p| = 128$ bits and $|\xi| = 2 \cdot 1024$ bits, the savings in network costs are 1–2 orders of magnitude (holding β constant).

4 Broadening the space of computations

GINGER extends to computations over floating-point fractional quantities and to a restricted general-purpose programming model that includes inequality tests, log-

ical expressions, conditional branching, etc. To do so, GINGER maps computations to the constraint-over-finite-field formalism (§2.1), and thus the core protocol in Section 3 applies. In fact, our techniques³ apply to the many protocols that use the constraint formalism or arithmetic circuits. Moreover, we have implemented a compiler (derived from Fairplay's [39]) that transforms high-level computations first into constraints and then into verifier and prover executables.

The challenges of representing computations as constraints over finite fields include: the “true answer” to the computation may live outside of the field; sign and ordering in finite fields interact in an unintuitive fashion; and constraints are simply equations, so it is not obvious how to represent comparisons, logical expressions, and control flow. To explain GINGER's solutions, we first present an abstract framework that illustrates how GINGER broadens the set of computations soundly and how one can apply the approach to further computations.

Framework to map computations to constraints. To map a computation Ψ over some domain D (such as the integers, \mathbb{Z} , or the rationals, \mathbb{Q}) to equivalent constraints over a finite field, the programmer or compiler performs three steps, as illustrated and described below:

$$\begin{array}{ccc}
 \Psi \text{ over } D & \xrightarrow{(C1)} & \Psi \text{ over } U & \xrightarrow{(C2)} & \theta(\Psi) \text{ over } \mathbb{F} \\
 & & & & \downarrow (C3) \\
 & & & & \mathcal{C} \text{ over } \mathbb{F}
 \end{array}$$

³We suspect that many of the individual techniques are known. However, when the techniques combine, the material is surprisingly hard to get right, so we will delve into (excruciating) detail, consistent with our focus on built systems.

- C1 *Bound the computation.* Define a set $U \subset D$ and restrict the input to Ψ such that the output and intermediate values stay in U .
- C2 *Represent the computation faithfully in a suitable finite field.* Choose a finite field, \mathbb{F} , and a map $\theta: U \rightarrow \mathbb{F}$ such that computing $\theta(\Psi)$ over $\theta(U) \subset \mathbb{F}$ is isomorphic to computing Ψ over U . (By “ $\theta(\Psi)$ ”, we mean Ψ with all inputs and literals mapped by θ .)
- C3 *Transform the finite field version of the computation into constraints.* Write a set of constraints over \mathbb{F} that are equivalent (in the sense of Section 2.1) to $\theta(\Psi)$.

4.1 Signed integers and floating-point rationals

We now instantiate C1 and C2 for integer and rational number computations; the next section addresses C3.

Consider $m \times m$ matrix multiplication over N -bit signed integers. For step C1, each term in the output, $\sum_{k=1}^m A_{ik}B_{kj}$, has m additions of $2N$ -bit subterms so is contained in $[-m \cdot 2^{2N-1}, m \cdot 2^{2N-1}]$; this is our set U .

For step C2, take $\mathbb{F} = \mathbb{Z}/p$ (the integers mod a prime p , to be chosen shortly) and define $\theta: U \rightarrow \mathbb{Z}/p$ as $\theta(u) = u \bmod p$. Observe that θ maps negative integers to $\{\frac{p+1}{2}, \frac{p+3}{2}, \dots, p-1\}$, analogous to how processors represent negative numbers with a 1 in the most significant bit (this technique is standard [17, 50]). Of course, addition and multiplication in \mathbb{Z}/p do not “know” when their operands are negative. Nevertheless, the computation over \mathbb{Z}/p is isomorphic to the computation over U , provided that $|\mathbb{Z}/p| > |U|$ (as shown in Appendix B [46]).⁴ Thus, for the given U , we require $p > m \cdot 2^{2N}$. Note that a larger p brings larger costs (see Figure 2), so there is a three-way trade-off among p, m, N .

We now turn to rational numbers. For step C1, we restrict the inputs as follows: when written in lowest terms, their numerators are $(N_a + 1)$ -bit signed integers, and their denominators are in $\{1, 2, 2^2, 2^3, \dots, 2^{N_b}\}$. Note that such numbers are (primitive) floating-point numbers: they can be represented as $a \cdot 2^{-q}$, so the decimal point floats based on q . Now, for $m \times m$ matrix multiplication, the computation does not “leave” $U = \{a/b: |a| < 2^{N'_a}, b \in \{1, 2, 2^2, 2^3, \dots, 2^{N'_b}\}\}$, for $N'_a = 2N_a + 2N_b + \log_2 m$ and $N'_b = 2N_b$ [46, Appendix B].

For step C2, we take $\mathbb{F} = \mathbb{Q}/p$, the quotient field of \mathbb{Z}/p . Take $\theta(\frac{a}{b}) = (a \bmod p, b \bmod p)$. For any $U \subset \mathbb{Q}$, there is a choice of p such that the mapped computation over \mathbb{Q}/p is isomorphic to the original computation over \mathbb{Q} [46, Appendix B]. For our U above, $p > (m + 1)^2 \cdot 2^{4(N_a + N_b)}$ suffices.

Limitations and costs. To understand the limitations of GINGER’s floating-point representation, consider the number $a \cdot 2^{-q}$, where $|a| < 2^{N_a}$ and $|q| \leq N_q$.

To represent this number, the IEEE standard requires roughly $N_a + \log N_q$ bits [29] while GINGER requires $2 \cdot (\max(N_a, N_q) + 1)$ bits [46, Appendix B]. As a result, GINGER’s range is vastly more limited: with 64 bits, the IEEE standard can represent numbers on the order of 2^{1023} and 2^{-1022} (with $N_a = 53$ bits of precision) while 64 bits buys GINGER only numbers on the order of 2^{32} and 2^{-32} (with $N_a = 32$). Moreover, unlike the IEEE standard, GINGER does not support a division operation or rounding.

However, comparing GINGER’s floating-point representation to its *integer* representation, the extra costs are not terrible. First, the prover and verifier take an extra pass over the input and output (for implementation reasons; see Appendix B [46] for details). Second, a larger prime p is required. For example, $m \times m$ matrix multiplication with 32-bit integer inputs requires p to have at least $\log_2 m + 64$ bits; if the inputs are rationals with $N_a = N_q = 32$, then p requires $2 \log_2(m + 1) + 256$ bits. Roughly speaking, the end-to-end costs are $3 \times$ those of the integers case (see Section 6.2). Of course, the actual numbers depend on the computation. (Our compiler computes suitable bounds with static analysis.)

4.2 General-purpose program constructs

Case study: branch on order comparison. We now illustrate C3 with a case study of a computation, Ψ , that includes a less-than test and a conditional branch; pseudocode for Ψ is in Figure 3. For clarity, we will restrict Ψ to signed integers; handling rational numbers requires additional mechanisms [46, Appendix C].

How can we represent the test $x_1 < x_2$ using constraint *equations*? The solution is to use special *range constraints* that decompose a number into its bits to test whether it is in a given range; in this case, $C_<$, depicted in Figure 3, tests whether $e = \theta(x_1) - \theta(x_2)$ is in the “negative” range of \mathbb{Z}/p (see Section 4.1). Now, under the input restriction $x_1 - x_2 \in U$, $C_<$ is satisfiable if and only if $x_1 < x_2$ [46, Appendix C]. Analogously, we can construct C_{\geq} that is satisfiable if and only if $x_1 \geq x_2$.

Finally, we introduce a 0/1 variable M that encodes a choice of branch, and then arrange for M to “pull in” the constraints of that branch and “exclude” those of the other. (Note that the prover need not execute the untaken branch.) Figure 3 depicts the complete set of constraints, C_Ψ ; these constraints are satisfiable if and only if the prover correctly computes Ψ [46, Appendix C].

Logical expressions and conditionals. Besides order comparisons and if-else, GINGER can represent $=$, $\&\&$, and $||$ as constraints. An interesting case is $! =$: we can represent $Z_1 \neq Z_2$ with $\{M \cdot (Z_1 - Z_2) - 1 = 0\}$ because this constraint is satisfiable when $(Z_1 - Z_2)$ has a multiplicative inverse and hence is not zero. These constructs and others are detailed in Appendix D [46].

⁴For space, Appendices B–E appear only in the extended version [46].

$$\begin{array}{l}
\Psi : \\
\text{if } (X_1 < X_2) \\
\quad Y = 3 \\
\text{else} \\
\quad Y = 4
\end{array}
\quad
\mathcal{C}_\zeta = \left\{ \begin{array}{l} B_0(1 - B_0) = 0, \\ B_1(2 - B_1) = 0, \\ \vdots \\ B_{N-2}(2^{N-2} - B_{N-2}) = 0, \\ \theta(X_1) - \theta(X_2) - (p - 2^{N-1}) - \sum_{i=0}^{N-2} B_i = 0 \end{array} \right\}
\quad
\mathcal{C}_\Psi = \left\{ \begin{array}{l} M\{\mathcal{C}_\zeta\}, \\ M(Y - 3) = 0, \\ (1 - M)\{\mathcal{C}_{>=}\}, \\ (1 - M)(Y - 4) = 0 \end{array} \right\}$$

Figure 3—Pseudocode for our case study of Ψ , and corresponding constraints \mathcal{C}_Ψ . Ψ 's inputs are signed integers x_1, x_2 ; per steps C1 and C2 (§4.1), we assume $x_1 - x_2 \in U \subset [-2^{N-1}, 2^{N-1}]$, where $p > 2^N$. The constraints \mathcal{C}_ζ test $x_1 < x_2$ by testing whether the bits of $\theta(x_1) - \theta(x_2)$ place it in $[p - 2^{N-1}, p)$. $M\{\mathcal{C}\}$ means multiplying all constraints in \mathcal{C} by M and then reducing to degree-2.

Limitations and costs. We compile a subset of SFDL, the language of the Fairplay compiler [39]. Thus, our limitations are essentially those of SFDL; notably, loop bounds have to be known at compile time.

How efficient is our representation? The program constructs above mostly have concise constraint representations. Consider, for instance, `comp1==comp2`; the equivalent constraint set \mathcal{C} consists of the constraints that represent `comp1`, the constraints that represent `comp2`, and an additional constraint to relate the outputs of `comp1` and `comp2`. Thus, \mathcal{C} is the same size as its two components, as one would expect.

However, two classes of computations are costly. First, inequality comparisons require variables and a constraint for every bit position; see Figure 3. Second, the constraints for `if-else` and `||`, as written, seem to be degree-3; notice, for instance, the $M\{\mathcal{C}\}$ in Figure 3. To be compatible with the core protocol, these constraints must be rewritten to be degree-2 (§2.1), which carries costs. Specifically, if \mathcal{C} has s variables and χ constraints, an equivalent degree-2 representation of $M\{\mathcal{C}\}$ has $s + \chi$ variables and $2 \cdot \chi$ constraints [46, Appendix D].

5 Parallelization and implementation

Many of GINGER's remaining costs are in the cryptographic operations in the commitment protocol (see Appendix A.1). To address these costs, we distribute the prover over multiple machines, leveraging GINGER's inherent parallelism. We also implement the prover and verifier on GPUs, which raises two questions. (1) Isn't this just moving the problem? Yes, and this is good: GPUs are optimized for the types of operations that bottleneck GINGER. (2) Why do we assume that the *verifier* has a GPU? Desktops are more likely than servers to have GPUs, and the prevalence of GPUs is increasing. Also, this setup models a future in which specialized hardware for cryptographic operations is common.

Parallelization. To distribute GINGER's prover, we run multiple copies of it (one per host), each copy receiving a fraction of the batch (Section 2.3). In this configuration, the provers use the Open MPI [2] message-passing library to synchronize and exchange data.

To further reduce latency, each prover offloads work to a GPU (see also [49] for an independent study of GPU

hardware in the context of [21]). We exploit three levels of parallelism here. First, the prover performs a ciphertext operation for each component in the commitment vector (§2.3); each operation is (to first approximation) separate. Second, each operation computes two independent modular exponentiations (the ciphertext of an ElGamal encryption has two elements). Third, modular exponentiation itself admits a parallel implementation (each input is a multiprecision number encoded in multiple machine words). Thus, in our GPU implementation, a group of CUDA [1] threads computes each exponentiation.

We also parallelize the verifier's encryption work during the commitment phase (§2.3), using the approach above plus an optimization: the verifier's exponentiations are fixed base, letting us memoize intermediate squares. We implement exponentiations for the prover and verifier with the `libgpcrypto` library of `SSLShader` [36], modified to implement the memoization.

Implementation details. Our compiler consists of two stages, which a future publication will detail. The front-end compiles a subset of Fairplay's SFDL [39] to constraints; it is derived from Fairplay and is implemented in 5294 lines of Java, starting from Fairplay's 3886 lines (per [51]). The back-end transforms constraints into C++ code that implements the verifier and prover and then invokes `gcc`; this component is 1105 lines of Python code.

For efficiency, PEPPER [45] introduced specialized PCP protocols for certain computations. For some experiments we use specialized PCPs in GINGER also; in these cases we write the prover and verifier manually, which typically requires a few hundred lines of C++. Automating the compilation of specialized PCPs is future work.

The verifier and prover are separate processes that exchange data using Open MPI [2]. GINGER uses the ElGamal cryptosystem [23] with 1024-bit keys.

6 Experimental evaluation

Our evaluation answers the following questions:

- What is the effect of the protocol refinements (§3)?
- What are the costs of supporting rational numbers and the additional program structures (§4)?
- What is GINGER's speedup from parallelizing (§5)?

Figure 4 summarizes the results.

| | |
|--|------|
| GINGER’s protocol refinements reduce per-instance network costs by 25–30× (to hundreds of KBs for the computations we study), prover CPU costs by about 10–14% (leaving them still high), and break-even batch size (β^*) by about 4×. | §6.1 |
| With accelerated encryption GINGER breaks even from outsourcing short computations at small batch sizes; for 400×400 matrix multiplication, the verifier gains from outsourcing at a batch size of 20 (tens of seconds of computation). | §6.1 |
| Rational arithmetic costs roughly 3× integer arithmetic under GINGER (but much more than native floating-point). | §6.2 |
| Parallelizing results in near-linear reduction in the prover’s latency. | §6.3 |

Figure 4—Summary of main evaluation results.

| computation (Ψ) | $O(\cdot)$ | input domain (see §4.1) | size of \mathbb{F} | s | n | default | local |
|-------------------------------|------------|------------------------------------|----------------------|----------------------------------|-------------------------------------|-----------|---------|
| matrix mult. | $O(m^3)$ | 32-bit signed integers | 128 bits | $2m^2$ | m^3 | $m = 200$ | 800 ms |
| matrix mult. (\mathbb{Q}) | $O(m^3)$ | rationals ($N_a = 32, N_b = 32$) | 320 bits | $2m^2$ | m^3 | $m = 100$ | 5.90 ms |
| deg-2 poly. eval. | $O(m^2)$ | 32-bit signed integers | 128 bits | m | m^2 | $m = 100$ | 0.40 ms |
| deg-3 poly. eval. | $O(m^3)$ | 32-bit signed integers | 192 bits | m | m^3 | $m = 200$ | 160 ms |
| m -Hamming dist. | $O(m^2)$ | 32-bit unsigned | 128 bits | $2m^2 + m$ | $2m^3$ | $m = 100$ | 0.90 ms |
| bisection method | $O(m^2)$ | rationals ($N_a = 32, N_b = 5$) | 320 bits | $16 \cdot (m + \mathcal{C}_<)$ | $256 \cdot (m + \mathcal{C}_<)^2$ | $m = 25$ | 180 ms |

Figure 5—Benchmark computations. s is the number of constraint variables; s affects n , which is the size of V ’s queries and of P ’s linear function π (see Figure 2). Only high-order terms are reported for n . The latter two columns give our experimental defaults and the cost of local computation (i.e., no outsourcing) at those defaults. In polynomial evaluation, V and P hold a polynomial; the input is values for the m variables. The latter two computations exercise the program constructs in Section 4.2. In m -Hamming distance, V and P hold a fixed set of strings; the input is a length m string, and the output is a vector of the Hamming distance between the input and the set of strings. Bisection method refers to root-finding via bisection: both V and P hold a degree-2 polynomial in m variables, the input is two m -element endpoints that bracket a root, and the output is a small interval that contains the root.

We use six benchmark computations, summarized in Figure 5 (Appendix E [46] has details). For bisection method and degree-2 polynomial evaluation, V and P were produced by our compiler; for the other computations, we use tailored encodings (see Section 5). We implemented and analyzed other computations (e.g., edit distance and circle packing) but found that V gained from outsourcing only at implausibly large batch sizes.

Method and setup. We measure latency and computing cycles used by the verifier and the prover, and the amount of data exchanged between them. We account for the prover’s cost in per-instance terms. Because the verifier amortizes costs over a batch (§2.3), we focus on the *break-even batch size*, β^* : the batch size at which the verifier’s CPU cost from GINGER equals the cost of computing the batch locally. We measure local computation using implementations built on the GMP library (except for matrix multiplication over rationals, where we use native floating-point).

For each result that we report, we run at least three experiments and take the averages (the standard deviations are always within 5% of the means). We measure CPU time using `getrusage`, latency using PAPI’s real time counter [3], and network costs by recording the number of application-level bytes transferred.

Our experiments use a cluster at the Texas Advanced Computing Center (TACC). Each machine is configured identically and runs Linux on an Intel Xeon processor E5540 2.53 GHz with 48GB of RAM. Experiments with GPUs use machines with an NVIDIA Tesla M2070. Each

GPU has 448 CUDA cores and 6GB of memory.

Validating the cost model. We will sometimes predict β^* , V ’s costs, and P ’s costs by using our cost model (Figure 2), so we now validate this model. We run microbenchmarks to quantify the model’s parameters— e is reported in this section; Appendix E [46] quantifies the other parameters—and then compare the parameterized model to GINGER’s measured performance. GINGER’s empirical results are at most 2%–15% more than are predicted by the model. However, local computation costs about 1.2–4.0 times more than is predicted; we think that the divergence results from adverse caching effects that increase the cost of a multiplication. Thus, we expect the verifier to break even at batch sizes that are about a factor of 1.2–4.0 smaller than predicted by the model.

6.1 The effect of GINGER’s protocol refinements

We begin with $m \times m$ matrix multiplication ($m = 100, 200$) and degree-3 polynomial evaluation ($m = 100, 200$), and batch size of $\beta = 5000$. We report *per-instance* network and CPU costs: the total network and CPU costs over the batch, divided by β .

Figure 6 depicts network costs. For matrix multiplication, these are about the same as the cost to send the inputs and receive the outputs; for polynomial evaluation, these are about 10 times the size of the inputs and outputs. Also, GINGER improves on PEPPER by 20–30×.

In this experiment, GINGER’s prover incurs about 10–14% less CPU time compared to PEPPER (estimated using a cost model from [45]) but still takes tens of minutes per-instance; this is obviously a lot, but we reduce

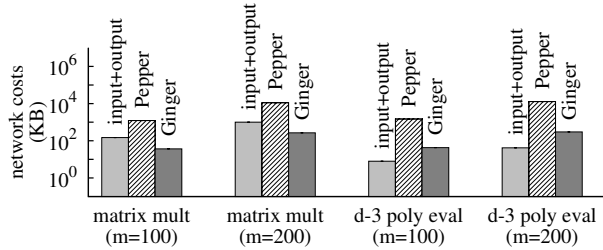


Figure 6—Per-instance network costs of GINGER and its base (PEPPER [45]), compared to the size of the inputs and outputs. At this batch size ($\beta = 5000$), GINGER’s refinements reduce per-instance network costs by a factor of 25–30 compared to PEPPER. GINGER’s network costs here are hundreds of KB or less. The y-axis is log-scaled.

| | PEPPER | GINGER | |
|---------------------|---------------------|--------|----------|
| local | 1.1 s | 1.1 s | |
| CPU | β^* | 13000 | 4100 |
| | verifier aggregate | 3.9 hr | 1.3 hr |
| | prover aggregate | 5.0 yr | 1.6 yr |
| | prover per-instance | 3.5 hr | 3.3 hr |
| GPU | β^* | 8700 | 2300 |
| | verifier aggregate | 2.7 hr | 43.4 min |
| | prover aggregate | 3.5 yr | 320 days |
| | prover per-instance | 3.5 hr | 3.3 hr |
| crypto hardware | β^* | 3900 | 20 |
| | verifier aggregate | 1.2 hr | 22.3 s |
| | prover aggregate | 1.6 yr | 2.8 days |
| prover per-instance | 3.5 hr | 3.3 hr | |

Figure 7—Break-even batch sizes (β^*) and predicted running times of prover and verifier at $\beta = \beta^*$, for matrix multiplication ($m = 400$), under three models of the encryption cost. The verifier’s per-instance work is not depicted because it equals the local running time, by definition of β^* . The local running time is high in part because the local implementation uses GMP.

latency by parallelizing (§6.3). For this computation and at this batch size ($\beta = 5000$), GINGER’s verifier takes a few hundreds of milliseconds per-instance, less than locally computing using our baseline of GMP.

Amortizing the verifier’s costs. Batching is both a limitation and a strength of GINGER: GINGER’s verifier *must* batch to gain from outsourcing but *can* batch to drive per-instance overhead arbitrarily low. Nevertheless, we want break-even batch sizes (β^*) to be as small as possible. But β^* mostly depends on e , the cost of encryption (Figure 2), because after our refinements the verifier’s main burden is creating $\text{Enc}(pk, r)$ (see §2.3), the cost of which amortizes over the batch.

What values of e make sense? We consider three scenarios: (1) the verifier uses a CPU for encryptions, (2) the verifier offloads encryptions to a GPU, and (3) the verifier has special-purpose hardware that can *only* perform encryptions. (See Section 5 for motivation.) Under scenario (1), we measure $e = 72.1\mu\text{s}$ on a 2.5 GHz CPU.

| | mat. mult. | mat. mult. (\mathbb{Q}) |
|-----------------------|------------|-----------------------------|
| local | 17.6 ms | 5.90 ms |
| verifier per-instance | 17.6 ms | 80.2 ms |
| verifier aggregate | 76.1 s | 5.7 min |
| prover per-instance | 3.1 min | 9.4 min |
| prover aggregate | 9.3 days | 28 days |

Figure 8—Predicted running times of GINGER’s verifier and prover for matrix multiplication ($m = 100$), under integer and floating-point inputs, at $\beta = 4300$ (the break-even batch size for this computation over integers). The “local” row refers to GMP arithmetic for \mathbb{Z} and native floating-point arithmetic for \mathbb{Q} . Handling rationals costs GINGER roughly $3\times$ more than handling integers, but both are still far from native.

| computation (Ψ) | # Boolean gates (est.) | # constraint vars. |
|------------------------|------------------------|--------------------|
| m -Hamming dist. | $1.3 \cdot 10^6$ | $2 \cdot 10^4$ |
| bisection method | $3.0 \cdot 10^8$ | 1528 |

Figure 9—GINGER’s constraints compared to Boolean circuits, for m -Hamming distance ($m = 100$) and bisection method ($m = 25$). The Boolean circuits are estimated using the unmodified Fairplay [39] compiler. GINGER’s constraints are not concise but are far more so than Boolean circuits.

Under scenario (3), we take $e = 0\mu\text{s}$. What about scenario (2)? Our cost model concerns *CPU* costs, so we need an exchange rate between GPU and CPU exponentiations. We make a crude estimate: we measure the number of encryptions per second achievable on an NVIDIA Tesla M2070 (which is 180,000) and on an Intel 2.5 GHz CPU (which is 13,700), normalize by the dollar cost of the chips, and obtain that their throughput-per-dollar ratio is $1.8\times$. We thus (very conservatively) take $e = 72.1/1.8 = 40\mu\text{s}$.

We plug these three values of e into the cost model in Figure 2, set the cost under GINGER equal to the cost of local computing, and solve for β^* . The values of β^* are 4150 (CPU), 2300 (crude GPU estimate), and 20 (crypto hardware). We also use the model to predict V ’s and P ’s costs at β^* , under PEPPER and GINGER. Figure 7 summarizes. GINGER is very sensitive to the value of e because its refinements have eliminated many of the other costs. Moreover, the aggregate verifier computing time drops significantly under all three cost models. The prover’s per-instance work is mostly unaffected, but as the batch size decreases, so does its aggregate work.

6.2 Evaluating GINGER’s computational model

To understand the costs of the floating-point representation (§4.1), we compare it to two baselines: GINGER’s signed integer representation and the computation executed locally, using the CPU’s floating point unit. Our benchmark application is matrix multiplication ($m = 100$). Figure 8 details the comparison.

We also consider GINGER’s general-purpose program constructs (§4). Our baseline is *Boolean* circuits (we are

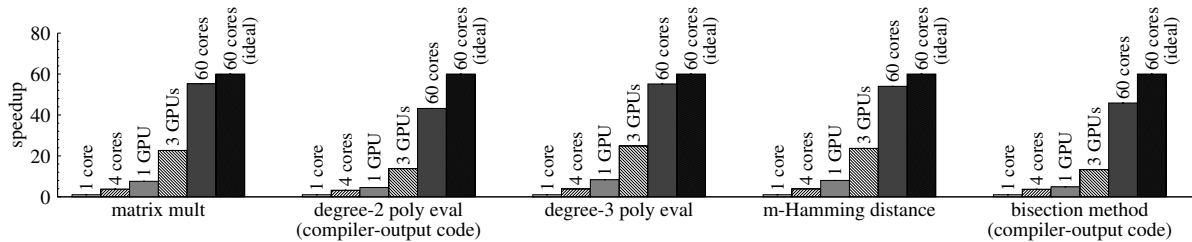


Figure 10—Latency speedup observed by GINGER’s verifier when the prover is parallelized. We run with $m = 100, \beta = 150$ for matrix multiplication and degree-3 polynomial evaluation; $m = 100, \beta = 1500$ for degree-2 polynomial evaluation; $m = 100, \beta = 15$ for m -Hamming distance; and $m = 25, \beta = 15$ for bisection method. GINGER’s prover achieves near-linear speedups except when the problem sizes are small and hence the overhead from parallelizing is significant (e.g., degree-2 polynomial evaluation).

unaware of efficient arithmetic representations of these constructs). We compare the number of Boolean circuit *gates* and the number of GINGER’s arithmetic constraint *variables*, since these determine the proving and verifying costs under the respective formalisms (see [5, 45]). Taken individually, GINGER’s constructs (\leq , $\&\&$, etc.) are the same cost or more than those of Boolean circuits (e.g., $\|\|$ introduces auxiliary variables). However, Boolean circuits are in general far more verbose: they represent quantities by their bits (which GINGER does only when computing inequalities). Figure 9 gives a rough end-to-end comparison.

6.3 Scalability of the parallel implementation

To demonstrate the scalability of GINGER’s parallelization, we run the prover using many CPU cores, many GPUs, and many machines. We measure end-to-end latency, as observed by the verifier. Figure 10 summarizes the results for various computations. In most cases, the speedup is near-linear.

7 Related work

A substantial body of work achieves two of our goals—it is general-purpose and practical—but it makes strong assumptions about the servers (e.g., trusted hardware). There is also a large body of work on protocols for special-purpose computation. We regard this work as orthogonal to our efforts; for a survey of this landscape, see [45]. Herein, we focus on approaches that are general-purpose and unconditional.

Homomorphic encryption and secure multi-party protocols. Homomorphic encryption (which enables computation over ciphertext) and secure multi-party protocols (in which participants compute over private data, revealing only the result [34, 39, 52]) provide only *privacy* guarantees, but one can build on them for verifiable computation. For instance, the Boneh-Goh-Nissim homomorphic cryptosystem [18] can be adapted to evaluate circuits, Groth uses homomorphic commitments to produce a zero-knowledge argument protocol [33], and Applebaum et al. use secure multi-party protocols for ver-

ifying computations [4]. Also, Gentry’s fully homomorphic encryption [27] has engendered protocols for verifiable non-interactive computation [20, 24, 26]. However, despite striking improvements [28, 42, 47], the costs of hiding inputs (among other expenses) prevent any of the aforementioned verified computation schemes from getting close to practical (even by our relaxed standards).

PCPs, argument systems, and interactive proofs. Applying proof systems to verifiable computation is standard in the theory community [5–7, 10, 15, 32, 37, 38, 41], and the asymptotics continue to improve [13, 14, 22, 43]. However, none of this work has paid much attention to building systems.

Very recently, researchers have begun to explore using this theory for practical verified outsourced computation. In a recent preprint, Ben-Sasson et al. [12] investigate when PCP protocols might be beneficial for outsourcing. Since many of the protocols require representing computations as constraints, Ben-Sasson et al. [11] study improved reductions to constraints from a RAM model of computation. And Gennaro et al. [25] give a new characterization of NP to provide asymptotically efficient arguments without using PCPs.

However, as far as we know, only two research groups have made serious efforts toward practical systems. Our previous work [44, 45] built upon the efficient argument system of Ishai et al. [35]. In contrast, Cormode, Mitzenmacher, and Thaler [21] (hereafter, CMT) built upon the protocol of Goldwasser et al. [31], and a follow-up effort studies a GPU-based parallel implementation [49].

Comparison of GINGER and CMT [21, 49]. We compared three different implementations: *CMT-native*, *CMT-GMP*, and GINGER. *CMT-native* refers to the code and configuration released by Thaler et al. [49]; it works over a small field and thereby exploits highly efficient machine arithmetic but restricts the inputs to the computation unrealistically (see Section 4.1). *CMT-GMP* refers to an implementation based on *CMT-native* but modified by us to use the GMP library for multi-precision arithmetic; this allows more realistic computation sizes and inputs, as well as rational numbers.

| m | domain | component | CMT-native | CMT-GMP | GINGER |
|-----|--------------|-----------|------------|---------|--------|
| 256 | \mathbb{Z} | verifier | 40 ms | 0.6 s | 0.3 s |
| | | prover | 22 min | 2.5 hr | 36 min |
| | | network | 88 KB | 0.3 MB | 1.1 MB |
| 128 | \mathbb{Q} | verifier | – | 260 ms | 190 ms |
| | | prover | – | 1.0 hr | 21 min |
| | | network | – | 1.8 MB | 1.4 MB |

Figure 11—CMT [21] compared to GINGER, in terms of *amortized* CPU and network costs (GINGER’s total costs are divided by a batch size of $\beta=5000$ instances), for $m \times m$ matrix multiplication. CMT-native uses native data types but is restricted to small problem sizes and domains. CMT-GMP uses the GMP library for multi-precision arithmetic (as does GINGER).

We perform two experiments using $m \times m$ matrix multiplication. Our testbed is the same as in Section 6. In the first one, we run with $m = 256$ and integer inputs. For CMT-GMP and GINGER, the inputs are 32-bit unsigned integers, and the prime (the field modulus) is 128 bits. For CMT-native, the prime is $2^{61} - 1$. In the second experiment, m is 128, the inputs are rational numbers (with $N_a = N_b = 32$; see Section 4.1), the prime is 320 bits, and we experiment only with CMT-GMP and GINGER.

We measure total CPU time and network cost; for CMT, we measure “network” traffic by counting bytes (the CMT verifier and prover run in the same process and hence the same machine). Each reported datum is an average over 3 sample runs; there is little experimental variation (less than 5% of the means).

Figure 11 depicts the results. CMT incurs a significant penalty when moving from native to GMP (and hence to realistic problem sizes). Comparing CMT-GMP and GINGER, the network and prover costs are similar (although network costs for CMT reflect high fixed overhead for their circuit). The *per-instance* verifier costs are also similar, but GINGER is batch verifying whereas CMT does not need to do so (a significant advantage).

A qualitative comparison is as follows. On the one hand, CMT does not require cryptography, has better asymptotic prover and network costs, and for some computations the verifier does not need batching to gain from outsourcing [49]. On the other hand, CMT applies to a smaller set of computations: if the computation is not efficiently parallelizable or does not naturally map to arithmetic circuits (e.g., it has order comparisons or conditionality), then CMT in its current form will be inapplicable or inefficient, respectively. Ultimately, GINGER and CMT should be complementary, as one can likely ease or eliminate some of the restrictions on CMT by incorporating the constraint formalism together with batching [48].

8 Summary and conclusion

This paper is a contribution to the emerging area of practical PCP-based systems for unconditional verifiable

computation. GINGER has combined theoretical refinements (slashing query costs and network overhead); a general computational model (including fractions and standard program constructs) with a compiler; and a massively parallel implementation that takes advantage of modern hardware. Together, these changes have brought us closer to a truly deployable system. Nevertheless, much work remains: the efficiency of the verifier depends on special hardware, the costs for the prover are still too high, and looping cannot yet be handled concisely.

Acknowledgments

Detailed attention from Edmund L. Wong substantially clarified this paper. Yuval Ishai, Mike Lee, Bryan Parno, Mark Silberstein, Chung-chieh (Ken) Shan, Sara L. Su, Justin Thaler, and the anonymous reviewers gave useful comments that improved this draft. The Texas Advanced Computing Center (TACC) at UT supplied computing resources. We thank Jane-ellen Long, of USENIX, for her good nature and inexhaustible patience. The research was supported by AFOSR grant FA9550-10-1-0073 and by NSF grants 1055057 and 1040083.

Our code and experimental configurations are available at <http://www.cs.utexas.edu/pepper>

References

- [1] CUDA (<http://developer.nvidia.com/what-cuda>).
- [2] Open MPI (<http://www.open-mpi.org>).
- [3] PAPI: Performance Application Programming Interface.
- [4] B. Applebaum, Y. Ishai, and E. Kushilevitz. From secrecy to soundness: efficient verification via secure computation. In *ICALP*, 2010.
- [5] S. Arora, C. Lund, R. Motwani, M. Sudan, and M. Szegedy. Proof verification and the hardness of approximation problems. *J. of the ACM*, 45(3):501–555, May 1998.
- [6] S. Arora and S. Safra. Probabilistic checking of proofs: a new characterization of NP. *J. of the ACM*, 45(1):70–122, Jan. 1998.
- [7] L. Babai, L. Fortnow, L. A. Levin, and M. Szegedy. Checking computations in polylogarithmic time. In *STOC*, 1991.
- [8] M. Bellare, D. Coppersmith, J. Håstad, M. Kiwi, and M. Sudan. Linearity testing in characteristic two. *IEEE Transactions on Information Theory*, 42(6):1781–1795, Nov. 1996.
- [9] M. Bellare, S. Goldwasser, C. Lund, and A. Russell. Efficient probabilistically checkable proofs and applications to approximations. In *STOC*, 1993.
- [10] M. Ben-Or, S. Goldwasser, J. Kilian, and A. Wigderson. Multi-prover interactive proofs: how to remove intractability assumptions. In *STOC*, 1988.
- [11] E. Ben-Sasson, A. Chiesa, D. Genkin, and E. Tromer. Fast reductions from RAMs to delegatable succinct constraint satisfaction problems. Feb. 2012. Cryptology eprint 071.
- [12] E. Ben-Sasson, A. Chiesa, D. Genkin, and E. Tromer. On the concrete-efficiency threshold of probabilistically-checkable proofs. *ECCC*, (045), Apr. 2012.
- [13] E. Ben-Sasson, O. Goldreich, P. Harsha, M. Sudan, and S. Vadhan. Robust PCPs of proximity, shorter PCPs and applications to coding. *SIAM J. on Comp.*, 36(4):889–974, Dec. 2006.

- [14] E. Ben-Sasson and M. Sudan. Short PCPs with polylog query complexity. *SIAM J. on Comp.*, 38(2):551–607, May 2008.
- [15] M. Blum and S. Kannan. Designing programs that check their work. *J. of the ACM*, 42(1):269–291, 1995.
- [16] M. Blum, M. Luby, and R. Rubinfeld. Self-testing/correcting with applications to numerical problems. *J. of Comp. and Sys. Sciences*, 47(3):549–595, Dec. 1993.
- [17] D. Boneh and D. M. Freeman. Homomorphic signatures for polynomial functions. In *EUROCRYPT*, 2011.
- [18] D. Boneh, E. J. Goh, and K. Nissim. Evaluating 2-DNF formulas on ciphertexts. In *TCC*, 2005.
- [19] G. Brassard, D. Chaum, and C. Crépeau. Minimum disclosure proofs of knowledge. *J. of Comp. and Sys. Sciences*, 37(2):156–189, 1988.
- [20] K.-M. Chung, Y. Kalai, and S. Vadhan. Improved delegation of computation using fully homomorphic encryption. In *CRYPTO*, 2010.
- [21] G. Cormode, M. Mitzenmacher, and J. Thaler. Practical verified computation with streaming interactive proofs. In *ITCS*, 2012.
- [22] I. Dinur. The PCP theorem by gap amplification. *J. of the ACM*, 54(3), June 2007.
- [23] T. ElGamal. A public key cryptosystem and a signature scheme based on discrete logarithms. *IEEE Transactions on Information Theory*, 31:469–472, 1985.
- [24] R. Gennaro, C. Gentry, and B. Parno. Non-interactive verifiable computing: Outsourcing computation to untrusted workers. In *CRYPTO*, 2010.
- [25] R. Gennaro, C. Gentry, B. Parno, and M. Raykova. Quadratic span programs and succinct NIZKs without PCPs. Apr. 2012. Cryptology eprint 215.
- [26] R. Gennaro and D. Wichs. Fully homomorphic message authenticators. May 2012. Cryptology eprint 290.
- [27] C. Gentry. *A fully homomorphic encryption scheme*. PhD thesis, Stanford University, 2009.
- [28] C. Gentry, S. Halevi, and N. Smart. Homomorphic evaluation of the AES circuit. In *CRYPTO*, 2012.
- [29] D. Goldberg. What every computer scientist should know about floating-point arithmetic. *ACM Computing Surveys*, 23(1):5–48, Mar. 1991.
- [30] O. Goldreich. *Foundations of Cryptography: II Basic Applications*. Cambridge University Press, 2004.
- [31] S. Goldwasser, Y. T. Kalai, and G. N. Rothblum. Delegating computation: Interactive proofs for muggles. In *STOC*, 2008.
- [32] S. Goldwasser, S. Micali, and C. Rackoff. The knowledge complexity of interactive proof systems. *SIAM J. on Comp.*, 18(1):186–208, 1989.
- [33] J. Groth. Linear algebra with sub-linear zero-knowledge arguments. In *CRYPTO*, 2009.
- [34] Y. Huang, D. Evans, J. Katz, and L. Malka. Faster secure two-party computation using garbled circuits. In *USENIX Security*, 2011.
- [35] Y. Ishai, E. Kushilevitz, and R. Ostrovsky. Efficient arguments without short PCPs. In *Conference on Computational Complexity (CCC)*, 2007.
- [36] K. Jang, S. Han, S. Han, S. Moon, and K. Park. SSLShader: Cheap SSL acceleration with commodity processors. In *NSDI*, 2011.
- [37] J. Kilian. A note on efficient zero-knowledge proofs and arguments (extended abstract). In *STOC*, 1992.
- [38] J. Kilian. Improved efficient arguments (preliminary version). In *CRYPTO*, 1995.
- [39] D. Malkhi, N. Nisan, B. Pinkas, and Y. Sella. Fairplay—a secure two-party computation system. In *USENIX Security*, 2004.
- [40] R. C. Merkle. Digital signature based on a conventional encryption function. In *CRYPTO*, 1987.
- [41] S. Micali. Computationally sound proofs. *SIAM J. on Comp.*, 30(4):1253–1298, 2000.
- [42] M. Naehrig, K. Lauter, and V. Vaikuntanathan. Can homomorphic encryption be practical? In *ACM Workshop on Cloud Computing Security*, 2011.
- [43] A. Polishchuk and D. A. Spielman. Nearly-linear size holographic proofs. In *STOC*, 1994.
- [44] S. Setty, A. J. Blumberg, and M. Walfish. Toward practical and unconditional verification of remote computations. In *HotOS*, 2011.
- [45] S. Setty, R. McPherson, A. J. Blumberg, and M. Walfish. Making argument systems for outsourced computation practical (sometimes). In *NDSS*, 2012.
- [46] S. Setty, V. Vu, N. Panpalia, B. Braun, A. J. Blumberg, and M. Walfish. Taking proof-based verified computation a few steps closer to practicality (extended version). Technical Report TR-12-14, Dept. of CS, UT Austin, June 2012.
- [47] N. Smart and F. Vercauteren. Fully homomorphic SIMD operations. Aug. 2011. Cryptology eprint 133.
- [48] J. Thaler. Personal communication, June 2012.
- [49] J. Thaler, M. Roberts, M. Mitzenmacher, and H. Pfister. Verifiable computation with massively parallel interactive proofs. In *USENIX HotCloud Workshop*, June 2012. Full paper at <http://arxiv.org/abs/1202.1350>, Feb. 2012.
- [50] C. Wang, K. Ren, J. Wang, and K. M. R. Urs. Harnessing the cloud for securely outsourcing large-scale systems of linear equations. In *Intl. Conf. on Dist. Computing Sys. (ICDCS)*, 2011.
- [51] D. A. Wheeler. SLOCCount.
- [52] A. C.-C. Yao. How to generate and exchange secrets. In *FOCS*, 1986.

A Efficient arguments with linear PCPs but no linearity tests

Whereas previous work [35, 45] established that the commitment protocol in phases 2 and 3 of PEPPER (§2.3) binds the prover to a particular function, there were no constraints on that function. The principal result of this section is that the prover is actually bound to a function that is linear, or very nearly so. As a consequence, we can eliminate linearity testing from the PCP protocol. Furthermore, the error bound from one run of this modified PCP protocol is far stronger (lower) than was known.

This section describes the base protocols (A.1), states the refinements and proves their soundness (A.2), and describes a few other optimizations (A.3).

A.1 Base protocols

GINGER uses a linear commitment protocol that is borrowed from PEPPER [45]; this protocol is depicted in Figure 12.⁵ As described in Section 2.3, PEPPER composes this protocol and a linear PCP; that PCP is depicted in Figure 13. The purpose of $\{\gamma_0, \gamma_1, \gamma_2\}$ in this figure is to make a maliciously constructed oracle unlikely to pass

⁵Like PEPPER, GINGER verifies in batches (§2.3), which changes the protocols a bit; see [45, Appendix C] for details.

Commit+Multidecommit

The protocol assumes an additive homomorphic encryption scheme ($\text{Gen}, \text{Enc}, \text{Dec}$) over a finite field, \mathbb{F} .

Commit phase

Input: Prover holds a vector $w \in \mathbb{F}^n$, which defines a linear function $\pi: \mathbb{F}^n \rightarrow \mathbb{F}$, where $\pi(q) = \langle w, q \rangle$.

1. Verifier does the following:
 - Generates public and secret keys $(pk, sk) \leftarrow \text{Gen}(1^k)$, where k is a security parameter.
 - Generates vector $r \in_R \mathbb{F}^n$ and encrypts r component-wise, so $\text{Enc}(pk, r) = (\text{Enc}(pk, r_1), \dots, \text{Enc}(pk, r_n))$.
 - Sends $\text{Enc}(pk, r)$ and pk to the prover.
2. Using the homomorphism in the encryption scheme, the prover computes $e \leftarrow \text{Enc}(pk, \pi(r))$ without learning r . The prover sends e to the verifier.
3. The verifier computes $s \leftarrow \text{Dec}(sk, e)$, retaining s and r .

Decommit phase

Input: the verifier holds $q_1, \dots, q_\mu \in \mathbb{F}^n$ and wants to obtain $\pi(q_1), \dots, \pi(q_\mu)$.

4. The verifier picks μ secrets $\alpha_1, \dots, \alpha_\mu \in_R \mathbb{F}$ and sends to the prover (q_1, \dots, q_μ, t) , where $t = r + \alpha_1 q_1 + \dots + \alpha_\mu q_\mu \in \mathbb{F}^n$.
5. The prover returns $(a_1, a_2, \dots, a_\mu, b)$, where $a_i, b \in \mathbb{F}$. If the prover behaved, then $a_i = \pi(q_i)$ for all $i \in [\mu]$, and $b = \pi(t)$.
6. The verifier checks: $b \stackrel{?}{=} s + \alpha_1 a_1 + \dots + \alpha_\mu a_\mu$. If so, it outputs (a_1, a_2, \dots, a_μ) . If not, it rejects, outputting \perp .

Figure 12—The commitment protocol of PEPPER [45], which generalizes a protocol of Ishai et al. [35]. q_1, \dots, q_μ are the PCP queries, and n is the size of the proof encoding. The protocol is written in terms of an additive homomorphic encryption scheme, but as stated elsewhere [35, 45], the protocol can be modified to work with a multiplicative homomorphic scheme, such as ElGamal [23].

the circuit test; to generate the $\{\gamma_i\}$, V multiplies each constraint by a random value and collects like terms, a process described in [5, 13, 35, 45]. The completeness and soundness of this PCP are explained in those sources, and our notation is borrowed from [45]. Here we just assert that the soundness error of this PCP is $\epsilon = (7/9)^\rho$; that is, if the proof π is incorrect, the verifier detects that fact with probability greater than $1 - \epsilon$. To make ϵ small, PEPPER takes $\rho = 70$.

A.2 Stronger soundness analysis and consequences

GINGER retains the (P, V) argument system of PEPPER [45] but uses a modified PCP protocol (depicted in Figure 14) that makes the following changes to the base PCP protocol (Figure 13):

- Remove the linearity queries and tests.
- Set $\rho = 1$.

Theorem A.1. The (P, V) described above is an argument system with soundness $\epsilon_G \approx \sqrt[6]{1/|\mathbb{F}|}$. (The exact value of ϵ_G depends on intermediate lemmas and will be given at the end of the section.)

We will prove this theorem at the end of this section. To build up to the proof, we first strengthen the definition of a linear commitment primitive. We note that only the third property (linearity) in the definition is new; the rest is taken from [45, Appendix B], which itself heavily borrows framing, notation, and text from Ishai et al. [35].

Definition A.1 (Commitment to a function with multiple decommitments (CFMD)). Define a two-phase experiment between two probabilistic polynomial time ac-

tors (S, R) (a sender and receiver, which correspond to our prover and verifier) in an environment \mathcal{E} that generates \mathbb{F} , w and $Q = (q_1, \dots, q_\mu)$. In the first phase, the *commit phase*, S has w , and S and R interact, based on their random inputs. In the *decommit phase*, \mathcal{E} gives Q to R , and S and R interact again, based on further random inputs. At the end of this second phase, R outputs $A = (a_1, \dots, a_\mu) \in \mathbb{F}^\mu$ or \perp . A CFMD meets the following properties:

- **Correctness.** At the end of the decommit phase, R outputs $\pi(q_i) = \langle w, q_i \rangle$ (for all i), if S is honest.
- **ϵ_B -Binding.** Consider the following experiment. The environment \mathcal{E} produces two (possibly distinct) μ -tuples of queries: $Q = (q_1, \dots, q_\mu)$ and $\hat{Q} = (\hat{q}_1, \dots, \hat{q}_\mu)$. R and a cheating S^* run the commit phase once and two independent instances of the decommit phase. In the two instances R presents the queries as Q and \hat{Q} , respectively. We say that S^* *wins binding* if R 's outputs at the end of the respective decommit phases are $A = (a_1, \dots, a_\mu)$ and $\hat{A} = (\hat{a}_1, \dots, \hat{a}_\mu)$, and for some i, j , we have $q_i = \hat{q}_j$ but $a_i \neq \hat{a}_j$. We say that the protocol meets the ϵ_B -Binding property if for all \mathcal{E} and for all efficient S^* , the probability of S^* winning binding is less than ϵ_B . The probability is taken over three sets of independent randomness: the commit phase and the two runnings of the decommit phase.
- **ϵ_L -Linearity.** Consider the same experiment above. We say that S^* *wins linearity* if R 's outputs at the end of the respective decommit phases are $A = (a_1, \dots, a_\mu)$ and $\hat{A} = (\hat{a}_1, \dots, \hat{a}_\mu)$, and for some i, j, k , we have $\hat{q}_k = q_i + q_j$ but $\hat{a}_k \neq a_i + a_j$. We say that

The linear PCP from [5]

Loop ρ times:

- Generate linearity queries: Select $q_1, q_2 \in_R \mathbb{F}^s$ and $q_4, q_5 \in_R \mathbb{F}^{s^2}$. Take $q_3 \leftarrow q_1 + q_2$ and $q_6 \leftarrow q_4 + q_5$.
- Generate quadratic correction queries: Select $q_7, q_8 \in_R \mathbb{F}^s$ and $q_{10} \in_R \mathbb{F}^{s^2}$. Take $q_9 \leftarrow (q_7 \otimes q_8 + q_{10})$.
- Generate circuit queries: Select $q_{12} \in_R \mathbb{F}^s$ and $q_{14} \in_R \mathbb{F}^{s^2}$. Take $q_{11} \leftarrow \gamma_1 + q_{12}$ and $q_{13} \leftarrow \gamma_2 + q_{14}$.
- Issue queries. Send q_1, \dots, q_{14} to oracle π , getting back $\pi(q_1), \dots, \pi(q_{14})$.
- Linearity tests: Check that $\pi(q_1) + \pi(q_2) = \pi(q_3)$ and that $\pi(q_4) + \pi(q_5) = \pi(q_6)$. If not, **reject**.
- Quadratic correction test: Check that $\pi(q_7) \cdot \pi(q_8) = \pi(q_9) - \pi(q_{10})$. If not, **reject**.
- Circuit test: Check that $(\pi(q_{11}) - \pi(q_{12})) + (\pi(q_{13}) - \pi(q_{14})) = -\gamma_0$. If not, **reject**.

If V makes it here, **accept**.

Figure 13—The linear PCP that PEPPER uses. It is from [5]. The notation $x \otimes y$ refers to the outer product of two vectors x and y (meaning the vector or matrix consisting of all pairs of components from the two vectors). The values $\{\gamma_0, \gamma_1, \gamma_2\}$ are described briefly in the text.

the protocol meets the ϵ_L -linearity property if for all \mathcal{E} and for all efficient S^* , the probability of S^* winning linearity is less than ϵ_L . As with the prior property, the probability is taken over three sets of independent randomness: the commit phase and the two runnings of the decommit phase.

Prior work proved that Commit+Multidecommit (Figure 12) meets the first two properties above [45]. We will now show that it also meets the third property.

Lemma A.1. Commit+Multidecommit meets the definition of ϵ_L -linearity, with $\epsilon_L = 1/|\mathbb{F}| + \epsilon_S$, where ϵ_S comes from the semantic security of the homomorphic encryption scheme.

Proof. We will show that if S^* can systematically cheat, then an adversary \mathcal{A} could use S^* to break the semantic security of the encryption scheme.

Let $r \in_R \mathbb{F}^n$ and $Z_1, Z_2 \in_R \mathbb{F}$ (we use \in_R to mean “drawn uniformly at random from”). Semantic security (see [30], definitions 5.2.2, 5.2.8 and Exercise 17) implies that for all PPT \mathcal{A} (\mathcal{A} can be non-uniform),

$$\Pr_{\text{Gen, Enc, } r, Z_1, Z_2} \{ \mathcal{A}(pk, \text{Enc}(pk, r), r + Z_1 q, r + Z_2 q) = Z_1 \} < 1/|\mathbb{F}| + \epsilon_S. \quad (1)$$

This holds for all $q \in \mathbb{F}^n$.⁶

⁶We are being loose here. Under the actual definition of semantic security, (a) ϵ_S should be replaced with a negligible function of n , and (b) the claim holds only for n sufficiently large.

GINGER’s PCP protocol

- Generate quadratic correction queries: Select $q_1, q_2 \in_R \mathbb{F}^s$ and $q_4 \in_R \mathbb{F}^{s^2}$. Define $q_3 \leftarrow (q_1 \otimes q_2 + q_4)$. Note that q_3 will not travel, as P can derive it.
- Generate circuit queries: Take $q_5 \leftarrow \gamma_1 + q_1$. Take $q_6 \leftarrow \gamma_2 + q_4$.
- Issue queries. Send $(q_1, q_2, q_4, q_5, q_6)$ to oracle π , getting back $\pi(q_1), \pi(q_2), \pi(q_3), \pi(q_4), \pi(q_5), \pi(q_6)$.
- Quadratic correction test: Check that $\pi(q_1) \cdot \pi(q_2) = \pi(q_3) - \pi(q_4)$. If not, **reject**.
- Circuit test: Check that $(\pi(q_5) - \pi(q_1)) + (\pi(q_6) - \pi(q_4)) = -\gamma_0$. If so, **accept**.

Figure 14—GINGER’s PCP protocol, which refines PEPPER’s protocol (Figure 13). This protocol eliminates linearity testing and repetition, and recycles queries [9].

Now, assume to the contrary that Commit+Multidecommit does not meet the definition of ϵ_L -linearity. Then there exists an environment \mathcal{E} producing $q_i, q_j, i, j, k, Q, \hat{Q}, S^*$ (where Q has q_i, q_j in the i th and j th positions and \hat{Q} has $q_i + q_j$ in the k th position) such that $\Pr_{\text{all 3 phases}} \{ S^* \text{ wins linearity under } \mathcal{E} \} > 1/|\mathbb{F}| + \epsilon_S$. Let $q' \triangleq \hat{q}_k = q_i + q_j$.

We now describe an algorithm \mathcal{A} that, when given input $I = (pk, \text{Enc}(pk, r), r + Z_1 q', r + Z_2 q')$, can recover Z_1 with probability more than $1/|\mathbb{F}| + \epsilon_S$. \mathcal{A} has $Q, \hat{Q}, q_i, q_j, i, j, k$ hard-wired (because it is working under environment \mathcal{E}) and works as follows:

- (a) \mathcal{A} gives $(pk, \text{Enc}(pk, r))$ to S^* and ignores the reply.
- (b) \mathcal{A} randomly generates $\alpha_1, \dots, \alpha_\mu$ and sends to S^* the input $(Q, r + \alpha_1 q_1 + \dots + (\alpha_i + Z_1) q_i + \dots + (\alpha_j + Z_1) q_j + \dots + \alpha_\mu q_\mu)$. \mathcal{A} is able to construct this input because \mathcal{A} was given $r + Z_1 q' = r + Z_1 q_i + Z_1 q_j$. In response, S^* returns $(b, a_1, \dots, a_i, \dots, a_j, \dots, a_\mu)$.
- (c) \mathcal{A} randomly generates $\hat{\alpha}_1, \dots, \hat{\alpha}_\mu$. \mathcal{A} sends to S^* the input $(\hat{Q}, r + \hat{\alpha}_1 \hat{q}_1 + \dots + Z_2 \hat{q}_k + \dots + \hat{\alpha}_\mu \hat{q}_\mu)$. \mathcal{A} is able to construct this input because \mathcal{A} was given $r + Z_2 q' = r + Z_2 \hat{q}_k$. \mathcal{A} gets back $(\hat{b}, \hat{a}_1, \dots, \hat{a}_k, \dots, \hat{a}_\mu)$.

At this point, \mathcal{A} assumes that the responses from S^* pass the decommitment phase; that is, \mathcal{A} acts as if $b = s + \alpha_1 a_1 + \dots + (\alpha_i + Z_1) a_i + \dots + (\alpha_j + Z_1) a_j + \dots + \alpha_\mu a_\mu$ and $\hat{b} = s + \hat{\alpha}_1 \hat{a}_1 + \dots + Z_2 \hat{a}_k + \dots + \hat{\alpha}_\mu \hat{a}_\mu$. \mathcal{A} can write

$$K_1 = Z_2 \hat{a}_k - Z_1 (a_i + a_j), \quad (2)$$

where \mathcal{A} can derive $K_1 = \hat{b} - b - \sum_{\iota \neq k} \hat{\alpha}_\iota \hat{a}_\iota + \sum_\iota \alpha_\iota a_\iota$. Now, let $t = r + Z_1 q'$ and let $\hat{t} = r + Z_2 q'$ (both of these were supplied as input to \mathcal{A}). These two equations concern vectors. However, by choosing an index ι in the vector q' where q' is not zero (if the vector is zero everywhere, then r is revealed), \mathcal{A} can derive

$$K_2 = Z_2 - Z_1, \quad (3)$$

where $K_2 = (\hat{t}^{(\iota)} - t^{(\iota)})/q^{(\iota)}$.

Now, observe that if $\hat{a}_k \neq a_i + a_j$ (as happens when S^* wins), then \mathcal{A} can recover Z_1 by solving equations (2) and (3). Thus,

$$\begin{aligned} & \Pr_{\text{Gen, Enc, } r, Z_1, Z_2, \vec{\alpha}, \vec{\alpha}} \{\mathcal{A}(I) = Z_1\} \\ & \geq \Pr_{\text{Gen, Enc, } r, Z_1, Z_2, \vec{\alpha}, \vec{\alpha}} \{S^* \text{ wins linearity under } \mathcal{E}\} \\ & = \Pr_{\text{all 3 phases}} \{S^* \text{ wins linearity under } \mathcal{E}\} \\ & > 1/|\mathbb{F}| + \epsilon_S. \end{aligned} \quad (4)$$

The equality holds because the distribution of $(\alpha_1, \dots, \alpha_i + Z_1, \dots, \alpha_j + Z_1, \dots, \alpha_\mu)$ and $(\hat{\alpha}_1, \dots, Z_2, \dots, \hat{\alpha}_\mu)$ is equivalent to the distribution from which R selects in the decommit phases of the three-phase experiment, under Commit+Multidecommit. Meanwhile, inequality (4) contradicts inequality (1). \square

The lemmas ahead show that, under Commit+Multidecommit, S is bound to a *nearly linear* function, $\tilde{f}(\cdot)$; specifically, $\tilde{f}(\cdot)$ is δ^* -close to linear for small δ^* . By contrast, previous work [35, 45] showed only that S was bound to *some* function $\hat{f}(\cdot)$.

We now give some notation and restate two claims from [45]. Let ζ be the event that R 's output is a vector (a_1, \dots, a_μ) ; equivalently, ζ is the event that R 's output is non- \perp . Below, we sometimes write $\Pr_{\text{comm}}\{\cdot\}$ or $\Pr_{\text{decomm}}\{\cdot\}$ to mean the probability over the random choices of the commit or decommit phases.

Lemma A.2 (Existence of an extractor function [45]).

Let (S, R) be a CFMD protocol with binding error ϵ_B . Let $\epsilon_C = \mu \cdot 2 \cdot (2\sqrt[3]{9/2} + 1) \cdot \sqrt[3]{\epsilon_B}$. Let $v = (v_{S^*}, v_R)$ represent the views of S^* and R after the commit phase (v captures the randomness of the commit phase). For every efficient S^* and for every v , there exists a function $\tilde{f}_v : \mathbb{F}^\mu \rightarrow \mathbb{F}$ such that the following holds.⁷ For any environment \mathcal{E} , the output of R at the end of the decommit phase is, except with probability ϵ_C , either \perp or satisfies $a_i = \tilde{f}_v(q_i)$ for all $i \in [\mu]$, where (q_1, \dots, q_μ) are the decommitment queries generated by \mathcal{E} , and the probability is over the random inputs of S^* and R in both phases.

Lemma A.3. Let $\epsilon_3 = (2\sqrt[3]{9/2} + 1) \cdot \sqrt[3]{\epsilon_B}$. Label the i th query in Q as q_i and the i th response as a_i . For all Q, i , we have $\Pr_{\text{comm, decomm}} \{\zeta \cap \{a_i \neq \tilde{f}_v(q_i)\}\} < 2\epsilon_3$.

Proof. Follows from a claim in [45] (Claim B.4). \square

Lemma A.4. For all $q_1, q_2 \in \mathbb{F}^\mu$, $\Pr_{\text{comm}} \{\tilde{f}_v(q_1) + \tilde{f}_v(q_2) \neq \tilde{f}_v(q_1 + q_2)\} < \epsilon_F \triangleq \epsilon_L + 6\epsilon_3$.

⁷Note that after the commit phase, $\tilde{f}_v(\cdot)$ is deterministic. $\hat{f}_v(\cdot)$ is defined [35, 45] to map q to the value that R is most likely to successfully output in the decommit phase.)

Proof. Assume otherwise. Then for some q_1 and q_2 , we have $\Pr_{\text{comm}} \{\tilde{f}_v(q_1) + \tilde{f}_v(q_2) \neq \tilde{f}_v(q_1 + q_2)\} \geq \epsilon_F$, which implies $\Pr_{\text{all 3 phases}} \{\tilde{f}_v(q_1) + \tilde{f}_v(q_2) \neq \tilde{f}_v(q_1 + q_2)\} \geq \epsilon_F$, since we can “add coin flips that don’t matter”, namely those of the two decommit phases.

Now, consider the game in the definition of ϵ_L -linearity, and set $Q = (q_1, q_2, \dots)$ and $\hat{Q} = (q_1 + q_2, \dots)$. Let η be the event that S^* wins in this game. Let ν be the event that the outputs a_1, a_2, \hat{a}_1 are given by the function $\tilde{f}_v(\cdot)$. Then $\Pr_{\text{all 3 phases}} \{\neg \nu\} < 6\epsilon_3$, by Lemma A.3, by the union bound, and by (again) “adding coin flips that don’t matter” to get from a probability over two phases to one over three phases. Now, note that $\Pr_{\text{all 3 phases}} \{\eta | \nu\} \geq \epsilon_F$, by the contrary hypothesis. This implies that $\Pr_{\text{all 3 phases}} \{\eta\} \geq \epsilon_F - 6\epsilon_3 = \epsilon_L$, which contradicts the definition of ϵ_L -linearity. \square

Lemma A.4 almost talks about a linearity test [16]! But linearity testing theory [8] relates (a) the probability *over randomly chosen queries* that the test fails and (b) the closeness-to-linearity of the tested function. Thus, to apply the theory, we line up Lemma A.4 and (a).

Lemma A.5. With probability greater than $1 - \sqrt{\epsilon_F}$ over the commit phase, the fraction of (q_1, q_2) pairs that cause $\tilde{f}_v(\cdot)$ to fail the linearity test is $\leq \sqrt{\epsilon_F}$.

Proof. Let I_{v, q_1, q_2} be an indicator random variable that equals 1 if, in view v (that is, given the randomness of the commit phase), $\tilde{f}_v(q_1 + q_2) \neq \tilde{f}_v(q_1) + \tilde{f}_v(q_2)$. The lemma is equivalent to the statement

$$\Pr_{\text{comm}} \{ \Pr_{q_1, q_2} \{I_{v, q_1, q_2} = 1\} > \sqrt{\epsilon_F} \} < \sqrt{\epsilon_F}.$$

Now, define a random variable $Y_v = \frac{1}{Q} \sum_{q_1, q_2} I_{v, q_1, q_2}$, where $Q = |\mathbb{F}^\mu|$ is the number of possibilities for each of q_1 and q_2 . By linearity of expectation, $E_{\text{comm}}[Y_v] = \frac{1}{Q^2} \cdot (E[I_{v, 1}] + \dots + E[I_{v, Q}])$, where $E[I_{v, i}]$ is the probability, over the commit phase, that a particular (q_j, q_k) pair causes $\tilde{f}_v(\cdot)$ to fail the linearity test. Lemma A.4 implies that $E[I_{v, i}] < \epsilon_F$ for all i ; hence, $E_{\text{comm}}[Y_v] < \epsilon_F$. We now apply a Markov bound to Y_v :

$$\Pr_{\text{comm}} \{Y_v > \sqrt{\epsilon_F}\} < \frac{E_{\text{comm}}[Y_v]}{\sqrt{\epsilon_F}} < \frac{\epsilon_F}{\sqrt{\epsilon_F}} = \sqrt{\epsilon_F}.$$

But Y_v is equivalent to $\Pr_{q_1, q_2} \{I_{v, q_1, q_2} = 1\}$; making this substitution immediately above yields the lemma. \square

Lemma A.6. Let δ^* be the lesser root of $6\delta^2 - 3\delta + \sqrt{\epsilon_F} = 0$. If $\sqrt{\epsilon_F} < \frac{2}{9}$, then with probability greater than $1 - \sqrt{\epsilon_F}$ over the commit phase, $\tilde{f}_v(\cdot)$ is δ^* -close to linear.

Proof. We use the linearity testing results of Bellare et al. [8, 9] and the terminology of [8]. Define $\text{Dist}(f, g)$ to be the fraction of inputs on which f and g disagree.

Define $\text{Dist}(f)$ to be the fraction of inputs on which f disagrees with its “closest linear function” [8]. Define $\text{Rej}(f)$ to be the probability, over uniformly random choices of x and y from the domain of f , that $f(x)+f(y) \neq f(x+y)$; $\text{Rej}(f)$ is the probability that f fails the linearity test. As stated by Bellare et al. [8]:

- If $\text{Dist}(f) = \delta$, then $\text{Rej}(f) \geq 3\delta - 6\delta^2$.
- If $\text{Dist}(f) \geq \frac{1}{4}$, then $\text{Rej}(f) \geq \frac{2}{9}$.

The above implies the following claim: for all $\delta' \in \{\delta' \mid 3\delta' - 6\delta'^2 < \frac{2}{9} \text{ and } 0 \leq \delta' \leq \frac{1}{4}\}$, if $\text{Rej}(f) \leq 3\delta' - 6\delta'^2$, then $\text{Dist}(f) \leq \delta'$. (To see this, fix δ' . Assume to the contrary that $\delta = \text{Dist}(f) > \delta'$. There are two cases, and both contradict the given. If $\delta < \frac{1}{4}$, then $\text{Rej}(f) \geq 3\delta - 6\delta^2 > 3\delta' - 6\delta'^2$. If $\delta \geq \frac{1}{4}$, then $\text{Rej}(f) \geq \frac{2}{9} > 3\delta' - 6\delta'^2$.)

From lemma A.5, the probability is greater than $1 - \sqrt{\epsilon_F}$ over the commit phase that $\text{Rej}(\tilde{f}_v) \leq \sqrt{\epsilon_F}$. We call such commit phases *usual*. Under a *usual* commit phase, we can apply the claim just above. To do so, we assume that $\sqrt{\epsilon_F} < \frac{2}{9}$, and we set δ^* so that $\sqrt{\epsilon_F} = 3\delta^* - 6\delta^{*2}$ and $\delta^* \leq \frac{1}{4}$ (such a δ^* is guaranteed to exist because the parabola is symmetric about $\delta = \frac{1}{4}$). The claim implies that $\text{Dist}(\tilde{f}_v) \leq \delta^*$, or that \tilde{f}_v is δ^* -close to linear. \square

Lemma A.7. If the PCP oracle π is known to be δ^* -close to linear, then the linear PCP (Section A.1) with linearity testing removed has soundness error $\kappa > \max\{4\delta^* + \frac{2}{|\mathbb{F}|}, 4\delta^* + \frac{1}{|\mathbb{F}|}\}$.

Proof. This follows from the proof flow that establishes the soundness of linear PCPs, as in [5]. (A self-contained example is in Appendix D of [45].) Those proofs first establish that if the linearity test passes with probability higher than the soundness error, then π is δ -close to linear, for some δ . However, if we are *given* that π is δ^* -close to linear, then we can start those proofs midway and obtain the soundness of π as κ . \square

Proof of Theorem A.1. Lemma A.2 implies that there exists an extractor function that determines a (possibly incorrect) oracle $\tilde{\pi}$ such that, if V' does not reject during decommit, then with all but probability ϵ_C , V' receives back $\tilde{\pi}(q_1), \dots, \tilde{\pi}(q_\mu)$. We can thus “pay” probability ϵ_C in the union bound (below) to assume that V' hears back from $\tilde{\pi}$ itself. This allows us to apply Lemma A.6, at which point we can “pay” $\sqrt{\epsilon_F}$ more probability (again in the union bound below) to get that $\tilde{\pi}$ is δ^* -close to linear. (Applying the lemma requires that $\sqrt{\epsilon_F} < \frac{2}{9}$, and we will verify below that this bound holds.) Now, we can apply Lemma A.7 to ρ runs of the PCP protocol, giving a PCP soundness error of κ^ρ . Thus, the probability that V' wrongly accepts a proof is bounded from above by:

$$\epsilon_G = \epsilon_C + \sqrt{\epsilon_F} + \kappa^\rho.$$

By inspection (of the lemmas), the dominant contributor to ϵ_G , namely $\sqrt{\epsilon_F}$, is proportional to $\sqrt[6]{1/|\mathbb{F}|}$. \square

We compute a bound on ϵ_G as follows.

- ϵ_C is given in Lemma A.2. We take $\mu = 6$ (per Figure 14). We also take $\epsilon_B = 1/|\mathbb{F}|$ (following [45]; this amounts to ignoring the error from the semantic security of the homomorphic encryption scheme) and $|\mathbb{F}| = 2^{128}$, giving $\epsilon_C < 7.4 \cdot 10^{-12}$.
- $\epsilon_F = \epsilon_L + 6\epsilon_3$ (from Lemma A.4). ϵ_3 is given in Lemma A.3. We set $\epsilon_L = 1/|\mathbb{F}|$ (which again amounts to ignoring ϵ_S). Again taking $|\mathbb{F}| = 2^{128}$, we get $\sqrt{\epsilon_F} < 1.9 \cdot 10^{-6}$. Thus, $\sqrt{\epsilon_F} < 2/9$, as required.
- $\kappa = 4\delta^* + \frac{2}{|\mathbb{F}|}$, where δ^* is the lesser root of $6\delta^2 - 3\delta + \sqrt{\epsilon_F}$. This gives $\delta^* = 6.4 \cdot 10^{-7}$ and $\kappa = 2.6 \cdot 10^{-6}$.

Since κ and $\sqrt{\epsilon_F}$ are roughly the same, there is not much point to taking $\rho > 1$. Thus, we take $\rho = 1$, giving $\epsilon_G < 4.5 \cdot 10^{-6}$ when $|\mathbb{F}| = 2^{128}$. When $|\mathbb{F}| = 2^{192}$, we get $\epsilon_G < 2.8 \cdot 10^{-9}$.

A.3 Optimizing out queries

GINGER’s PCP protocol includes two further refinements. First, the protocol reuses q_4 and q_1 from test to test. This reuse is sound because the PCP soundness lemma [5] is of the form, “if all tests pass with probability greater than X , then the proof oracle π has a certain desired property”; meanwhile, as Bellare et al. [9] observe, the tests need not be independent! One can observe the savings by comparing Figure 13 (minus the linearity queries) to Figure 14. The protocol goes from 8 queries (the original 14 minus 6 linearity queries) to 6 queries, though the real savings for the prover is in reducing the 4 high-order queries (that is, queries to the $\mathbb{F}^{\mathbb{F}^2}$ component of π) to 3. Moreover, the verifier saves because it goes from generating pseudorandomness for 3 high-order queries (including γ_2) to 2. Second, V avoids transmitting a query (q_3) that P can generate for itself. This optimization offsets the consistency query, which is computed over \mathbb{Z} not \mathbb{Z}/p (owing to the details of our use of ElGamal [45, Appendix E]) and thus has roughly twice as many bits as a PCP query.