

Jolt: SNARKs for Virtual Machines via Lookups

Arasu Arun*

Srinath Setty[†]

Justin Thaler[‡]

Abstract

Succinct Non-interactive Arguments of Knowledge (SNARKs) allow an untrusted prover to establish that it correctly ran some “witness-checking procedure” on a witness. A zkVM (short for zero-knowledge Virtual Machine) is a SNARK that allows the witness-checking procedure to be specified as a computer program written in the assembly language of a specific instruction set architecture (ISA).

A *front-end* converts computer programs into a lower-level representation such as an arithmetic circuit or generalization thereof. A SNARK for circuit-satisfiability can then be applied to the resulting circuit.

We describe a new front-end technique called **Jolt** that applies to a variety of ISAs. **Jolt** arguably realizes a vision called the *lookup singularity*, which seeks to produce circuits that only perform lookups into pre-determined lookup tables. The circuits output by **Jolt** primarily perform lookups into a gigantic lookup table, of size more than 2^{128} , that depends only on the ISA. The validity of the lookups are proved via a new *lookup argument* called Lasso described in a companion work (Setty, Thaler, and Wahby, e-print 2023). Although size- 2^{128} tables are vastly too large to materialize in full, the tables arising in **Jolt** are structured, avoiding costs that grow linearly with the table size.

We describe performance and auditability benefits of **Jolt** compared to prior zkVMs, focusing on the popular RISC-V ISA as a concrete example. The dominant cost for the **Jolt** prover applied to this ISA (on 64-bit data types) is cryptographically committing to about six 256-bit field elements per step of the RISC-V CPU. This compares favorably to prior zkVM provers, even those focused on far simpler VMs.

1 Introduction

A SNARK (Succinct non-interactive argument of knowledge) is a cryptographic protocol that lets an untrusted prover \mathcal{P} prove to a verifier \mathcal{V} that they know a witness w satisfying some property. A trivial proof is for \mathcal{P} to explicitly send w to \mathcal{V} , who can then directly check on its own that w satisfies the claimed property. We refer to this trivial verification procedure as *direct witness checking*. A SNARK achieves the same effect, but with better costs to the verifier. Specifically, the term *succinct* roughly means that the proof should be shorter than this trivial proof (i.e., the witness w itself), and verifying the proof should be much faster than direct witness checking.

As an example, the prover could be a cloud service provider running an expensive computation on behalf of its client who acts as the verifier. A SNARK gives the verifier confidence that the prover ran the computation honestly. Alternatively, in a blockchain setting, the witness could be a list of valid digital signatures authorizing several blockchain transactions. A SNARK can be used to prove that one *knows* the signatures, so that the signatures themselves do not have to be stored and verified by all blockchain nodes. Instead, only the SNARK needs to be stored and verified on-chain.

1.1 SNARKs for Virtual Machine abstractions

A popular approach to SNARK design today is to prove the correct execution of *computer programs*. This means that the prover proves that it correctly ran a specified computer program Ψ on a witness. In the example above, Ψ might take as input a list of blockchain transactions and associated digital signatures authorizing each of them, and verify that each of the signatures is valid.

*New York University

[†]Microsoft Research

[‡]a16 crypto research and Georgetown University

Many projects today accomplish this via a CPU abstraction (in this context, also often called a *Virtual Machine (VM)*). Here, a VM abstraction entails fixing a set of *primitive instructions*, known as an instruction set architecture (ISA), analogous to assembly instructions in processor design. A full specification of the VM also includes the number of registers and the type of memory that is supported. The computer program that the prover proves it ran correctly must be specified in this assembly language.

To list a few examples, several so-called zkEVM projects seek to achieve “byte-code level compatibility” with the Ethereum Virtual Machine (EVM), which means that the set of primitive instructions is the 141 opcodes available on the EVM. Other zkEVMs do not aim for byte-code level compatibility, instead aiming to offer SNARKs for high-level smart contract languages such as Solidity (without first compiling the solidity to EVM bytecode).

Still other so-called zkVM projects take a similar approach but do not target the EVM instruction set, nor high-level languages like Solidity that are often compiled to EVM bytecode. These projects typically choose (or design) ISAs for their purported “SNARK-friendliness”, or for surrounding infrastructure and tooling, or for a combination thereof. For example, Cairo-VM is a very simple virtual machine designed specifically for compatibility with SNARK proving [GPR21, AGL⁺22]. The VM has 3 registers, memory that is read-only (each cell can only be written to once) and must be “continuous”, and the primitive instructions are roughly addition and multiplication over a finite field, jumps, and function calls.¹

Another example is the RISC Zero project, which uses the RISC-V instruction set. RISC-V is popular in the computer architecture community, and comes with a rich ecosystem of compiler tooling to transform higher-level programs into RISC-V assembly. Other zkVM projects include Polygon Miden,² Valida,³ and many others.

Front-end, back-end paradigm. SNARKs are built using protocols that perform certain probabilistic checks, so to apply SNARKs to program executions, one must express the execution of a program in a specific form that is amenable to probabilistic checking (e.g., as arithmetic circuits or generalizations thereof). Accordingly, most SNARKs consist of a so-called *front-end* and *back-end*: the front-end transforms a witness-checking computer program Ψ into an equivalent circuit-satisfiability instance, and the back-end allows the prover to establish that it knows a satisfying assignment to the circuit.

Typically, the circuit will “execute” each step of the compute program one at a time (with the help of untrusted “advice inputs”). Executing a step of the CPU conceptually involves two tasks: (1) identify which primitive instruction should be executed at this step, and (2) execute the instruction and update the CPU state appropriately. Existing front-ends implement these tasks by carefully devising gates or so-called constraints that implement each instruction. This is time-intensive and potentially error-prone. As we show in this work, it also leads to circuits that are substantially larger than necessary.

Pros and cons of the zkVM paradigm. One major benefit of zkVMs that use pre-existing ISAs is that they can exploit extant compiler infrastructure and tooling. This applies, for example, to the RISC-V and EVM instruction set, and leads to a developer-friendly toolchain without building the infrastructure from scratch. One can directly invoke existing compilers that transform witness-checking programs written in high-level languages down to assembly code for the ISA, and benefit from prior audits or other verification efforts of these compilers.

Another benefit of zkVMs is that a single circuit can suffice for running all programs up to a certain time bound, whereas alternative approaches may require re-running a front-end for every program (see the discussion in Section 1.6 of other front-end approaches). Finally, frontends for VM abstractions output circuits with repeated

¹The Cairo toolchain allows programmers to write programs in a higher-level language called Cairo 1.0, and these programs are compiled into primitive instructions for the Cairo-VM. Even the high-level language only exposes write-once (also known as immutable) memory to the programmer and does not offer signed integer data types. See <https://www.cairo-lang.org/> for information on the high-level language and [GPR21, AGL⁺22] and <https://github.com/lambdaclass/cairo-vm> for information on the Virtual Machine.

²<https://polygon.technology/polygon-miden>

³<https://github.com/valida-xyz/valida-compiler/issues/2>

structure. For a given circuit size, backends targeting circuits with repeated structure [Set20, BSBHR19, WTS⁺18] can be much faster than backends that do not leverage repeated structure [CHM⁺20, GWC19, Gro16].

However, zkVMs also have downsides that render them inappropriate for some applications. First, circuits implementing a VM abstraction are often much larger than circuits that do not. This means that zkVM provers are often much slower end-to-end than SNARK provers that do not impose upon themselves a VM abstraction.

For example, implementing certain important operations in a zkVM (e.g., cryptographic operations such as Keccak hashing or ECDSA signature verification) is extremely expensive—e.g., ECDSA signature verification takes up to 100 microseconds to verify on real CPUs, which translates to millions of RISC-V instructions.⁴ This is why zkVM projects contain so-called gadgets or built-ins, which are hand-optimized circuits and lookup tables computing specific functionalities.

A second downside is that, in order to expose a high-level programming language to developers, zkVMs require a compiler that transforms such high-level computer programs into assembly code for the VM. These compilers represent a large attack surface. Any bug in the compiler can render the system insecure: proving that one correctly ran assembly code does not guarantee knowledge of a valid witness if the assembly code fails to correctly implement the intended witness-checking procedure.

The conventional wisdom on zkVMs. The prevailing viewpoint today is that simpler VMs can be turned into circuits with fewer gates per step of the VM. This is most apparent in the design of particularly simple and ostensibly SNARK-friendly VMs such as the Cairo-VM. However, this comes at a cost, because primitive operations that are standard in real-world CPUs require many primitive instructions to implement on the simple VM.

In part to minimize the overheads in implementing standard operations on such limited VMs, many projects have designed domain specific languages (DSLs) that are exposed to the programmer who writes the witness-checking program. The proliferation of DSLs places a burden on the programmer, who is responsible both for learning the DSL and writing correct programs in it (with catastrophic security consequences if a program is incorrect).

Moreover, existing zkVMs remain expensive for the prover, even for very simple ISAs. For example, the prover for Cairo-VM programs described in [GPR21, AGL⁺22] cryptographically commits to 51 field elements per step of the Cairo-VM. This means that a single primitive instruction for the Cairo-VM may cause the prover to execute millions of instructions on real CPUs. This severely limits the applicability of SNARKs for VM abstractions, to applications involving only very simple witness-checking procedures.

1.2 Jolt: A new paradigm for zkVM design

In this work, we introduce a new paradigm in zkVM design. The result is zkVMs with much faster provers, as well as substantially improved auditability and extensibility (i.e., a simple workflow for adding additional primitive instructions to the VM). Our techniques are general. As a concrete example, we instantiate them for the RISC-V instruction set (with multiplication extension [And17]), a popular open-source ISA developed by the computer architecture community without SNARKs in mind.

Our results upend the conventional wisdom that simpler instruction sets necessarily lead to smaller circuits and associated faster provers. First, our prover is faster per step of the VM than existing SNARK provers for much simpler VMs. Second, the complexity of our prover primarily depends on the size (i.e., number of bits) of the inputs to each instruction. This holds so long as all of the primitive instructions satisfy a natural notion of structure, called *decomposability*. Roughly speaking, decomposability means that one can evaluate the instruction on a given pair of inputs (x, y) by breaking x and y up into smaller chunks, evaluating a small number of functions of each chunk, and combining the results. A primary contribution of our work is to show that decomposability is satisfied by all instructions in the RISC-V instruction set.

⁴See <https://github.com/risc0/risc0/tree/v0.16.0/examples/ecdsa>.

Lookup arguments and Lasso. In a lookup argument, there is a predetermined “table” T of size N , meaning that $T \in \mathbb{F}^N$. An (*unindexed*) lookup argument allows the prover to commit to any vector $a \in \mathbb{F}^m$ and prove that every entry of a resides somewhere in the table. That is, for every $i \in \{1, \dots, m\}$, there exists some k such that $a_i = T[k]$. In an *indexed* lookup argument, the prover commits not only to $a \in \mathbb{F}^m$, but also a vector $b \in \mathbb{F}^m$, and the prover proves that for every i , $a_i = T[b_i]$. In this setting, we call a the vector of *lookups* and b the vector of associated *indices*.

In a companion paper, we describe a new lookup argument called **Lasso** (which applies to both indexed and unindexed lookups). One distinguishing feature of **Lasso** is that it applies even to tables that are far too large for anyone to materialize in full, so long as the table satisfies the *decomposability* condition mentioned earlier.

Jolt. Say \mathcal{P} claims to have run a certain computer program for m steps, and that the program is written in the assembly language for a VM. Today, front-ends produce a circuit that, for each step of the computation: (1) identifies what instruction to execute at that step, and then (2) executes that instruction.

Lasso lets one replace Step 2 with a single lookup. For each instruction, the table stores the entire evaluation table of the instruction. If instruction f operations on two 64-bit inputs, the table stores $f(x, y)$ for every pair of inputs $(x, y) \in \{0, 1\}^{64} \times \{0, 1\}^{64}$. This table has size 2^{128} . In this work, we show that all RISC-V instructions are decomposable.

1.3 Costs of Jolt

1.3.1 Background and context

Polynomial commitments and MSMs. A central component of most SNARKs is a cryptographic protocol called a *polynomial commitment scheme* (see Definition 2.2). Such a scheme allows an untrusted prover to succinctly commit to a polynomial p and later reveal an evaluation $p(r)$ for a point r chosen by the verifier (the prover will also return a *proof* that the claimed evaluation is indeed equal to the committed polynomial’s evaluation at r). In **Jolt**, as with most SNARKs, the bottleneck for the prover is the polynomial commitment scheme.

Many popular polynomial commitments are based on multi-exponentiations (also known as multi-scalar multiplications, or MSMs). This means that the commitment to a polynomial p (with n coefficients c_0, \dots, c_{n-1} over an appropriate basis) is

$$\prod_{i=0}^{n-1} g_i^{c_i},$$

for some public generators g_1, \dots, g_n of a multiplicative group \mathbb{G} . Examples include KZG [KZG10], Bulletproofs/IPA [BCC⁺16, BBB⁺18], Hyrax [WTS⁺18], and Dory [Lee21].⁵

The naive MSM algorithm performs n group exponentiations and n group multiplications (note that each group exponentiation is about $400\times$ slower than a group multiplication). But Pippenger’s MSM algorithm saves a factor of about $\log(n)$ relative to the naive algorithm. This factor can be well over $10\times$ in practice.

Working over large fields, but committing to small elements. If all exponents appearing in the multi-exponentiation are “small”, one can save another factor of $10\times$ relative to applying Pippenger’s algorithm to an MSM involving random exponents. This is analogous to how computing $g_i^{2^{16}}$ is $10\times$ faster than computing $g_i^{2^{160}}$: the first requires 16 squaring operations, while the second requires 160 such operations.

In other words, if one is promised that all field elements (i.e., exponents) to be committed via an MSM are in the set $\{0, 1, \dots, K\} \subset \mathbb{F}$, the number of group operations required to compute the MSM depend only on K and not on the size of \mathbb{F} .⁶

⁵In Hyrax and Dory, the prover does \sqrt{n} MSMs each of size \sqrt{n} .

⁶Of course, the cost of each group operation depends on the size of the group’s base field, which is closely related to that of the scalar field \mathbb{F} . However, the *number* of group operations to compute the MSM depends only on K , not on \mathbb{F} .

Quantitatively, if all exponents are upper bounded by some value K , with $K \ll n$, then Pippenger’s algorithm only needs (about) one group *operation* per term in the multi-exponentiation. More generally, with any MSM-based commitment scheme, Pippenger’s algorithm allows the prover to commit to roughly $k \cdot \log(n)$ -bit field elements (meaning field elements in $\{0, 1, \dots, n\}$) with only k group *operations* per committed field element. So for size- n MSMs, one can commit to $\log(n)$ bits with a *single* group operation.

Polynomial evaluation proofs. In any SNARK, the prover not only has to commit to one or more polynomials, but also reveal to the verifier an evaluation of the committed polynomials at a point of the verifier’s choosing. This requires the prover to compute a so-called evaluation proof, which establishes that the returned evaluation is indeed consistent with the committed polynomial. For some polynomial commitment schemes, such as Bulletproofs/IPA [BCC⁺16, BBB⁺18], evaluation proofs are quite slow and this cost can bottleneck the prover. However, for others, evaluation proof computation is a low-order cost [WTS⁺18, BBHR18, Lee21]. Moreover, evaluation proofs exhibit excellent batching properties, whereby the prover can commit to many polynomials and only produce a single evaluation proof across all of them [BGH19, Lee21, KST22, BDFG20]. So in many contexts, computing opening proofs is not a bottleneck even when a scheme such as Bulletproofs/IPA is employed. For these reasons, our accounting in this work ignores the cost of polynomial evaluation proofs.

1.3.2 Costs of Jolt

Prover costs. For RISC-V instructions on 64-bit data types (with the multiply extension), Jolt’s \mathcal{P} commits to under 60 field elements per step of the RISC-V CPU. Only six of those field elements are larger than 2^{25} , and none of them are larger than 2^{64} . With MSM-based polynomial commitment, the Jolt prover costs are roughly that of committing to 6 arbitrary (256-bit) field elements per CPU step.

One caveat is that we handle six RISC-V instructions (all in the multiplication extension) via several “pseudoinstructions”. For example, we handle the division with remainder instruction by having \mathcal{P} provide the quotient and remainder as untrusted advice, and they are checked for correctness by applying multiplication and addition instructions. Another caveat is that some load and store instructions have modestly higher costs than those listed above. Conversely, many instructions (those involving addition, subtraction, shifts, jumps, loads, and stores) can be handled with *fewer than* five committed 256-bit field elements.

Comparison of prover costs to prior works. A detailed experimental comparison of Jolt to existing zkVMs will have to wait until a full implementation is complete, but some crude comparisons to prior works are illustrative. Recall that, when using an MSM-based multilinear polynomial commitment scheme (such as multilinear analogs of KZG, like Zeromorph [KT23]) we estimate the cost of the Jolt prover as being roughly that of committing to five arbitrary 256-bit field elements per step of the RISC-V CPU.

Plonk [GWC19] is a popular backend that can prove statements about certain generalizations of arithmetic circuit satisfiability. When Plonk is applied to an arithmetic circuit (i.e., consisting of addition and multiplication gates of fan-in two), the Plonk prover commits to 11 field elements per gate of the circuit, and 7 of these 11 field elements are random. Thus, the Jolt prover costs are roughly equivalent to applying the Plonk backend to an arithmetic circuit with only about one gate per step of the RISC-V CPU.

A more apt comparison is to the RISC Zero project⁷, which currently targets the RISC-V ISA on 32-bit data types (with the multiplication extension). A direct comparison is complicated, in part because RISC Zero uses FRI as its (univariate) polynomial commitment scheme, which is based on FFTs and Merkle-hashing, avoiding the use of elliptic curve groups. Jolt can use related polynomial commitment schemes (Jolt can use any commitment scheme for multilinear polynomials). However, we choose to focus on elliptic-curve-based schemes, because Jolt’s property of having the prover commit only to elements in $\{0, \dots, b\}$ for some $b \ll |\mathbb{F}|$ benefits those commitment schemes more than hashing-based ones.⁸ Still, a crude comparison can be made by comparing how many field elements the RISC Zero prover commits to, vs. the Jolt prover.

⁷<https://www.risczero.com/>

⁸This property would also benefit hashing-based commitment schemes that operate over an extension field of a relatively small base field, owing to all committed elements in Lasso being in the base field.

The RISC Zero prover commits to at least 275 31-bit field elements per CPU step [Tom23]. This is roughly equivalent to committing to about $275 \cdot 32/256 \approx 34$ different 256-bit field elements per CPU step: at least on small instances, the prover bottleneck is Merkle-hashing the result of various FFTs [Tom23], and one can hash 8 different 31-bit field elements with the same cost as hashing one 256-bit field element.

A final comparison point is to the SNARK for the Cairo-VM described in the Cairo whitepaper [GPR21]. The prover in that SNARK commits to about 50 field elements per step of the Cairo Virtual Machine, using FRI as the polynomial commitment scheme. StarkWare currently works over a 251-bit field.⁹ This field size may be larger than necessary (it is chosen to match the field used by certain ECDSA signatures), but the provided arithmetization of Cairo-VM *requires* a field of size at least 2^{63} . So the commitment costs for the prover are at least equivalent to committing to $50 \cdot 64/256 \approx 13$ 256-bit field elements.¹⁰ Jolt’s prover costs per CPU compare favorably to this, despite the RISC-V instruction set being vastly more complicated than the Cairo-VM (and with the Cairo-VM instruction set specifically designed to be ostensibly “SNARK-friendly”).

Verifier costs of Jolt. For RISC-V programs running for at most T steps, the dominant costs for the Jolt verifier are performing $O(\log(T) \log \log(T))$ hash evaluations and field operations,¹¹ plus checking one evaluation proof from the chosen polynomial commitment scheme (when applied to a multilinear polynomial over at most $O(\log T)$ variables).

Verifier costs can be further reduced, and the SNARK rendered zero-knowledge, via composition with a zero-knowledge SNARK with smaller proof size. For example, see the recent work Testudo for a related approach (Testudo instantiates Spartan [Set20] with a variant of PST polynomial commitments [PST13] (an analog of KZG commitments [KZG10] for multilinear rather than univariate polynomials) and composes this with Groth16 [Gro16]).

1.4 The lookup singularity

In a research forum post in 2022, Barry Whitehat articulated a goal of designing front-ends that produce circuits that *only* perform lookups [Whi]. Whitehat terms this the *lookup singularity* and sketches how achieving this would help address a key issue (the potential for security bugs, and difficulty of auditability) that must be addressed for long-term and large-scale adoption of SNARKs. Circuits that only perform lookups (and the lookup arguments that enable them) should be much simpler to understand and formally verify than circuits consisting of many gates that are often hand-optimized.

Whitehat’s post acknowledges that current lookup arguments are expensive, but predicts that lookup arguments will get more performative with time. Arguably, Jolt realizes the vision of the lookup singularity. The bulk of the prover work in Jolt lies in the lookup argument, Lasso. The Jolt front-end does output some constraints that effectively implement the task of the RISC-V CPU figuring out, at each step of the computation, which instruction to execute. These constraints are simple and easily captured in R1CS.

1.5 Technical details: CPU instructions as structured polynomials

Lasso is most efficient when applied to lookup tables satisfying a property called *decomposability*. Intuitively, this refers to tables t such that one lookup into t of size N can be answered with a small number (say, about c) of lookups into much smaller tables t_1, \dots, t_ℓ , each of size $N^{1/c}$. Furthermore, if a certain polynomial \tilde{t}_i associated with each t_i can be evaluated at any desired point r using, say, $O(\log(N)/c)$ field operations,¹² then no one needs to cryptographically commit to any of the tables (neither to t itself, nor to t_1, \dots, t_ℓ).

⁹See, for example, https://github.com/starkware-libs/starkex-contracts/blob/master/audit/EVM_STARK_Verifier_v4.0_Audit_Report.pdf.

¹⁰Furthermore, in order to control proof size, StarkWare currently uses a “FRI blowup factor” of 16, compared to RISC Zero’s choice of 4. This adds at least an extra factor of 4 to the prover time per field element committed, relative to RISC Zero’s.

¹¹As described in Appendix B.3, Lasso can use any so-called *grand product argument*. The $O(\log(T) \log \log(T))$ verifier cost are due to the choice of grand product argument from [SL20, Section 6]. Other choices of lookup argument offer different tradeoffs between commitment costs for the prover, versus proof size and verifier time.

¹²The Lasso verifier has to evaluate \tilde{t}_i at a random point r on its own, so we need this computation to be fast enough that we are satisfied with the resulting verifier runtime. For all tables arising in Jolt, the verifier can compute all necessary \tilde{t}_i polynomial evaluations in $O(\log(N))$ total field operations.

Specifically, \tilde{t}_i can be any so-called *low-degree extension* polynomial of t_i . In Jolt, we will exclusively work with a specific low-degree extension of t_i , called the *multilinear extension*, and denoted \tilde{t}_i .

Hence, to take full advantage of Lasso, we must show two things:

- The evaluation table t of each RISC-V instruction has is decomposable in the above sense. That is, one lookup into t , which has size N , can be answered with a small number of lookups into much smaller tables t_1, \dots, t_ℓ , each of size $N^{1/c}$. For most RISC-V instructions, ℓ equals one or two, and about c lookups are performed into each table.
- For each of the small tables t_i , the multilinear extension \tilde{t}_i is evalutable at any point, using just $O(\log(N)/c)$ field operations.

Establishing the above is the main technical contribution of our work. It turns out to be quite straightforward for certain instructions (e.g., bitwise AND), but more complicated for others (e.g., bitwise shifts, comparisons).

Decomposable instructions. Suppose that table t contains all evaluations of some primitive instruction $f: \{0, 1\}^n \rightarrow \mathbb{F}$. Decomposability of the table t is equivalent to the following property of f : for any n -bit input x to f , x can be decomposed into c “chunks”, X_0, \dots, X_{c-1} , each of size n/c , and such that there following holds. There are ℓ functions $f_0, \dots, f_{\ell-1}$ such that $f(x)$ can be derived in a relatively simple manner from $f_i(x_j)$ as i ranges over $0, \dots, \ell-1$ and j ranges over $0, \dots, c-1$. Then the evaluation table t of f is decomposable: one lookup into t can be answered with c total lookups into $\ell \cdot c$ lookups into the evaluation tables of $f_0, \dots, f_{\ell-1}$.

Bitwise AND is a clean example by which to convey intuition for why the evaluation tables of RISC-V instructions are decomposable. Suppose we have two field elements a and b in \mathbb{F} , both in $\{0, \dots, 2^{64} - 1\}$. We refer to a and b as 64-bit field elements (we clarify here that “64 bits” does *not* refer to the size of the field \mathbb{F} , which may, for example, be a 256-bit field. Rather to the fact that a and b are both in the much smaller set $\{0, \dots, 2^{64} - 1\} \subset \mathbb{F}$, no matter how large \mathbb{F} may be).

Our goal is to determine the 64-bit field element c whose binary representation is given by the bitwise AND of the binary representations of a and b . That is, if $a = \sum_{i=0}^{63} 2^i \cdot a_i$ and $b = \sum_{i=0}^{63} 2^i \cdot b_i$ for $(a_0, \dots, a_{63}) \in \{0, 1\}^{64}$ and $(b_0, \dots, b_{63}) \in \{0, 1\}^{64}$, then $c = \sum_{i=0}^{63} 2^i \cdot a_i \cdot b_i$.

One way to compute c is as follows. Break a and b into 8 chunks of 8 bits each compute the bitwise AND of each chunk, and concatenate the results to obtain c . Equivalently, we can express

$$c = \sum_{i=0}^7 2^{8 \cdot i} \cdot \text{AND}(a'_i, b'_i), \quad (1)$$

where each $a'_i, b'_i \in \{0, \dots, 2^8 - 1\}$ is such that $a = \sum_{i=0}^7 2^{8 \cdot i} \cdot a'_i$ and $b = \sum_{i=0}^7 2^{8 \cdot i} \cdot b'_i$. These a'_i 's and b'_i 's represent the decomposition of a and b into 8-bit limbs.¹³

In this way, one lookup into the evaluation table of bitwise-AND, which has size 2^{128} , can be answered by the prover providing $a'_1, \dots, a'_8, b'_1, \dots, b'_8 \in \{0, \dots, 2^8 - 1\}$ as untrusted advice, and performing 8 lookups into the size- 2^{16} table t_1 containing all evaluations of bitwise-AND over pairs of 8-bit inputs. The results of these 8 lookups can easily be collated into the result of the original lookup, via Equation (1). No party has to commit to the size- 2^{16} table t_1 because for any input $(r'_0, \dots, r'_7, r''_0, \dots, r''_7) \in \mathbb{F}^{16}$,

$$\tilde{t}_1(r'_0, \dots, r'_7, r''_0, \dots, r''_7) = \sum_{i=0}^{15} 2^i \cdot r'_i \cdot r''_i,$$

which can be evaluated directly by the verifier with only 32 field operations.

¹³Just as “digits” refers a base-10 decomposition of an integer or field element, “limbs” refer to a decomposition into a different base, in this case base 8.

Challenges for other instructions. One may initially expect that correct execution of RISC-V operations capturing 64-bit addition and multiplication would be easy prove, because large prime-order fields come with addition and multiplication operations that behave like integer addition and multiplication until the result of the operation overflows the field characteristic. Unfortunately, the RISC-V instructions capturing addition and multiplication have specified behavior upon overflow that differs from that of field addition and multiplication. Resolving this discrepancy is one key challenge that we overcome.

1.6 Other front-end approaches

As with other zkVM projects, Jolt produces a so-called *universal circuit*, meaning one circuit works for all RISC-V programs running up to some time bound T . This has the benefit that the circuit-generation process only needs to be run once.

Other front-end approaches do not implement a Virtual Machine abstraction (i.e., they do not produce circuits that repeatedly execute the transition function of a specific ISA). These approaches typically output a different circuit for every computer program, such as Buffet [WSR⁺15], Bellman, Circom, Zokrates, Noir, etc. The circuits produced by these approaches can also be made smaller and proved faster using our techniques, though we leave this to future work.

2 Technical Preliminaries

2.1 Multilinear extensions

An ℓ -variate polynomial $p: \mathbb{F}^\ell \rightarrow \mathbb{F}$ is said to be *multilinear* if p has degree at most one in each variable. Let $f: \{0, 1\}^\ell \rightarrow \mathbb{F}$ be any function mapping the ℓ -dimensional Boolean hypercube to a field \mathbb{F} . A polynomial $g: \mathbb{F}^\ell \rightarrow \mathbb{F}$ is said to *extend* f if $g(x) = f(x)$ for all $x \in \{0, 1\}^\ell$. It is well-known that for any $f: \{0, 1\}^\ell \rightarrow \mathbb{F}$, there is a unique *multilinear* polynomial $\tilde{f}: \mathbb{F}^\ell \rightarrow \mathbb{F}$ that extends f . The polynomial \tilde{f} is referred to as the *multilinear extension* (MLE) of f .

Multilinear extensions of vectors. Given a vector $u \in \mathbb{F}^m$, we will often refer to the *multilinear extension of u* and denote this multilinear polynomial by \tilde{u} . \tilde{u} is obtained by viewing u as a function mapping $\{0, 1\}^{\log m} \rightarrow \mathbb{F}$ in the natural way¹⁴: the function interprets its $(\log m)$ -bit input $(i_0, \dots, i_{\log m-1})$ as the binary representation of an integer i between 0 and $m-1$, and outputs u_i . \tilde{u} is defined to be the multilinear extension of this function.

Lagrange interpolation. An explicit expression for the MLE of any function is given by the following standard lemma (see [Tha22, Lemma 3.6]).

Lemma 1. *Let $f: \{0, 1\}^\ell \rightarrow \mathbb{F}$ be any function. Then the following multilinear polynomial \tilde{f} extends f :*

$$\tilde{f}(x_0, \dots, x_{\ell-1}) = \sum_{w \in \{0, 1\}^\ell} f(w) \cdot \chi_w(x_0, \dots, x_{\ell-1}), \quad (2)$$

where, for any $w = (w_0, \dots, w_{\ell-1})$, $\chi_w(x_0, \dots, x_{\ell-1}) := \prod_{i=0}^{\ell-1} (x_i w_i + (1 - x_i)(1 - w_i))$. Equivalently,

$$\chi_w(x_0, \dots, x_{\ell-1}) = \widetilde{\mathbf{EQ}}(x_0, \dots, x_{\ell-1}, w_0, \dots, w_{\ell-1}).$$

The polynomials $\{\chi_w: w \in \{0, 1\}^\ell\}$ are called the *Lagrange basis polynomials* for ℓ -variate multilinear polynomials. The evaluations $\{\tilde{f}(w): w \in \{0, 1\}^\ell\}$ are sometimes called the coefficients of \tilde{f} in the *Lagrange basis*, terminology that is justified by Equation (2).

¹⁴All logarithms in this paper are to base 2.

SNARKs. We adapt the definition provided in [KST22].

Definition 2.1. Consider a relation \mathcal{R} over public parameters, structure, instance, and witness tuples. A non-interactive argument of knowledge for \mathcal{R} consists of PPT algorithms $(\mathcal{G}, \mathcal{P}, \mathcal{V})$ and deterministic \mathcal{K} , denoting the generator, the prover, the verifier and the encoder respectively with the following interface.

- $\mathcal{G}(1^\lambda) \rightarrow \text{pp}$: On input security parameter λ , samples public parameters pp .
- $\mathcal{K}(\text{pp}, \mathbf{s}) \rightarrow (pk, vk)$: On input structure \mathbf{s} , representing common structure among instances, outputs the prover key pk and verifier key vk .
- $\mathcal{P}(pk, u, w) \rightarrow \pi$: On input instance u and witness w , outputs a proof π proving that $(\text{pp}, \mathbf{s}, u, w) \in \mathcal{R}$.
- $\mathcal{V}(vk, u, \pi) \rightarrow \{0, 1\}$: On input the verifier key vk , instance u , and a proof π , outputs 1 if the instance is accepting and 0 otherwise.

A non-interactive argument of knowledge satisfies completeness if for any PPT adversary \mathcal{A}

$$\Pr \left[\mathcal{V}(vk, u, \pi) = 1 \mid \begin{array}{l} \text{pp} \leftarrow \mathcal{G}(1^\lambda), \\ (\mathbf{s}, (u, w)) \leftarrow \mathcal{A}(\text{pp}), \\ (\text{pp}, \mathbf{s}, u, w) \in \mathcal{R}, \\ (pk, vk) \leftarrow \mathcal{K}(\text{pp}, \mathbf{s}), \\ \pi \leftarrow \mathcal{P}(pk, u, w) \end{array} \right] = 1.$$

A non-interactive argument of knowledge satisfies knowledge soundness if for all PPT adversaries \mathcal{A} there exists a PPT extractor \mathcal{E} such that for all randomness ρ

$$\Pr \left[\begin{array}{l} \mathcal{V}(vk, u, \pi) = 1, \\ (\text{pp}, \mathbf{s}, u, w) \notin \mathcal{R} \end{array} \mid \begin{array}{l} \text{pp} \leftarrow \mathcal{G}(1^\lambda), \\ (\mathbf{s}, u, \pi) \leftarrow \mathcal{A}(\text{pp}; \rho), \\ (pk, vk) \leftarrow \mathcal{K}(\text{pp}, \mathbf{s}), \\ w \leftarrow \mathcal{E}(\text{pp}, \rho) \end{array} \right] = \text{negl}(\lambda).$$

A non-interactive argument of knowledge is succinct if the verifier's time to check the proof π and the size of the proof π are at most polylogarithmic in the size of the statement proven.

Polynomial commitment schemes. We adapt the definition from [BFS20]. A polynomial commitment scheme for multilinear polynomials is a tuple of four protocols $\text{PC} = (\text{Gen}, \text{Commit}, \text{Open}, \text{Eval})$:

- $\text{pp} \leftarrow \text{Gen}(1^\lambda, \ell)$: takes as input ℓ (the number of variables in a multilinear polynomial); produces public parameters pp .
- $\mathcal{C} \leftarrow \text{Commit}(\text{pp}, \mathcal{G})$: takes as input a ℓ -variate multilinear polynomial over a finite field $\mathcal{G} \in \mathbb{F}[\ell]$; produces a commitment \mathcal{C} .
- $b \leftarrow \text{Open}(\text{pp}, \mathcal{C}, \mathcal{G})$: verifies the opening of commitment \mathcal{C} to the ℓ -variate multilinear polynomial $\mathcal{G} \in \mathbb{F}[\ell]$; outputs $b \in \{0, 1\}$.
- $b \leftarrow \text{Eval}(\text{pp}, \mathcal{C}, r, v, \ell, \mathcal{G})$ is a protocol between a PPT prover \mathcal{P} and verifier \mathcal{V} . Both \mathcal{V} and \mathcal{P} hold a commitment \mathcal{C} , the number of variables ℓ , a scalar $v \in \mathbb{F}$, and $r \in \mathbb{F}^\ell$. \mathcal{P} additionally knows a ℓ -variate multilinear polynomial $\mathcal{G} \in \mathbb{F}[\ell]$. \mathcal{P} attempts to convince \mathcal{V} that $\mathcal{G}(r) = v$. At the end of the protocol, \mathcal{V} outputs $b \in \{0, 1\}$.

Definition 2.2. A tuple of four protocols $(\text{Gen}, \text{Commit}, \text{Open}, \text{Eval})$ is an extractable polynomial commitment scheme for multilinear polynomials over a finite field \mathbb{F} if the following conditions hold.

- **Completeness.** For any ℓ -variate multilinear polynomial $\mathcal{G} \in \mathbb{F}[\ell]$,

$$\Pr \left\{ \begin{array}{l} \text{pp} \leftarrow \text{Gen}(1^\lambda, \ell); \mathcal{C} \leftarrow \text{Commit}(\text{pp}, \mathcal{G}); \\ \text{Eval}(\text{pp}, \mathcal{C}, r, v, \ell, \mathcal{G}) = 1 \wedge v = \mathcal{G}(r) \end{array} \right\} \geq 1 - \text{negl}(\lambda)$$

- **Binding.** For any PPT adversary \mathcal{A} , size parameter $\ell \geq 1$,

$$\Pr \left\{ \begin{array}{l} \text{pp} \leftarrow \text{Gen}(1^\lambda, \ell); (\mathcal{C}, \mathcal{G}_0, \mathcal{G}_1) = \mathcal{A}(\text{pp}); \\ b_0 \leftarrow \text{Open}(\text{pp}, \mathcal{C}, \mathcal{G}_0); b_1 \leftarrow \text{Open}(\text{pp}, \mathcal{C}, \mathcal{G}_1): \\ b_0 = b_1 \neq 0 \wedge \mathcal{G}_0 \neq \mathcal{G}_1 \end{array} \right\} \leq \text{negl}(\lambda)$$

- **Knowledge soundness.** Eval is a succinct argument of knowledge for the following NP relation given $\text{pp} \leftarrow \text{Gen}(1^\lambda, \ell)$:

$$\mathcal{R}_{\text{Eval}}(\text{pp}) = \{ \langle (\mathcal{C}, r, v), (\mathcal{G}) \rangle : \mathcal{G} \in \mathbb{F}[\mu] \wedge \mathcal{G}(r) = v \wedge \text{Open}(\text{pp}, \mathcal{C}, \mathcal{G}) = 1 \}.$$

2.2 Polynomial IOPs and polynomial commitments

Modern SNARKs are constructed by combining a type of interactive protocol called a *polynomial IOP* [BFS20] with a cryptographic primitive called a *polynomial commitment scheme* [KZG10]. The combination yields a succinct *interactive* argument, which can then be rendered non-interactive via the Fiat-Shamir transformation [FS86], yielding a SNARK.

Roughly, a polynomial IOP is an interactive protocol where, in one or more rounds, the prover may “send” to the verifier a very large polynomial g . Because g is so large, one does not wish for the verifier to read a complete description of g . Instead, in any efficient polynomial IOP, the verifier only “queries” g at one point (or a handful of points). This means that the only information the verifier needs about g to check that the prover is behaving honestly is one (or a few) evaluations of g .

In turn, a polynomial commitment scheme enables an untrusted prover to succinctly *commit* to a polynomial g , and later provide to the verifier any evaluation $g(r)$ for a point r chosen by the verifier, along with a proof that the returned value is indeed consistent with the committed polynomial. Essentially, a polynomial commitment scheme is exactly the cryptographic primitive that one needs to obtain a succinct argument from a polynomial IOP. Rather than having the prover send a large polynomial g to the verifier as in the polynomial IOP, the argument system prover instead cryptographically commits to g and later reveals any evaluations of g required by the verifier to perform its checks.

Whether or not a SNARK requires a trusted setup, as well as whether or not it is plausibly post-quantum secure, is determined by the polynomial commitment scheme used. If the polynomial commitment scheme does not require a trusted setup, neither does the resulting SNARK, and similarly if the polynomial commitment scheme is plausibly binding against quantum adversaries, then the SNARK is plausibly post-quantum sound.

Lasso can make use of any commitment schemes for *multilinear* polynomials g .¹⁵ Here an ℓ -variate multilinear polynomial $g: \mathbb{F}^\ell \rightarrow \mathbb{F}$ is a polynomial of degree at most one in each variable.

2.3 Lookup arguments

Lookup arguments allow a prover to commit to two vectors $a \in \mathbb{F}^m$ and $b \in \mathbb{F}^m$ (with a polynomial commitment scheme) and prove that each entry a_i of vector a resides in index b_i of a pre-determined lookup table $T \in \mathbb{F}^N$. That is, For each $i = 1, \dots, m$, $a_i = T[b_i]$. Here, to emphasize the interpretation of T as a table, we use square brackets $T[i]$ to denote the i ’th entry of T . Here, if $b_i \notin \{1, \dots, N\}$, then $T[b_i]$ is undefined, and hence $a_i \neq T[b_i]$. We refer to a as the vector of *looked-up values* and b as the vector of *indices*.

Definition 2.3 (Lookup arguments, indexed variant). *Let $PC = (\text{Gen}, \text{Commit}, \text{Open}, \text{Eval})$ be an extractable polynomial commitment scheme for multilinear polynomials over \mathbb{F} . A lookup argument (for indexed lookups) for table $T \in \mathbb{F}^N$ is a SNARK for the relation*

$$\{ (\text{pp}, \mathcal{C}_1, \mathcal{C}_2, w = (a, b)) : a, b \in \mathbb{F}^m \wedge a_i = T[b_i] \text{ for all } i \in \{1, \dots, n\} \wedge \text{Open}(\text{pp}, \mathcal{C}_1, \tilde{a}) = 1 \wedge \text{Open}(\text{pp}, \mathcal{C}_2, \tilde{b}) = 1 \}.$$

Here $w = (a, b) \in \mathbb{F}^m \times \mathbb{F}^m$ is the witness, while pp , \mathcal{C}_1 , and \mathcal{C}_2 are public inputs.

¹⁵Any univariate polynomial commitment scheme can be transformed into a multilinear one, though the transformations introduce some overhead (see, e.g., [CBBZ23, BCHO22, ZXZS20]).

Definition 2.3 captures so-called *indexed* lookup arguments (this terminology was introduced in our companion work [STW23]). Other works consider *unindexed* lookup arguments, in which only the vector $a \in \mathbb{F}^m$ of looked-up values is committed, and the prover claims that *there exists* a vector b of indices such that $a_i = T[b_i]$ for all $i = 1, \dots, m$.

Definition 2.4 (Lookup arguments, unindexed variant). *Let $PC = (\text{Gen}, \text{Commit}, \text{Open}, \text{Eval})$ be an extractable polynomial commitment scheme for multilinear polynomials over \mathbb{F} . A lookup argument (for indexed lookups) for table $T \in \mathbb{F}^N$ is a SNARK for the relation*

$$\{(\text{pp}, \mathcal{C}_1, \mathcal{C}_2, a) : a \in \mathbb{F}^m \wedge \text{for all } i \in \{1, \dots, n\}, \text{ there exists a } b_i \text{ such that } a_i = T[b_i] \wedge \text{Open}(\text{pp}, \mathcal{C}_1, \tilde{a}) = 1\}.$$

Here $a \in \mathbb{F}^m \times \mathbb{F}^m$ is the witness, while pp and \mathcal{C}_1 are public inputs.

Jolt primarily requires indexed lookups. However, a few instructions (namely **ADVICE** and **MOVE**) require range checks, which are naturally handled by unordered lookups (to prove that a value is in the range $\{0, \dots, 2^L - 1\}$, perform an unordered lookup into the table T with $T[i] = i$ for $i = \{0, \dots, 2^L - 1\}$).

There are natural reductions in both directions, i.e., unindexed lookup arguments can be transformed into index lookup arguments and vice versa. To obtain an unindexed lookup argument from an indexed one, \mathcal{P} separately commits to the index vector b and applies the indexed lookup argument. Obtaining an indexed lookup argument from an unindexed one is slightly more complicated and is detailed in our companion paper [STW23, Appendix A]. Our companion work, **Lasso**, described below, directly yields an indexed lookup argument, and hence does not require this transformation.

A companion work: Lasso. Our companion work [STW23] introduces a family of lookup arguments called **Lasso**. The lookup arguments in this family are the first that do not require any party to cryptographically commit to the table vector $T \in \mathbb{F}^N$, so long as T satisfies one of the two structural properties defined below.

Definition 2.5 (MLE-structured tables). *We say that a vector $T \in \mathbb{F}^N$ is MLE-structured if for any input $r \in \mathbb{F}^{\log(N)}$, $\tilde{T}(r)$ can be evaluated with $O(\log(N))$ field operations.*

Definition 2.6 (Decomposable tables). *Let $T \in \mathbb{F}^N$. We say that T is c -decomposable if there exist a constant k and $\alpha \leq kc$ tables T_1, \dots, T_α each of size $N^{1/c}$ and each MLE-structured, as well as a multilinear α -variate polynomial g such that the following holds. As in Section 2.1, let us view T as a function mapping $\{0, 1\}^{\log N}$ to \mathbb{F} in the natural way, and view each T_i as a function mapping $\{0, 1\}^{\log(N)/c} \rightarrow \mathbb{F}$. Then for any $r \in \{0, 1\}^{\log N}$, writing $r = (r_1, \dots, r_c) \in \{0, 1\}^{\log(N)/c}$,*

$$T[r] = g(T_1[r_1], \dots, T_k[r_1], T_{k+1}[r_2], \dots, T_{2k}[r_2], \dots, T_{\alpha-k+1}[r_c], \dots, T_\alpha[r_c]).$$

We refer to T_1, \dots, T_α as sub-tables.

For any constant $c > 0$ and any c -decomposable table, our companion paper gives a lookup argument called **Lasso**, in which the prover commits to roughly $3cm + cN^{1/c}$ field elements. Moreover, all of these field elements are *small*, meaning that they are all in $\{0, \dots, m\}$ (specifically, they are counts for the number of times each entry of each subtable is read), or are elements of the subtables T_1, \dots, T_α . The verifier performs $O(\log(m) \log \log(m))$ hash evaluations and field operations, processes one evaluation proof from the polynomial commitment scheme applied to a multilinear polynomial in $\log m$ variables, and evaluates $\tilde{T}_1, \dots, \tilde{T}_\alpha$ each at a single randomly chosen point.

Our companion paper also describes a lookup argument called **Generalized-Lasso**, which applies to any MLE-structured table, not just decomposable ones.¹⁶ The main disadvantage of **Generalized-Lasso** relative to **Lasso** is that cm out of the $3cm + cN^{1/c}$ field elements committed by the **Generalized-Lasso** prover are random rather than small. As described in Section 1.3.1, such field elements can take an order of magnitude more work to commit to than small field elements.

¹⁶In fact, **Generalized-Lasso** applies to any table with *some* low-degree extension, not necessarily its multilinear one, that is evaluable in logarithmic time.

The relationship between MLE-structured and decomposable tables. For any decomposable table $T \in \mathbb{F}^N$, there is some low-degree extension \hat{T} of T (namely, an extension of degree at most k in each variable) that can be evaluated in $O(\log N)$ time. Specifically, the extension polynomial is

$$\hat{T}(r) = g(\tilde{T}_1(r_1), \dots, \tilde{T}_\alpha(r_c)).$$

In general, \hat{T} is not necessarily multilinear, so a table being decomposable does not necessarily imply that it is MLE-structured. But **Generalized-Lasso** actually applies to any table with a low-degree extension that is evaluable in logarithmic time. In this sense, decomposability (the condition required to apply **Lasso**) is a strictly stronger condition than what is necessary to apply **Generalized-Lasso**.

In Jolt, we show *all* lookup tables used are *both* c -decomposable (for any integer $c > 0$) as well as MLE-structured. We choose to apply **Lasso** rather than **Generalized-Lasso** due to its superior efficiency (which comes from the prover only committing to small field elements, avoiding the need to commit to random field elements). On the other hand, we believe that there would be meaningful improvements in simplicity of implementation if Jolt used **Generalized-Lasso** rather than **Lasso**. Arguably, the performance loss from using **Generalized-Lasso** in place of **Lasso** is justified by the simplicity benefits. See Section 7 for further discussion.

2.4 Offline Memory Checking

Any SNARK for VM execution has to perform *memory-checking*. This means that the prover must be able to commit to an execution trace for the VM (that is, a step-by-step record of what the VM did over the course of its execution), and the verifier has to find a way to confirm that the prover maintained memory correctly throughout the entire execution trace. In other words, the value purportedly returned by any read operation in the execution trace must equal the value most recently written to the appropriate memory cell. We use the term *memory-checking argument* to refer to a SNARK for the above functionality. Note that a lookup table $T \in \mathbb{F}^N$ can be viewed as a read-only memory of size N , with memory cell i initialized to $T[i]$. Hence, a lookup argument for indexed lookups (Definition 2.3) is equivalent to a memory-checking argument for read-only memories.

A variety of memory-checking arguments have been described in the research literature [ZGK⁺18, BCG⁺18, STW23, BFR⁺13, BSCGT13] (with the underlying techniques rediscovered multiple times). The most efficient are based on lightweight fingerprinting techniques for the closely related problem of *offline memory checking* [Lip89, BEG⁺91]. In this work, we use such an argument due to Spice [SAGL18], but optimize it using **Lasso**. For completeness, we can provide overview of other memory-checking arguments in Appendix B, and Spice’s in particular in Appendix B.3.

2.5 R1CS

The Jolt prover convinces the verifier that it correctly ran a VM for some number of steps on a specified input. Most of the VM’s work is verified in Jolt via a lookup argument and a memory-checking arguments. The remaining checks are simple and can be captured via any natural constraint system. One such option is the standard notion of rank-one constraint systems, defined next.

Definition 2.7. An *R1CS instance* is a tuple $(\mathbb{F}, A, B, C, M, M', N, x)$, where $A, B, C \in \mathbb{F}^{M \times M'}$, $M' \geq |x| + 1$, x denotes the public input and output, and there are at most N non-zero entries in each matrix. A vector $z = (w, 1, x) \in \mathbb{F}^{M'}$ is said to satisfy the instance if $A \cdot z \circ B \cdot z = C \cdot z$, where \cdot denotes matrix-vector product and \circ denotes Hadamard (i.e., entrywise) product.

3 An Overview of RISC-V and Jolt’s Approach

This section first provides a brief overview of the RISC-V instruction set architecture considered in this work. Our goal is to convey enough about the architecture that readers who have not previously encountered it can follow this paper. However, a complete specification is beyond the scope of this work, and can be found

at [And17].¹⁷ We also stick to regular control flow and do not support external events and other unusual run-time conditions like exceptions, traps, interrupts and CSR registers.

Informally, the RISC-V ISA consists of a CPU and a read-write memory, collectively called the *machine*.

Definition 3.1 (Machine State). *The machine state consists of $(PC, \mathcal{R}, \mathcal{M})$. \mathcal{R} denotes the 32 integer registers, each of W bits, where W is 32 or 64. \mathcal{M} is a linear read-write byte-addressable array consisting of a fixed number of total locations (such as 2^{20}) with each location storing 1 byte. The PC, also of W bits, is a separate register that stores the memory location of the instruction to be executed.*

Assembly programs consist of a sequence of instructions, each of which operate on the machine state. The instruction to be executed at a step is the one stored at the address pointed to by the PC. Unless specified by the instruction, the PC is advanced to the next memory location after executing the instruction. The RISC-V ISA specifies that all instructions are 32 bits long (i.e., 4 bytes), so advancing the PC to the next memory location entails incrementing PC by 4.

While RISC-V uses multiple formats to store instructions in memory, we can abstract away the details and represent all instructions in the following 5-tuple format.

Definition 3.2 (5-tuple RISC-V Instruction Format). *Any RISC-V instruction can be written in the following format: $[opcode, rs1, rs2, rd, imm]$. That is, each instruction specifies an operation code uniquely identifying its function, at most two source registers $rs1, rs2$, a destination register rd , and a constant value imm (standing for “immediate”) provided in the program code itself.*

Figure 1 provides a brief schematic of the CPU state change and instruction format. Operations read the source registers, perform some computation, and can do any or all of the following: read from memory, write to memory, store a value in rd , or update the PC. For example, the logical left-shift instruction “(SLL, $r5, r8, r2, -$)” reads the value stored in register #5, performs a logical left shift on the value by the length stored in register #8, and stores the result in register #2 (and does not involve any immediates).

As another example, the branch instruction “(BEQ, $r5, r8, -, imm$)” sets PC to be $PC + imm$ if the values stored in registers #5 and #8 are equal, or increments PC by 4, otherwise. (The destination register is not involved).

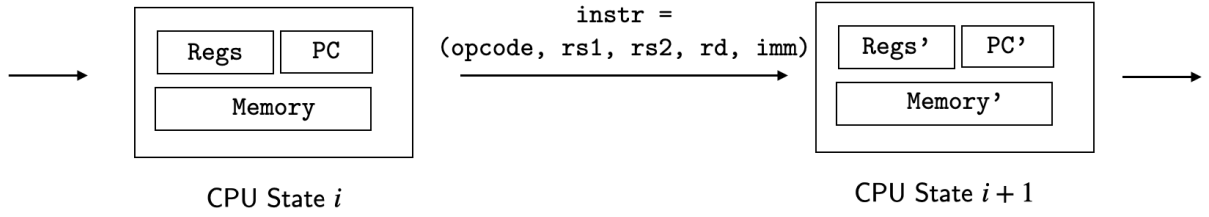
Unsigned and signed data types. For the RISC-V ISA, data in registers has no type. A register simply stores W bits. However, different instructions can be conceptualized as interpreting register values in different ways. Specifically, some instructions operate upon unsigned data types, while others operate over signed data types. All RISC-V instructions involving signed data types interpret the bits in a register as an integer via two’s complement representation.¹⁸ For many instructions (such as ADD and SUB), the use of two’s complement has the consequence that the instruction operates identically regardless of whether or not the inputs are interpreted as signed or unsigned. See Appendix C for more information on two’s complement notation and arithmetic.

For some instructions, like multiplication MUL, and integer comparison, the desired input/output behavior differs depending on whether the inputs are interpreted as signed or unsigned. In these cases, there will be two different RISC-V instructions, one for each interpretation. For example, there are MUL and MULU instructions, with the former interpreting its inputs as signed, and the latter interpreting its inputs as unsigned. Similarly, there are two integer comparison operations, SLT and SLTU.

Let z be an W -bit data type with constituent bits $[z_{W-1}, \dots, z_0]$ such that $z = \sum_{i=0}^{W-1} 2^i \cdot z_i$. When discussing instructions interpreting their W -bit inputs as signed data types represented in twos-complement format (e.g., Section 5.3), we refer to z_{W-1} as the sign bit of z , and denote this by z_s . (Concretely, the sign bit of a 64-bit register value z will be $z_s = z_{63}$.) We use $z_{<s}$ to refer to $[z_{W-2}, \dots, z_0] \in \{0, 1\}^{W-1}$.

¹⁷Another helpful resource for interested readers is Lectures 5-8 at <https://inst.eecs.berkeley.edu/~cs61c/resources/su18 lec/>.

¹⁸See https://en.wikipedia.org/wiki/Two%27s_complement for an overview of how two’s complement maps bit vectors in $\{0, 1\}^L$ to integers in $\{-2^L, \dots, 2^L - 1\}$ and vice versa.



(a) The CPU state and instruction formats.

CPU Step Transition:

1. Read the instruction at location PC in Program Code.
Parse instruction as `[opcode, rs1, rs2, rd, imm]`.
2. Read the W -bit values stored in registers `rs1, rs2`.
3. If required, write to or read from memory.
The value written and memory location accessed are derived from the values stored in `rs1, rs2, imm`.
4. Perform the instruction's function on the values read from registers and `imm` to get `result`.
Examples of functions are arithmetic, logical and comparison operations.
5. Store `result` to register `rd`.
Only a few instructions, like STOREs, do not involve `rd`.
6. Update PC.
PC is usually incremented by 4, but instructions like jumps and branches update PC in other ways.

(b) The broad stages of a CPU step transition.

Figure 1: A model of RISC-V's CPU state and transition function. Note that the transition function is deterministic and all information required, such as the location of memory accessed, is derived from the CPU state and `instr`.

Sign and Zero Extensions. A “sign-extension” of an L -bit value z to W bits (where $L < W$) is the W -bit value $z_{\text{sign-ext}}$ with bits $[z_s, \dots, z_s, z_{L-1}, \dots, z_0]$. That is, the sign bit of z is replicated to fill the higher-order bits of z until it reaches length W . A “zero-extension” is when, instead of the sign bit, the 0 bit is used. This results in W -bit $z_{\text{zero-ext}}$ with bits $[0, \dots, 0, z_{L-1}, \dots, z_0]$.

3.1 Performing instruction logic using lookups

As described in Section 2.3, the Jolt paradigm avoids the complexity of implementing each instruction's logic as constraints by encapsulating instruction execution into a lookup table. Specifically, we identify an “evaluation table” for each operation `opcode`, $T_{\text{opcode}}[x \parallel y] = r$, that contains the required result for all possible inputs x, y . Jolt combines the tables for all instructions into one table and thus makes only one lookup query per step to this table as $T_{\text{risc-v}}[\text{opcode} \parallel x \parallel y] = r$ (see Section 7 for details). Given a processor and instruction set, this table is fixed and independent of the program or inputs. The key contribution of Jolt is to design these enormous tables with a certain structure (see Definition 2.6) that allows for efficient lookup arguments using Lasso.

Preparing operands and the lookup query. The main responsibility of the constraint system is to prepare the appropriate operands x, y at each step before the lookup. This is efficient to do as the operands only come from the set $\{\text{value in } \text{rs1}, \text{value in } \text{rs2}, \text{imm}, \text{PC}\}$. This means, for example, that the instructions

ADD and ADDI are expressed by the same lookup table as they only differ in whether the second operand comes from register `rs2` or is `imm`, respectively. With the operands prepared, the lookup query is then committed to by the prover and fed to the lookup argument for verification. The query is of the form `opcode || z` where z is generally $x || y$ or $(x + y)$ or $(x \times y)$, making it either $2 \cdot W$ or $W + 1$ bits in length. The prover provides as advice the claimed entry, `result`, in the lookup table corresponding to the query.

The trace of all lookup queries and entries is sent to **Lasso**. As described in Definition 2.6, **Lasso** requires the query to be split into “chunks” which are fed into different subtables. The prover provides these chunks as advice, which are c in number for some small constant c , and hence approximately W/c or $2W/c$ bits long, depending on the structure of z . The constraint system must verify that the chunks correctly constitute z , but need not perform any range checks as the **Lasso** algorithm itself later implicitly enforces these on the chunks.

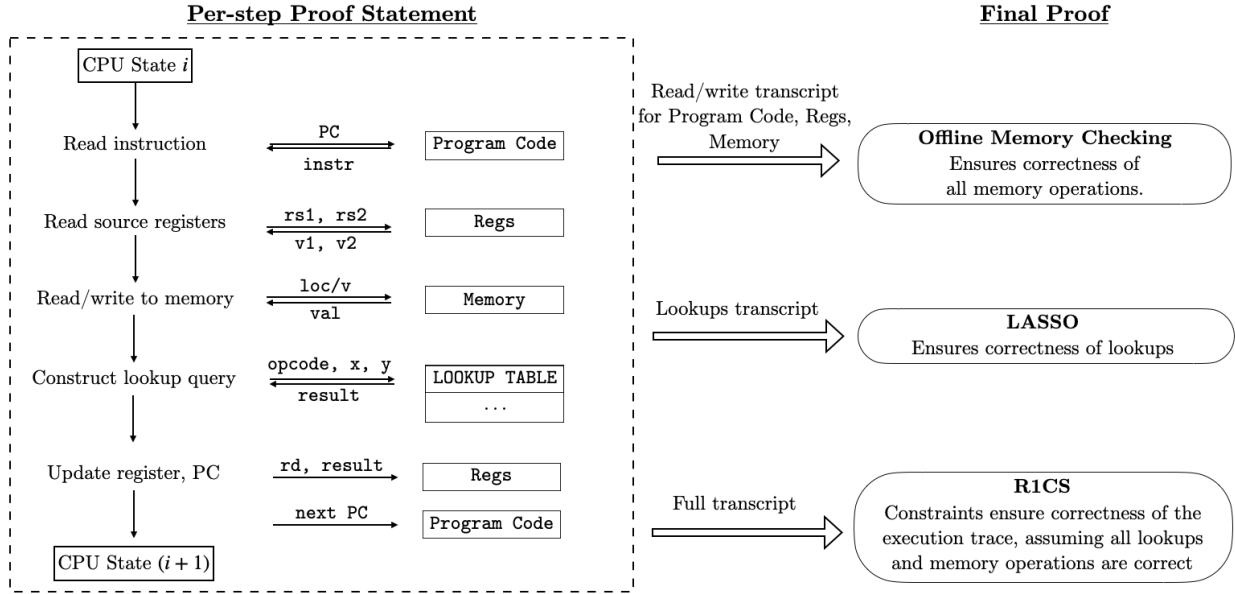


Figure 2: Proving the correctness of CPU execution using offline memory checking (Section 3.2) and lookups (Section 3.1).

3.2 Using Memory-Checking

The machine state transition involves reading from and writing to three conceptually separate parts of memory: (1) the program code, (2) the registers and (3) the random access memory. As discussed in Section 2.4, the most efficient way to enforce correct reads and writes is by using the offline memory checking techniques. Unlike other operations, loads and stores do not involve lookups to a large table to perform their core function.

As is standard in zkVM design, Jolt conceptualizes the memory-checking procedure as a black box that guarantees correctness of all the memory reads and writes required by the CPU execution, and hence the proof proceeds assuming these operations are correct. For example, suppose a value v was written to location k at step t of the CPU’s execution. Later, when location k needs to be read, the prover sends as advice (v', t') claiming that the most recent write to location k was done at step t' and the value written was v' . All of these read and write (k, v, t) tuples are committed to during the execution of Jolt and later fed to the memory-checking argument, which will only pass if every read was consistent with the latest write. The main job of the constraint system here is to prevent any cheating by enforcing range checks on the time values provided by the prover (t' in the above example). These are done efficiently using **Lasso**. See Appendix B.3 for more details.

Supporting byte-addressable memory. RISC-V requires that memory be byte-addressable (as opposed to word-addressable). A load or store operation may read up to $W/8$ (which equals four and eight for 32-bit and 64-bit processors, respectively) bytes in a given instruction. Thus, when writing a W -bit value v , the prover must provide its byte-decomposition $[v_1 \dots v_{W/8}]$ as each byte is stored in a separate address in memory. Jolt enforces range-checks on the provided bytes through lookups performed using Lasso. See Appendix A.2 for more details.

Furthermore, certain load instructions also require the values read from memory to be sign-extended to W total bits before stored in the register. This requires only a short lookup query using the highest order byte to a small table to obtain the sign bit.

3.3 Formatting assembly code

Before the proof starts, the assembly code is formatted into the 5-tuple form of Definition 3.2: (`opcode`, `rs1`, `rs2`, `rd`, `imm`). Additionally, each instruction also comes with 14 one-bit “flags” `opflags[14]` that guide the constraint system. For example, `opflag[5]` is 1 for only Jump instructions, and `opflags[7]` is 1 if and only if the lookup’s result is to be stored in `rd`. Note that these flags are fixed for any given instruction. See Appendix A.1 for a list of all the flags used in Jolt. Load and store instructions (that is, those involving memory) involve additional flags of their own.

In RISC-V, instructions may need to sign-extend or zero-extend `imm` to W bits. This is a deterministic choice that depends only on the instruction (and is independent of the rest of the program or inputs). Thus, when formatting the instruction to the 5-tuple format, the immediate is appropriately extended to W bits.

Putting this together, before the proof starts, the prover and verifier convert the RISC-V assembly code and store it in a manner accessible to the constraint system. For the purposes of memory-checking, program code is in a different address-space than regular random-access memory. It is read-only and is initialized as follows: the original instruction at location `PC` is converted to the form (`opcode`, `rs1`, `rs2`, `rd`, `imm`, `opflags[14]`) and the elements of this tuple are respectively stored at locations $[6 \cdot \text{PC}, \dots, 6 \cdot \text{PC} + 5]$. The constraint system performs six memory-checking reads per CPU step to obtain these entries. As the program code is read-only, the prover simply commits to the elements of the tuple and only one extra timestamp that we know to be less than T , the number of steps up to that point in the program (see Appendix B.3 for more details).

4 Analyzing MLE-structure and Decomposability

This section illustrates the process of designing MLE-structured tables and decomposing them as per Definition 2.6 required by Lasso. We first establish notation and then design the tables for three important functions that are used as building blocks for the tables of many RISC-V instructions: equality, less than, and shifts.

4.1 Notation

Associating field elements with bit-vectors and vice versa. Let z be a field element in $\{0, 1, \dots, 2^W - 1\} \subset \mathbb{F}$. We denote the binary representation of z as $\text{bin}(z) = [z_{W-1}, \dots, z_0] \in \{0, 1\}^W$. Here, z_0 is the least significant bit (LSB), while z_{W-1} is the most significant bit (MSB). That is, $z = \sum_{i=0}^{W-1} 2^i z_i$. We refer to the “sign-bit” of z as $z_s = z_{W-1}$.

We use $z_{<i}$ to refer to the subsequence $[z_{i-1}, \dots, z_0]$. Analogously, $z_{>i}$ refers to the subsequence $[z_{W-1}, \dots, z_{i+1}]$. Similarly, given a vector $z = [z_{W-1}, \dots, z_0] \in \{0, 1\}^W$, we denote the associated field element as $\text{int}(z) = \sum_{i=0}^{W-1} 2^i \cdot z_i$.

Remark 1. In the above paragraphs, we used an italicized z to denote both a field element in $\{0, \dots, 2^W - 1\}$ and a vector in $\{0, 1\}^W$. Throughout the paper, which of the two sets any variable z resides in will be clear from context.

Concatenation of bit vectors. Given two bit vectors $x, y \in \{0, 1\}^W$, we use $x \parallel y$ to refer to the number whose binary representation is the concatenation $[x_{W-1}, \dots, x_0 \parallel y_{W-1}, \dots, y_0]$. Under this definition, it holds that $\text{int}(x \parallel y) = \text{int}(x) \cdot 2^W + \text{int}(y)$.

Decomposing bit vectors into chunks. For a constant c , and any $x \in \{0, 1\}^L$, we divide the bits of input x naturally into chunks

$$x = [x_{W-1} \dots x_0] = X_{c-1} \parallel \dots \parallel X_2 \parallel X_0, \quad (3)$$

with each $X_i \in \{0, 1\}^{W/c}$.

Throughout the following description of tables and decompositions, we assume c divides W for simplicity. However, this is not necessary. In fact, it is more efficient in practice to set $c = 3$ for $W = 32$ and $c = 6$ for $W = 64$, resulting in some chunks being length 10 and others being length 11.

4.2 Three instructive functions and associated lookup tables

Let field \mathbb{F} be a prime order field of size at least 2^W (for concreteness, let us fix W to be 64). Let x and y denote field elements that are guaranteed to be in the set $\{0, 1, \dots, 2^W - 1\}$.

4.2.1 The Equality function

MLE-structured. The equality function EQ takes as inputs two vectors $x, y \in \{0, 1\}^W$ of identical length and outputs 1 if they are equal, and 0 otherwise. We will use a subscript to clarify the number of bits in each input to EQ , e.g., EQ_W denotes the equality function defined over domain $\{0, 1\}^W \times \{0, 1\}^W$. It is easily confirmed that the multilinear extension of EQ_W is as follow:

$$\widetilde{\text{EQ}}_W(x, y) = \prod_{j=0}^{W-1} (x_j y_j + (1 - x_j)(1 - y_j)). \quad (4)$$

Indeed, the right hand side is clearly a multilinear polynomial in x and y , and if $x, y \in \{0, 1\}^W$, it equals 1 if and only if $x = y$. Hence, the right hand side must equal the unique multilinear extension of the equality function. Clearly, it can be evaluated at any point $(x, y) \in \mathbb{F}^W \times \mathbb{F}^W$ with $O(W)$ field operations.

Decomposability. To determine whether two W -bit inputs $x, y \in \{0, 1\}^W$ are equal, one can decompose x and y into c chunks of length W/c , compute equality of each chunk, and multiply the results together.

Let $x = [X_{c-1}, \dots, X_0]$ and $y = [Y_{c-1}, \dots, Y_0]$ denote the decomposition of x and y into c chunks each, as per Equation (3). Let EQ_W denote the “big” table of size $N = 2^{2W}$ indexed by pairs (x, y) with $x, y \in \{0, 1\}^W$, such that $\text{EQ}_W[x \parallel y] = \widetilde{\text{EQ}}_W(x, y)$. Let $\text{EQ}_{W/c}$ denote the “small” table of size $N^{2W/c}$ indexed by pairs (X, Y) of chunks $X, Y \in \{0, 1\}^{W/c}$, such that $\text{EQ}_{W/c}[X \parallel Y] = 1$ if $X = Y$ and $\text{EQ}_{W/c}[X \parallel Y] = 0$ otherwise. The table below asserts that evaluating the equality function on x and y is equivalent to evaluating the equality function on each chunk $X_i \parallel Y_i$ and multiplying the results.

CHUNKS	SUBTABLES	FULL TABLE
$C_i = X_i \parallel Y_i$	$\text{EQ}_{W/c}[X_i \parallel Y_i] = \widetilde{\text{EQ}}_{W/c}(X_i, Y_i)$	$\text{EQ}_W[x, y] = \prod_{i=0}^{c-1} \text{EQ}_{W/c}[X_i \parallel Y_i]$

The (lone) subtable $\text{EQ}_{W/c}$ is MLE-structured by Equation (4).

4.2.2 Less Than comparison

MLE-structured. The comparison of two unsigned data types $x, y \in \{0, 1, \dots, 2^{W-1}\}$ is involved in many instructions. For example, **LTU** outputs 1 if $x < y$ and 0 otherwise, where the inequality interprets x and y as integers in the natural way. Note that the inequality computed here is strict. Consider the following $2W$ -variate multilinear polynomial (**LTU** below stands for “less than unsigned”):

$$\widetilde{\text{LTU}}_i(x, y) = (1 - x_i) \cdot y_i \cdot \widetilde{\text{EQ}}_{W-i-1}(x_{>i}, y_{>i}). \quad (5)$$

Clearly, this polynomial satisfies the following two properties:

- (1) Suppose $x \geq y$. Then $\widetilde{\text{LTU}}_i(x, y) = 0$ for all i .
- (2) Suppose $x < y$. Let k be the first index (starting from the MSB of x and y) such that $x_k = 0$ and $y_k = 1$. Then $\widetilde{\text{LTU}}_k(x, y) = 1$ and $\widetilde{\text{LTU}}_i(x, y) = 0$ for all $i \neq k$.

Based on the above properties, it is easy to check that

$$\widetilde{\text{LTU}}(x, y) = \sum_{i=0}^{W-1} \widetilde{\text{LTU}}_i(x, y). \quad (6)$$

Indeed, the right hand side is clearly multilinear, and by the two properties above, it equals $\widetilde{\text{LTU}}(x, y)$ whenever $x, y \in \{0, 1\}^W$. It is not difficult to see that the right hand side of Equation (6) can be evaluated at any point $(x, y) \in \mathbb{F}^W \times \mathbb{F}^W$ with $O(W)$ field operations as the set $\{\widetilde{\text{EQ}}_{W-i}(x_{>i}, y_{>i})\}_{i=0}^{W-1}$ can be computed in $O(W)$ total steps using the recurrence relation

$$\widetilde{\text{EQ}}_{W-i-1}(x_{>i}, y_{>i}) = \widetilde{\text{EQ}}_{W-i-2}(x_{>(i+1)}, y_{>(i+1)}) \cdot \widetilde{\text{EQ}}(x_i, y_i). \quad (7)$$

See [Tha22, Figure 3.3] for a depiction of this procedure.

Decomposing $\widetilde{\text{LTU}}$. A similar reasoning to the derivation of Equation (6) reveals the following. As usual, break x and y into c chunks, $X_{c-1} \parallel \dots \parallel X_0$ and $Y_{c-1} \parallel \dots \parallel Y_0$. Let $\text{LTU}_{W/c}[X_i \parallel Y_i] = \widetilde{\text{LTU}}_{W/c}(X_i, Y_i)$ denote the subtable with entry 1 if $X_i < Y_i$ when interpreted as unsigned (W/c) -bit data types, and 0 otherwise. Then

$$\text{LTU}_W[x \parallel y] = \sum_{i=0}^{c-1} \text{LTU}_{W/c}[X_i \parallel Y_i] \cdot \text{EQ}_{W/c}[X_{>i} \parallel Y_{>i}] = \sum_{i=0}^{c-1} (\text{LTU}_{W/c}[X_i \parallel Y_i] \cdot \prod_{j<i} \text{EQ}_{W/c}(X_j \parallel Y_j)).$$

Thus, evaluating $\text{LTU}(x, y)$ can be done by evaluating $\text{LTU}_{W/c}$ and $\text{EQ}_{W/c}$ on each chunk (X_i, Y_i) ($\text{EQ}_{W/c}$ need not be evaluated on the lowest-order chunk (X_c, Y_c)). This is summarized in the table below.

CHUNKS	SUBTABLES	FULL TABLE
$C_i = X_i \parallel Y_i$	$\text{LTU}_{W/c}[X_i \parallel Y_i], \text{EQ}_{W/c}[X_i \parallel Y_i]$	$\text{LTU}_W[x \parallel y] = \sum_{i=0}^{c-1} \text{LTU}_{W/c}[X_i \parallel Y_i] \cdot \prod_{j<i} \text{EQ}_{W/c}[X_j \parallel Y_j]$

The two subtables **LTU** and **EQ** are MLE-structured by Equations (4) and (6).

4.2.3 Shift Left Logical

MLE-structured. **SLL** takes an W -bit integer x and a $\log(W)$ -bit integer y , and shifts the binary representation of x to the left by length y . Bits shifted beyond the MSB of x are ignored, and the vacated lower bits are filled with zeros.¹⁹ For a constant k , let

$$\widetilde{\text{SLL}}_k(x) = \sum_{j=k}^{W-1} 2^j \cdot x_{j-k}. \quad (8)$$

¹⁹For $L = 32$ -bit data types, the RISC-V manual says that the “shift amount is encoded in the lower $5 = \log(W)$ bits”.

It is straightforward to check that the right hand side of Equation (8) is multilinear (in fact, linear) function in x , and that when evaluated at $x \in \{0, 1\}^W$, it outputs the unsigned W -bit data type whose binary representation is the same as that of the output of the SLL instruction on inputs x and k , $\text{SLL}(x, k)$.

Now consider

$$\widetilde{\text{SLL}}(x, y) = \sum_{k \in \{0, 1\}^{\log W}} \widetilde{eq}(y, k) \cdot \widetilde{\text{SLL}}_k(x). \quad (9)$$

It is straightforward to check that the right hand side of Equation (9) is multilinear in (x, y) , and that, when evaluated at $x \in \{0, 1\}^W \times \{0, 1\}^{\log W}$, it outputs the unsigned W -bit data type $\text{SLL}(x, y)$.

Decomposability. We split the value to be shifted, x , into c chunks, X_1, \dots, X_c , each consisting of $W' = W/c$ bits. y has only one chunk, Y_0 , consisting of the lowest order $\log W$ bits. As explained below, we decompose a lookup into the evaluation table of SLL into a lookup into c different subtables, each of size $2^{W' + \log W}$. For $W = 64$, a reasonable setting of c would be 4 (instead of the usual $c = 6$ for most other instructions), ensuring that $2^{W' + \log W} = 2^{20}$.

Conceptually, each chunk X_i of X needs to determine how many of its input bits goes “out of range” after the shift of length y . By out of range, we mean that shifting x left by y bits causes those bits to overflow the MSB of x and hence not contribute to the output of the instruction.

For chunks $i = 0, \dots, (c - 1)$ and shift length $k \in \{0, 1\}^{\log W}$, define:

$$m_{i,k} = \min\{W', \max\{0, (\text{int}(k) + W' \cdot (i + 1)) - W\}\}$$

Here, $m_{i,k}$ equals the number of bits from the i 'th chunk that go out of range. Let $m'_{i,k} = W' - m_{i,k} - 1$ denote the index of the highest-order bit within the i 'th chunk that does *not* go out of range. Then the evaluation table of SLL decomposes into c smaller tables $\text{SLL}_0, \dots, \text{SLL}_{c-1}$ as follows.

CHUNKS	SUBTABLES	FULL TABLE
$C_i = X_i \parallel Y_0$	$\text{SLL}_i[X_i \parallel Y_0] = \sum_{k \in \{0, 1\}^{\log W}} \widetilde{\text{EQ}}(Y_0, k) \cdot \left(\sum_{j=0}^{m'_{i,k}} 2^{j + \text{int}(k)} \cdot X_{i,j} \right)$	$\text{SLL}[x \parallel y] = \sum_{i=0}^{c-1} 2^{i \cdot W'} \cdot \text{SLL}_i[X_i \parallel Y_c]$

Note that each SLL_i can be evaluated at any input $(x, y) \in \mathbb{F}^{W'} \times \mathbb{F}^{\log W}$ in $O(W')$ field operations. Indeed, the set $\{\widetilde{\text{EQ}}(Y_0, k)\}_{k \in \{0, 1\}^{\log W}}$ can be computed in $O(W)$ field operations via the recurrence in Equation (7). Similarly, the set $\{2^{j + \text{int}(k)}\}_{i \in \{0, \dots, c-1\}, k \in \{0, 1\}^{\log W}}$ can be computed with $O(W)$ field operations. It follows that $\text{SLL}_0(x \parallel y), \dots, \text{SLL}_{c-1}(x \parallel y)$ can be evaluated in $O(W)$ field operations in total.

4.3 The Cost of a Lookup

We briefly state the costs incurred by the prover when making a lookup query. As usual, this is analyzed in terms of the bit-lengths of the elements to be committed to. Most lookup queries require separating $x \parallel y \in \{0, 1\}^{2W}$ into c chunks which are then sent to c subtables within Lasso. This involves the Jolt prover committing to $3c$ elements: (1) c are the chunks themselves, which are of $2W/c$ bits, (2) c are the entries in the subtables involved, which are up to W/c bits long, (3) c elements are “access counts” for the subtables, which can go up to T , the current step count.

With the parameter settings $(W = 32, c = 3)$ and $(W = 64, c = 6)$, the first two bit-lengths involved are $2W/c \approx 22$ and $W/c \approx 11$. Many instructions (such as SLL), involve smaller lookup queries of closer to W bits and can be split into fewer chunks, leading to a proportional reduction in the number of elements committed to in each group.

An interesting case is that of the LTU subtable, which uses c chunks but involves $2c$ total subtables (recall that each chunk is sent to both LTU and EQ). This involves committing to c extra elements in both

groups (2) and (3) above. While most instructions do not incur this extra cost, we report costs for this instruction in Section 8, as it captures the worst-case scenario (excluding instructions that are not handled directly, but are rather transformed into a short sequence of other instructions, such as division, see Section 6.1).

5 Evaluation Tables for the Base Instruction Set

We now consider each of the RISC-V instructions one at a time, and analyze the MLE-structure and the decomposability of their evaluation tables. Conceptually, Jolt’s “one giant lookup table” is obtained by simply concatenating each of the evaluation tables for all instructions (see Section 7.3). For most instructions, we first present an MLE-structured table and describe how it can be decomposed into subtables, as required by Lasso.

5.1 Logical instructions

Each instruction performs the corresponding operation bitwise over the W bits of x and y and stores the W -bit result in **rd**. The lookup tables here have a row for each possible $x \parallel y$ with the entry being the desired output to be stored in **rd**.

OP	INDEX	FULL MLE
AND	$x \parallel y \in \{0, 1\}^W \times \{0, 1\}^W$	$\sum_{i=0}^{W-1} 2^i \cdot (x_i \cdot y_i)$
OR	$x \parallel y \in \{0, 1\}^W \times \{0, 1\}^W$	$\sum_{i=0}^{W-1} 2^i \cdot (x_i + y_i - x_i \cdot y_i)$
XOR	$x \parallel y \in \{0, 1\}^W \times \{0, 1\}^W$	$\sum_{i=0}^{W-1} 2^i \cdot (x_i \cdot (1 - y_i) + y_i \cdot (1 - x_i))$

Decomposition. These MLEs can further be decomposed in the natural way, requiring only one subtable per instruction. For example, a bitwise AND on two W -bit inputs can be decomposed into c bitwise ANDs, each on two (W/c) -bit inputs. That is, if $D_0, \dots, D_{c-1} \in \{0, 1\}^{W/c}$ denote the results of these “smaller” bitwise AND operations, with $C_i = \text{AND}_{W/c}(X_i, Y_i)$ then the result of the W -bit operation is simply $\sum_{i=0}^{c-1} 2^{(W/c) \cdot i} D_i$. The decompositions for OR and XOR follow analogously.

5.2 Arithmetic instructions

Addition. For $x, y \in \{0, 1\}^W$, $\text{ADD}(x, y)$ returns the lowest W bits of (the binary representation of) the sum $\text{int}(x) + \text{int}(y)$. We need not specify whether the inputs to these instructions are signed versus unsigned because in RISC-V, signed data types are represented via two’s complements. When the inputs and outputs are viewed as strings in $\{0, 1\}^W$, the behavior of ADD and SUB is identical for signed data types as for unsigned ones.

As finite field addition costs just one constraint in R1CS, we cheaply compute $z = x + y$ in the circuit, where here, addition is performed over finite field \mathbb{F} . However, this sum can be $W + 1$ bits long, i.e., z can be any field element in $\{0, \dots, 2^{W+1} - 2\}$. The prescribed behavior for the RISC-V instruction ADD in this event is for the “overflow bit” to be ignored.

To this end, the lookup table for the ADD instructions contains an entry for all possible $(W + 1)$ -bit vectors $z \in \{0, 1\}^{W+1}$, with each z ’s entry equal to the field element $\sum_{i=0}^{W-1} 2^i \cdot z_i$ equal to $\text{int}(z_{<W})$. Note that this lookup table has size only 2^{W+1} , which is less than the tables of size 2^{2W} that we identify for most RISC-V instructions.

Subtraction. Due to RISC-V's use of two's complement representation of signed data types, subtraction can be performed using addition. Specifically, $\text{SUB}(x, y)$ outputs the same W -bit string as

$$\text{ADD}(x, \text{bin}(2^W - \text{int}(y))).$$

In words, subtracting y from x is equivalent to adding the two's complement of y to x .²⁰

OP	INDEX	FULL MLE
$\text{ADD}(x, y)$	$z = \text{bin}(x + y) \in \{0, 1\}^{W+1}$	$\sum_{i=0}^{W-1} 2^i z_i$
$\text{SUB}(x, y)$	$z = \text{bin}(x + (2^W - y)) \in \{0, 1\}^{W+1}$	$\sum_{i=0}^{W-1} 2^i z_i$

Decomposition. The decomposition of the above lookup table is simple (and essentially equivalent to the case of range checks considered in our companion paper on Lasso [STW23]).

5.3 Set Less Than

SLTU and SLT return 1 if $x < y$ and 0 otherwise, where x, y are unsigned and two's complement signed W -bit numbers, respectively.

The tables for SLTU is equivalent to LTU derived in Section 4.2.2. For SLT, we must additionally take into consideration the sign bits of the two numbers and resort to a comparison of the remaining bits only when the sign bits are the same.

OP	INDEX	FULL MLE
SLTU	$x \parallel y$	See Section 4.2.2
SLT	$x \parallel y$	$\widetilde{\text{LTS}} = x_s \cdot (1 - y_s) + \widetilde{\text{EQ}}(x_s, y_s) \cdot \widetilde{\text{LTU}}(x_{<s}, y_{<s})$

Decomposition. The decomposition of SLTU was discussed in Section 4.2.2 and requires two subtables of size $2^{W/c}$. The decomposition of SLT uses the same decomposition (applied to $x_{<s}, y_{<s}$, which has $W - 1$ rather than W bits), but additionally devotes more two subtables to compute the $x_s(1 - y_s)$ and $\widetilde{\text{EQ}}(x_s, y_s)$ terms. This brings the total to four types of subtables and $2c + 1$ total subtables used for SLT.

5.4 Shifts

SLL(x, y) (Shift Left Logical), SRL(x, y) (Shift Right Logical) and SRA(x, y) (Shift Right Arithmetic) are the three shift operations. All shift operations take a W -bit input x , shift it by a length defined by the lowest $\log W$ bits of y and return W -bit values. Bits shifted beyond the MSB or LSB are ignored. In logical shifts, vacated bits are filled by zeros, and in arithmetic shifts, the vacated bits are filled by the sign bit of the original input x .

The MLE-structured table for SLL and its decomposition were presented in Section 4.2.3. The tables for SRL and SRA are presented below. Let $W' = W/c$. For chunks $i = 0, \dots, (c - 1)$ and shift length $k \in \{0, 1\}^{\log W}$, define:

$$m_{i,k}^r = \max\{0, \min\{16, \text{int}(k) - W' \cdot i\}\}.$$

²⁰In a system implementing arithmetic on W -bit unsigned data types, the quantity 2^W cannot be represented. Hence, the two's complement of y needs to be computed in two steps, as $(2^W - 1 - y) + 1$, with the expression in parenthesis computed first, and then one added to the result. See https://en.wikipedia.org/wiki/Two%27s_complement for details. In Jolt, the quantity $2^W - \text{int}(y)$ will be computed directly in the field \mathbb{F} , which we assume to have characteristic more than 2^W . Hence, the quantity 2^W can be represented, and the two's complement of y is (the binary representation of) the field element $2^W - y$.

Here, $m_{i,k}^r$ equals the number of bits from the i 'th chunk that go out of range (that is, to the “right” of bit 0). Note that $m_{i,k}^r$ is also the index of the lowest-order bit within the i 'th chunk that does *not* go out of range. That is, i 'th chunk X_i will have subsequence $[X_{i,0}, X_{i,m_{i,k}^r-1}]$ go out of range and the remaining values will now be present in indices $[W' \cdot (i-1) - \text{int}(k) + m_{i,k}^r, \dots, W' \cdot (i-1) - \text{int}(k) + W' - 1]$ of the final output. The evaluation table of SLL decomposes into c smaller tables $\text{SLL}_0, \dots, \text{SLL}_{c-1}$ as follows.

CHUNKS	SUBTABLES	FULL TABLE
$C_i = X_i \parallel Y_0$	$\text{SRL}_i[X_i \parallel Y_0] = \sum_{k \in \{0,1\}^{\log W}} \widetilde{\text{EQ}}(Y_0, k) \cdot \left(\sum_{j=m_{i,k}^r}^{W'-1} 2^{W' \cdot (i-1) - \text{int}(k) + j} \cdot X_{i,j} \right)$	$\text{SRL}[x \parallel y] = \sum_{i=0}^{c-1} \text{SRL}_i[C_i]$

The SRA instruction uses an extra subtable to perform sign extension. This subtable takes as its input chunk $C_c = x_s \parallel Y_0$, where x_s is the sign bit of the input.

CHUNKS	SUBTABLES	FULL TABLE
$i \in [0, c-1]: C_i = X_i \parallel Y_0$ and $C_c = x_s \parallel Y_0$	For $i \in [0, c-1]$, $\text{SRA}_i[C_i] = \text{SRL}_i[C_i]$ and $\text{SRA}_c(x_s \parallel Y_0) = \sum_{i=W-\text{int}(k)}^{W-1} \widetilde{\text{EQ}}(Y_0, k) \cdot 2^i \cdot x_s$	$\text{SRA}[x \parallel y] = \sum_{i=0}^c \text{SRA}_i[C_i]$

5.5 Immediate Loads

$\text{AUIPC}(x, y)$ takes the 20-bit immediate (operand y here), adds it to PC (operand x here) and stores the output in the destination register rd , but does *not* change the PC. LUI takes the 20-bit immediate (operand y , here) and loads it into the upper 20 bits of the destination register rd .

In both of these instructions, the 20-bit immediate is formatted (see Section 3.3) into a W -bit value with the 20 significant bits stored in the *higher positions* of imm in program code (as opposed to the lower positions, as done for all other instructions). With this pre-processing, AUIPC can be treated identically to ADD with the two operands being PC and imm .

As the above pre-processing does most of the work, the only task for LUI in the circuit is to store the given imm as provided into rd . This does not require a lookup table.

OP	INDEX	FULL MLE
AUIPC	$z = x + y$	$\sum_{i=0}^{W-1} 2^i z_i \quad // \text{ identical to ADD}$

Decomposability. This table can be decomposed just like the table for ADD .

5.6 Jumps

JAL sets $\text{PC} \leftarrow \text{PC} + \text{imm}$ and stores the address of the memory location after this new PC (obtained by incrementing it by 4) into rd . JALR similarly sets PC to be the sum $\text{PC} + \text{imm}$ but with the LSB set to 0. It sets rd to be memory location after this new PC.

For both jump instructions, the sum $z \leftarrow \text{PC} + \text{imm} + 4$ is calculated using constraints. The lookup table for JAL is identical to that of ADD and returns the lower W bits of z . The table for JALR does the same but then sets the LSB to 0, as well. In both instructions, PC is set to be the lookup's result minus 4. This subtraction is performed with constraints. If it results in an underflow the memory-checking will fail when reading PC in the next step.

OP	INDEX	FULL MLE
JAL	$z = x + y + 4$	$\sum_{i=0}^{W-1} 2^i z_i$
JALR	$z = x + y + 4$	$\sum_{i=1}^{W-1} 2^i z_i$

Decomposability. These tables can be decomposed just like the table for ADD. Note that the lookup queries here are also $W + 1$ bits long.

5.7 Branches

The B[COND] instructions set $PC \leftarrow PC + \text{imm}$ if $\text{COND}(x, y) = \text{true}$. If false, they resort to the default change in PC.

The new PC is computed using only constraints and not with a lookup (as the lookups here are used to test the branching condition). When imm is positive, the sum $(PC + \text{imm})$ obtained is correct as is. However, when imm is negative, we must perform $(PC - \text{imm})$ directly without using two's complement subtraction (as that might result in an overflow). To choose which to perform, the sign of imm is stored in the program code itself during formatting as one of the **opflags** discussed in Section 3.3. If this subtraction results in an underflow the memory-checking algorithm will fail when reading PC in the next step.

Now, the lookup is used to perform the comparisons to decide whether to use this new shifted PC or not. The MLE for doing both signed and unsigned strict less than comparisons were discussed in Section 4.2.2 and used in SLT/SLTU. We use the same MLEs here, along with the $\widetilde{\text{EQ}}$ MLE.

OP	INDEX	FULL MLE
BEQ	$x \parallel y$	$1 - \widetilde{\text{EQ}}(x, y)$
BNE	$x \parallel y$	$\widetilde{\text{EQ}}(x, y)$
BLTU	$x \parallel y$	$\widetilde{\text{LTU}}(x, y) // \text{ as used in SLTU}$
BLT	$x \parallel y$	$\widetilde{\text{LTS}}(x, y) // \text{ as used in SLT}$
BGEU	$x \parallel y$	$1 - \widetilde{\text{LTU}}(x, y)$
BGE	$x \parallel y$	$1 - \widetilde{\text{LTS}}(x, y)$

Decomposability. These MLE-structured tables can be decomposed with the same techniques used for the EQ, LTU and LTS tables.

5.8 Memory Loads and Stores

RISC-V uses a byte-addressable memory system that can be accessed by only the following variants of the load and store instructions:

- LD reads a 64-bit value from memory and stores it into **rd**. L[W/H/B] are similar, but they read only the lowest 32/16/8 bits of the value in memory and store it sign-extended to W bits into **rd**. L[W/H/B]U are identical to their signed counterparts but do not perform any sign-extension.
- SD takes a 64-bit operand, y , and stores it into a specified memory location. S[W/H/B] store only the lowest 32/16/8 bits of the operand (without any sign-extension).

These operations are performed using offline memory-checking techniques, as discussed in Section 3.2. They thus do not require a large lookup table and can be handled efficiently using just constraints and simple

range checks in *Lasso*. Note that the memory location involved in these operations is obtained as the sum of the value in `rs1` and `imm`, calculated using constraints just like in the Branch instructions (see Section 5.7).

6 Evaluation Tables for the Multiplication Extension

The M extension adds multiplication, division and remainder operations to the RISC-V ISA. These instructions are generally more complex than the base ones covered so far and involve new techniques in *Jolt* to handle - namely the addition of “virtual” registers and instructions.

6.1 Virtual Instructions and Virtual Registers

Jolt splits certain complex assembly instructions (such as `MULH`) into a sequence of instructions that are executed in the ZKVM in place of the original instruction. The CPU state transition guarantee that should hold for the original assembly instruction now holds after the entire sequence is executed. Note that the splitting of instructions is done in the assembly code during formatting and is independent of the input or even the rest of the program code.

To avoid jumbling with the base registers, *Jolt* introduces new “virtual” registers that virtual instructions use to store intermediate values. These registers have addresses outside the standard set of base registers but are otherwise read from and stored to identically. To ensure safety, the only time the “real” CPU state is changed is when the last virtual instruction of the sequence stores the final result in the “real” destination register of the original instruction.

Reflecting the program counter. It is common for programs to use the value stored in the PC register directly in program logic, such as through Jumps or the `AUIPC` instruction. In these cases, it may be required to maintain two program counters in the ZKVM. The first is the normal one used for keeping tracking of the next address in the VM including virtual instructions. The second is used for reflecting the true value of “real” PC and is not updated during virtual instructions. As the PC is a relatively small value to commit to, this incurs negligible overhead to the prover’s costs and the constraint system. For simplicity, we assume the standard one PC model in the rest of this paper’s discussion but incorporate the cost of a second PC in Section 8.

6.1.1 ASSERT Instructions

Asserts are a type of virtual instruction that add circuit constraints on an instruction’s result. For example, an `ASSERT-[COND]` constraint uses the lookup table for the branch instruction `B[COND]` but additionally adds a constraint that the lookup must return 1. Assert instructions do not have a destination register. On top of the conditional checks seen in the Set-Less-Than and Branch instructions, *Jolt* supports the following assert instructions:

`ASSERT-LT-ABS` takes two W -bit two’s complement signed inputs and outputs $|x| < |y|$.

`ASSERT-EQ-SIGNS` takes two W -bit two’s complement signed inputs and outputs $x_s == y_s$.

OP	INDEX	FULL MLE
ASSERT-LT-ABS	$x \parallel y$	$\text{LTU}(x_{<s}, y_{<s}) \quad // \text{ ignore sign bits}$
ASSERT-EQ-SIGNS	$x \parallel y$	$\widetilde{\text{EQ}}(x_s, y_s)$

6.1.2 ADVICE and MOVE Instructions

`ADVICE v` : stores a special W -bit non-deterministic circuit input into virtual register v .

`MOVE v_1, v_2` : copies the value in register v_1 into register v_2 (either could be virtual).

The advice instruction allows the prover to store non-deterministic advice into virtual registers. The lookup query’s function here is to act as a range check on the advice and thus, uses the range check table. The “non-deterministic” part of these instructions is that their lookup’s query isn’t derived in the circuit (such as through registers, memory or `imm`) but comes from advice passed into the CPU step circuit. Thus, unlike the immediate `imm`, these values aren’t fixed in the assembly code and can be set by the prover at proving time. `ADVICE` has no source register or immediate and only specifies a destination register.

The MLEs of these instructions are identical to the that of the range checks.

OP	INDEX	MLE
ADVICE	x	$\sum_{i=0}^{W-1} 2^i \cdot x_i \quad // \text{ range check}$
MOVE	x	$\sum_{i=0}^{W-1} 2^i \cdot x_i \quad // \text{ range check}$

6.2 The M-Extension Tables

As before, for some new instructions, we give MLE-structured tables for each instruction and describe how it can be decomposed. For other new instructions, we provide the “virtual” sequence of previously-defined instructions that result in the same CPU state change.

6.2.1 Unsigned or Lower Multiplication

The following instructions take two W -bit operands x and y .

`MUL` returns the lower W bits of $x \times y$ where the operands are treated as signed two’s complement numbers.

`MULU` returns the lower W bits of $x \times y$ where the operands are treated as unsigned W -bit numbers.

`MULHU` returns the higher W bits of $x \times y$ where the operands are treated as unsigned W -bit numbers.

Similar to `ADD`, `Jolt` performs the core multiplication operation in the circuit as computing $z = x \times y$ costs just one constraint. The circuit then queries z in the lookup tables of the instructions, which have a row for every possible $2W$ -bit z with the entry being the desired bits. Note that while `MUL` is a signed operation, performing unsigned multiplication returns the same lower bits.

OP	INDEX	FULL TABLE MLE
MUL	$z = x \times y$	$\sum_{i=0}^{W-1} 2^i \cdot z_i \quad // \text{ lower } L \text{ bits}$
MULU	$z = x \times y$	$\sum_{i=0}^{W-1} 2^i \cdot z_i \quad // \text{ lower } L \text{ bits}$
MULHU	$z = x \times y$	$\sum_{i=L}^{2W-1} 2^i \cdot z_i \quad // \text{ higher } L \text{ bits}$

Decomposability. These MLEs can be decomposed in a manner similar to the `AND` table and effectively like the tables for range checks.

6.2.2 Signed and Higher MUL

`MULH` returns the higher W bits of $x \times y$ where the operands are treated as signed two’s complement numbers.

`MULHSU` returns the higher W bits of $x \times y$ where only x is signed but y is unsigned.

These instructions are more complicated than the others as they require signed multiplication which means the operands are sign-extended to $2W$ bits before performing the multiplication. This leads to the result having $4W$ and $3W$ total bits in MULH and MULHSU, respectively. As this is too large to handle with a lookup query, we instead compute the desired bits in stages.

For a number x , let s_x be $\sum_{i=0}^{W-1} 2^i x_s$ such that $[s_x \parallel x]$ is the sign-extension of x to $2W$ bits. The signed multiplication algorithm performs the following $2W \times 2W$ -bit multiplication and returns the highest $2W$ bits: $[s_x \parallel x] \times [s_y \parallel y]$. As the instructions above are only interested in the higher W bits of this result, we can represent the required bits as the *lower* W bits of the sum of the following three values each computed using only unsigned multiplication:

$$[\text{higher } W \text{ bits of } x \times y] + [\text{lower } W \text{ bits of } s_x \times y] + [\text{lower } W \text{ bits of } s_y \times x]$$

Given s_x, s_y , the above terms can be obtained using MULH, MULU instructions and the sum computed using ADD. To get s_x, s_y , we define a new instruction, MOVSIGN, which takes an W -bit input x and stores the W -bit number with x_s as all of its binary coefficients in the destination register.

OP	INPUT	FULL MLE
MOVSIGN	x	$\sum_{i=0}^{W-1} 2^i \cdot x_s \quad // \text{ place sign bit in all positions}$

Decomposability. This table can be **decomposed** naturally using one subtable function.

We can now split MULH, MULHSU into virtual instructions following the above procedure. We use r_x, r_y to denote the two operand registers. We use “v” to name virtual registers. (In actual formatted assembly code, these are replaced by a free numbered virtual register.)

Original	Virtual Sequence (OPCODE, rs1, rs2, imm, rd)
MULH r_x, r_y, rd	<ol style="list-style-type: none"> 1. MOVSIGN $r_x, -, -, v_{s_x}$ <i>// store s_x in a virtual register</i> 2. MOVSIGN $r_y, -, -, v_{s_y}$ <i>// store s_y</i> 3. MULHU $r_x, r_y, -, v_0$ <i>// get higher bits of $x \times y$</i> 4. MULU $v_{s_x}, r_y, -, v_1$ <i>// get lower bits of $s_x \times y$</i> 5. MULU $v_{s_y}, r_x, -, v_2$ <i>// get lower bits of $s_y \times x$</i> 6. ADD $v_0, v_1, -, v_3$ 7. ADD $rd, v_2, -, rd$
MULH r_x, r_y, rd	<ol style="list-style-type: none"> 1. MOVSIGN $r_x, -, -, v_{s_x}$ 2. MULHU $r_x, r_y, -, v_1$ 3. MULU $v_{s_x}, v_y, -, v_2$ 4. ADD $v_1, v_2, -, rd$

The correctness of the output can be seen by inspection as the steps follow the natural binary multiplication algorithm. It can also be seen that the “real” CPU state is only modified in the final steps of each sequence, when the result is stored into **rd**.

6.3 Division and Remainder

In RISC-V, division and remainder operations take two W -bit values read from registers. In Jolt, for both operations, the prover provides as non-deterministic advice the quotient q and remainder r using the ADVICE instruction introduced in Section 6.1.2. The correctness of this advice is verified using a sequence of virtual

instructions, as shown below. As both DIV and REM instructions perform the same checks, they have nearly identical virtual instructions with only the last instruction differing based on the desired value (q or r).

Unsigned versions. In unsigned division, both operands x, y and quotient q and remainder r are all treated as unsigned W -bit numbers. DIVU/REMU require x to be equal to $q \times y + r$ such that $r < y$ and $q \times y \leq x$.

Original	Virtual Sequence (OPCODE, rs1, rs2, imm, rd)
DIVU r_x, r_y, rd	<ol style="list-style-type: none"> 1. ADVICE $-, -, -, v_q$ // store non-deterministic advice q into v_q 2. ADVICE $-, -, -, v_r$ // store non-deterministic advice r into v_r 3. MULU $v_q, r_y, -, v_{qy}$ // compute $q \times y$ 4. ASSERT_LTU $v_r, r_y, -, -$ // verify that $r < y$ 5. ASSERT_LTE $v_{qy}, r_x, -, -$ // assert $q \times y \leq x$ 6. ADD $v_{qy}, v_r, -, v_0$ // compute $q \times y + r$ 7. ASSERT_EQ $v_0, r_x, -, -$ 8. MOVE $v_q, -, -, rd$ // store q in rd
REMU r_x, r_y, rd	<ol style="list-style-type: none"> 1-7. same as above 8. MOVE $v_r, -, -, rd$ // store r in rd

Signed versions. In signed division, both operands x, y and quotient q and remainder r are all treated as signed 2's complement W -bit numbers.

DIVU/REMU requires $x = q \times y + r$ such that $|r| < |y|$ and r, y have the same sign.

Original	Virtual Sequence
DIV r_x, r_y, rd	<ol style="list-style-type: none"> 1. ADVICE $-, -, -, v_q$ // store non-deterministic advice q into v_q 2. ADVICE $-, -, -, v_r$ // store non-deterministic advice r into v_r 3. ASSERT_LT_ABS $v_r, r_y, -, -$ // verify that $r < y$ 4. ASSERT_EQ_SIGNS $v_r, r_y, -, -$ // require r to have the sign of y 5. MUL $v_q, r_y, -, v_{qy}$ // compute $q \times y$ 6. ADD $v_{qy}, v_r, -, v_0$ // compute $q \times y + r$ 7. ASSERT_EQ $v_0, x, -, -$ 8. MOVE $v_q, -, -, rd$ // store q in rd
REM r_x, r_y, rd	<ol style="list-style-type: none"> 1-7. same as above 8. MOVE $v_r, -, -, rd$ // store r in rd

As with the splitting of the multiplication instructions, the correctness of the output can be seen by inspection as the steps follow the straightforward verification of division. Also, the real CPU state is only modified in the final steps.

7 Putting It all Together: a SNARK for RISC-V Emulation

The overall architecture of Jolt is depicted in Figure 2. The prover begins by cryptographically committing to the execution trace z of the VM on the appropriate input (or more precisely, its multilinear extension polynomial \tilde{z} , using any multilinear polynomial commitment scheme).

The R1CS instance and the cost of proving its satisfaction. It is straightforward to identify R1CS constraint matrices A, B, C (see Definition 2.7) such that z is a valid execution trace if and only if z satisfies $Az \circ Bz = Cz$ assuming z satisfies memory-consistency and relevant entries of z are indeed in the relevant lookup tables capturing evaluation of the RISC-V instructions. The constraint system is sketched in Appendix A.3. A preliminary implementation of the constraint system in Zokrates [ET18] requires about 250 lines of code. We apply Spartan to establish that z satisfies the constraint system.

The R1CS constraint matrices A, B, C are uniform, in the sense that the multilinear extensions \tilde{A}, \tilde{B} , and \tilde{C} can each be evaluated in $O(\log T)$ field operations, where T is the number of steps the prover runs the RISC-V program. Because of this, the Spartan prover need only commit cryptographically to \tilde{z} , and in particular no party has to commit to \tilde{A}, \tilde{B} , or \tilde{C} . The constraint system has several hundred constraints per step of the RISC-V CPU. When Spartan is applied to this R1CS instance, the number of field operations performed by the prover is a small constant factor larger than the number of constraints (rows of the constraint matrices) plus the number of variables (columns of the constraint matrices). Since a field operation is at least an order of magnitude faster than a group operation, if an MSM-based polynomial commitment scheme is used, this implies that the Spartan prover time is dominated by the polynomial commitment costs summarized in Section 8.

Memory checking. Memory checking is handled by the memory-checking argument from Spice [SAGL18] sketched in Appendix B.3. Details specific to the memory model of RISC-V are discussed in Appendix A.2. The cost of the memory-checking argument are included in Section 8.

Lookup argument. The Lasso family of lookups is used to confirm that the relevant entries of z are in the relevant lookup tables. The remaining issue to address is that Sections 4.3-6 explained that the evaluation table of each individual RISC-V instruction is both MLE-structured and decomposable. But Lasso is a lookup argument for a single decomposable table (and Generalized-Lasso is a lookup argument for a single MLE-structured table). The following subsections explain how to bridge this gap.

7.1 The case of Generalized-Lasso

This gap is simple to bridge if using Generalized-Lasso rather than Lasso. This is because the concatenation of several MLE-structured tables is also MLE-structured. So we just apply Generalized-Lasso to the concatenation of the evaluation tables of every RISC-V instruction.

Indeed, suppose that the lookup table is the union of ℓ different tables, where assume for simplicity that each table has size N . For $i = 1, \dots, \ell$, let $T_i \in \mathbb{F}^N$ denote the vector representing table i . and let $T \in \mathbb{F}^{\ell \cdot N}$ denote the concatenation of these vectors. Then viewing T as a function mapping $\{0, 1\}^{\log \ell} \times \{0, 1\}^{\log N} \rightarrow \mathbb{F}$, in the natural way, the following equation holds, with $x \in \{0, 1\}^{\log \ell}$ and $y \in \{0, 1\}^{\log N}$:

$$\tilde{T}(x, y) = \sum_{k \in \{0, 1\}^{\log \ell}} \tilde{E}Q(k, x) \cdot \tilde{T}_{\text{int}(k)}(y), \quad (10)$$

where recall that $\text{int}(k) = \sum_{i=1}^{\ell} 2^{i-1} \cdot k_i$ denotes the integer of which k is the binary representation. This is because the right hand side of Equation (10) is a multilinear polynomial in the variables of x and y and agrees with T at all inputs in $\{0, 1\}^{\log \ell} \times \{0, 1\}^{\log N}$. Hence, it is the unique multilinear polynomial extending T . Moreover, Equation (10) can be evaluated at any point $(r_1, r_2) \in \mathbb{F}^{\log \ell} \times \mathbb{F}^{\log N}$ in time $O(\log \ell)$ plus the amount of time required to evaluate each of T_1, \dots, T_ℓ at r_2 .

7.2 The case of Lasso

Addressing the gap discussed above is more involved if using Lasso. Closely related issues have been addressed in earlier work on zkVMs.

Specifically, the fact that the evaluation tables of different RISC-V instructions have different decompositions into subtables is analogous to the following issue dealt with in earlier approaches to front-end design for

zkVMs: different instructions are computed by different circuits. and while only one instruction is executed per step of the VM, in general it is not known until runtime which instruction will be executed at any given step. Very early front-ends for zkVMs such as TinyRAM dealt with this naively, by producing circuits that contained logic to execute *every possible instruction* for every step of the VM [BSCG⁺13a, BSCG⁺13b].

A reordering approach. vRAM [ZGK⁺18] proposed to avoid this overhead in a manner similar to the sorting-based memory-checking arguments sketched in Appendix B.2 (see also Arya [BCG⁺18] for related techniques). We sketch vRAM’s approach as an interactive argument, but it can be rendered non-interactive via the Fiat-Shamir transformation.

Specifically, after running the VM on the appropriate input, the prover tells the verifier how many times each of the primitive instructions was executed (if there are ℓ primitive instructions, call these numbers k_1, \dots, k_ℓ), and commits to a reordering of the execution trace in which the entries are not sorted by timestep, but rather grouped by which instruction is executed at that step. That this second execution trace is indeed a reordering of the first is confirmed via the standard randomized permutation-checking procedure described in Appendix B.2 (Equation (11)). The instruction counts k_1, \dots, k_ℓ uniquely specify a circuit \mathcal{C} that takes as input the reordered execution trace, and confirms that at each step the appropriate instruction was executed correctly. That is, \mathcal{C} contains circuitry executing instruction one k_1 times, instruction two k_2 times, and so on.

Any SNARK can then be applied to prove the reordered execution trace indeed passes all of the checks computed within \mathcal{C} . Conceptually, this approach uses the prover’s knowledge (after running the program) of which instruction was executed at each step to eliminate the extraneous circuitry in earlier work (devoted to executing instructions that were not actually executed by the VM). See [ZGK⁺18] for further details.

An identical approach applies to the use of Lasso in Jolt. After the prover tells the verifier k_1, \dots, k_ℓ and commits to a reordered execution trace grouped by executed instruction, Lasso is applied ℓ different times, i.e., to the first k_1 entries of the reordered execution trace using the evaluation table of the first instruction, to the second k_2 entries of the reordered execution trace using the evaluation table of the second instruction, and so on. Via standard batching techniques [Set20], this does not increase the verification costs relative to a single invocation of Lasso.

Avoiding reordering. Cairo and other zkVM projects address the issue by including constraints that do capture all instructions, but using a SNARK targeted at uniform computation that ensures the prover’s cryptographic work does not grow with the number of constraints (i.e., the prover only cryptographically commits to the solution vector to the constraint system, whose length is independent of the number of constraints. Note, however, that the number of field operations performed by the prover does grow with the number of constraints). We can take an analogous approach with Lasso.

Conceptually, this approach expresses the concatenation of the evaluation tables of each instruction (and of which we have shown to be decomposable) as itself decomposable, analogous to how the concatenation of MLE-structured tables is itself MLE-structured (Equation (10)).

To this end, it is convenient to treat each instruction as leading to $2c - 1$ different lookups into subtables ($2c - 1$ here comes from the maximum number of subtable lookups across all instructions, namely due to SLTU as described in Section 5.3). For instructions that require fewer than $2c - 1$ lookups into subtables, the extraneous lookup results can be set to 0, thereby avoiding any cryptographic work on the part of the prover if using an MSM-based commitment scheme. (we will explain below how to ensure that these extraneous subtable lookup results will be ignored by all subtables).

There will be a single collation polynomial g (Definition 2.6) for all instructions, but g will take as input not only the results of relevant subtable lookups, but also 8 additional variables that, when assigned values in $\{0, 1\}$, are interpreted as the bit-representation of the opcode. Denoting these variables as $w = (w_1, \dots, w_8)$, and letting $g_i(z)$ denote the collation polynomial for the i ’th instruction, and letting x denote a vector of $2c - 1$ variables, interpreted as specifying the results of $2c - 1$ subtable lookups, we define

$$g(w, x) = \sum_{y \in \{0, 1\}^8} \widetilde{\text{EQ}}(w, y) \cdot g_{\text{int}(y)}(x).$$

This definition ensures that for any instruction i , $g(\text{bin}(i), x) = g_i(x)$, i.e., collation for each instruction is performed correctly by g .

Ensuring that subtable lookups get “routed” to the correct subtables. As sketched in Appendix B.3, the core of **Lasso** is to invoke a grand product argument for each subtable (in order to compute the quantities given in Expressions (12) and (13) of Appendix B.3, to confirm they are equal). Applying a highly-optimized variant of the GKR protocol to a circuit computing this product [Tha13] yields a grand product argument in which the prover only performs field operations (the number of which is linear in the circuit size). Other related grand product arguments reduce the proof size at the cost of a low-order increase in commitment costs for the prover.

We can modify the circuit used for each subtable to take as input the bits of the opcode associated with each lookup. This way, the circuit can simply ignore any lookups associated with opcodes that do not access the subtable associated with the circuit. By ignore, we mean that the circuit sets to 1 the factor in the size- m product within Expressions (12) and (13) that is associated with this particular lookup.

To minimize the size of this circuit, rather than having the prover commit to the bits of the opcode, it may be preferable to instead have the prover commit to some additional Boolean flags (beyond the 14 Boolean flags already described in Appendix A, so that each subtable’s circuit only needs to inspect a handful of Boolean flags to determine whether or not a given lookup operation actually is intended to access the subtable.

7.3 Pros and cons of using **Lasso** vs. **Generalized-Lasso** within **Jolt**

We estimate that using **Lasso** in **Jolt** reduces commitment costs for the prover by $2 \times - 3 \times$, owing to **Generalized-Lasso**’s need for the prover to commit to c random field elements per lookup, where c is a parameter.²¹ **Generalized-Lasso** has the prover commit to the same number of field elements as **Lasso**, but all of them are small, meaning in the set $\{0, 1, \dots, m\}$ where m is the number of lookups.

However, a benefit of **Generalized-Lasso** is that it encapsulates inside of it all of the complexity regarding different instructions having different subtables. That is, for **Generalized-Lasso**, the prover and verifier can be implemented using only an oracle (i.e., function) that takes as input a point $r \in \mathbb{F}^\ell$ where ℓ is the logarithm of the evaluation table size, and evaluates the multilinear extension of each instruction’s evaluation table at a required point. This makes specifying the evaluation table of each instruction extremely simple. In contrast, with **Lasso**, the verifier needs to know how each evaluation table is decomposed into subtables, and how the results of lookups into subtables are collated into the result of the associated lookup into the “big” evaluation table.

One challenge here is that **Generalized-Lasso** involves an additional invocation of the sum-check protocol [LFKN90], and if the prover is only given an evaluation oracle for the multilinear extension of an instruction, it may require $O(m \log N)$ oracle queries for m lookups into a table of size N [CMT12]. Since, for all the evaluation tables used in **Jolt**, each oracle query can be answered in $O(\log N)$ field operations this translates to $O(m \log^2 N)$ field operations total, and the $\log^2 N$ factor may be large enough for this to bottleneck the prover. Fortunately, by leveraging more information about the lookup tables arising in **Jolt** than that their multilinear extensions can be evaluated in $O(\log N)$ time, it is in fact possible to implement the sparse-dense sum-check prover in $O(m \log N)$ or even $O(cm)$ field operations (our companion paper [STW23] works through the details of this for several of the lookup tables arising in **Jolt**). These sum-check prover implementations are more complicated than the $O(m \log^2 N)$ -time prover that only uses an evaluation oracle for the multilinear extensions. However, this does not affect auditability because only the **Jolt** verifier needs to be implemented correctly for the SNARK to be secure (and indeed, the verifier itself acts as an automatic auditor of the prover implementation).

In summary, while prover commitment costs are higher if **Jolt** is built using **Generalized-Lasso** as the lookup argument rather than **Lasso**, there are substantial auditability benefits to using **Generalized-Lasso** that may justify the performance price.

²¹Reasonable values of c are $c = 3$ or $c = 4$ for 32-bit data types, ensuring subtables of size 2^{22} or 2^{16} and $c = 6$ or $c = 8$ for 64-bit data types.

Per-step Commitment Costs for Non-Memory Operations

Bit-length	Number of Elements	In RV32	In RV64
		$W = 32, c = 3$	$W = 64, c = 6$
1	22	22	22
$[2, 12]$	$3 + 2c$	9	15
$(2W/c) \approx 22$	$1 + c$	4	7
$\log(T)$	$4 + 2c$	10	16
W	5	5	5
Total Elements	$35 + 5c$	50	65
In 256-bit equivalents:		≈ 5 elements	≈ 6 elements

Figure 3: An overview of the spread of elements committed to in **Jolt** in non-memory operations (i.e., excluding loads and stores which do not involve lookups) by their bit-length. The **Lasso** parameter c is 3 with RV32 ($W = 32$) and $c = 6$ with RV64 ($W = 64$). We approximate the per-step commitments costs in terms of the cost of committing to a 256-bit element when using Pippenger’s MSM algorithm, assuming that the program code is under 2^{22} bytes long, and the program finishes in under 2^{30} CPU steps.

8 Qualitative Cost Estimation

In this section, we provide a qualitative evaluation of the prover cost involved in **Jolt** when using the **Lasso** lookup argument. As discussed in Section 1.3.1, the dominating cost is in producing commitments to the elements fed to the constraint system and lookup argument every step. Table 3 provides an upper bound on the elements committed per step grouped by their bit-lengths. We analyze bit-length as that is the main factor determining the commitment cost when using Pippenger’s multi-scalar multiplication algorithm.²² We provide below a brief overview of the elements involved and leave a more detailed discussion to Appendix A.

Elements involved in CPU Execution. First, let’s look at the elements involved in satisfying the CPU step circuit’s constraints before looking at the elements needed for the **Lasso** argument. The smallest of these are the 1-bit ones constituting the bits of `opflags`[14] and `opcode`[8] and the 5-bit elements indexing the source (`rs1`, `rs2`) and destination (`rd`) registers read from the instruction. Slightly larger elements are the PC (which could be as large as $\log |\text{program_code}|$ bits) and the step counter. The latter starts small but grows over time: for example, it can reach 25 bits for a reasonably long execution of 30 million steps. Finally, the largest elements involved are the W -bit ones specifying the values stored in the two source registers, the sign-extended `imm` read from the program code, the lookup output, (which is generally stored in the destination register), and the advice element involved (only) in division and reminder operations. Note that most instructions do not involve all these elements (notably, the advice element is used rarely and the ones that do use advice never use `imm`) and thus, the numbers from Table 3 are a worst-case upper bound.

Elements involved in Lasso The cost of a lookup query with **Lasso** was discussed in Section 4.3. In brief, most lookups require the prover to commit to three groups of c elements each with the groups having bit-lengths W/c , $2W/c$ and $\log(T)$, respectively. The lookup costs figuring into Table 3 are for the worst-case scenario of lookups involving the less-than comparison table of Section 5.3 which commits to $2c$ more elements than normal. Thus, most RISC-V instructions are actually slightly cheaper than the costs reported. In fact, many, like **ADD** and those related to it, require fewer than $3c$ elements to be committed as their queries are smaller (only $W + 1$ bits long).

As a quick note, using **Generalized-Lasso** for lookups would increase the bit complexity of c of the committed

²²In Pippenger’s multi-scalar multiplication algorithm to commit to elements, committing to an N -bit element costs roughly $\text{ceil}(N/22)$ group operations. This makes committing to a 32-bit element cost two group operations while a 256-bit element costs 12 group operations.

Base Costs per Memory Instruction		Overhead per Byte	
Bit-length	Number of Elements	for Loads	for Stores
1	22	1	1
[2, 8]	3	1	1
$\log(T)/2$	3	4	4
$\log(T)$	5	6	7
W	3	-	-
Total Elements	36	12	13
In 256-bit equivalents (both RV32, RV64)	3.5 elements	1.5 elements	1.5 elements

Figure 4: The spread of elements committed to per memory operation with the extra overhead elements per byte of load or store. That is, a load and store of k bytes involves the prover committing to $36 + 12k$ elements, and $36 + 13k$ elements, respectively. We approximate the per-step commitments costs in terms of the cost of committing to a 256-bit element when using Pippenger’s MSM algorithm, assuming that the program code is under 2^{22} bytes long (placing it in the $\log(T)$ category), and the program finishes in under 2^{30} CPU steps.

field elements from $2W/c$ to $\log|\mathbb{F}|$.

8.1 Cost of Memory Operations

We analyze the cost of load and store operations separately as these do not involve large lookups to perform the core instruction logic. Rather, the main cost here is performing memory-checking operations, one for each byte of memory involved in the load/store. This can be up to four for 32-bit processors and eight for 64-bit processors. The elements involved on top of the non-lookup elements of the non-memory instructions are the actual bytes read/written, the timestamps involved in memory-checking, and the cost of range-checks and computing max function (via small lookup tables) with the timestamps. Memory operations also commit to two fewer W -bit elements as they never involve the advice (used only in division and remainder operations) and lookup output elements.

Firstly, a minor cost involved in memory operations are extra 1-bit memory “flags” which act as a unary vector indicating the exact number of bytes read/stores. For both loads/stores, for each byte read/written, the prover commits to the actual byte value and provides a timestamp `ts_read` to be used in offline memory-checking (indicating when the byte was last written to). This must be range-checked and verified to be less than the current step counter. This range-check is very efficient in **Lasso** and requires committing to a single element of value bounded by the number of steps up to that point. Verifying that `ts_read` is less than the current step counter uses the less-than lookup table (of Section 5.3) and requires committing to four elements of bit-length bounded by that of the step counter and four more of half those many bits.

Stores, which are memory “writes”, require 8-bit range checks of the bytes written. These range-checks are again very efficient in **Lasso** and only involve committing to a single element of value at most the number of steps up to that point.

Acknowledgements and Disclosures

Justin Thaler was supported in part by NSF CAREER award CCF-1845125 and by DARPA under Agreement No. HR00112020022. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the United States Government or DARPA.

Disclosures. Thaler is a Research Partner at a16z crypto and is an investor in various blockchain-based platforms, as well as in the crypto ecosystem more broadly (for general a16z disclosures, see <https://www.a16z.com/disclosures/>.)

References

- [AGL⁺22] Jeremy Avigad, Lior Goldberg, David Levit, Yoav Seginer, and Alon Titelman. A verified algebraic representation of cairo program execution. In *Proceedings of the 11th ACM SIGPLAN International Conference on Certified Programs and Proofs*, pages 153–165, 2022.
- [And17] Andrew Waterman¹, Krste Asanovic. The RISC-V instruction set manual. <https://riscv.org/wp-content/uploads/2017/05/riscv-spec-v2.2.pdf>, 2017.
- [BBB⁺18] Benedikt Bünz, Jonathan Bootle, Dan Boneh, Andrew Poelstra, Pieter Wuille, and Greg Maxwell. Bulletproofs: Short proofs for confidential transactions and more. In *Proceedings of the IEEE Symposium on Security and Privacy (S&P)*, 2018.
- [BBHR18] Eli Ben-Sasson, Iddo Bentov, Yinon Horesh, and Michael Riabzev. Fast Reed-Solomon interactive oracle proofs of proximity. In *Proceedings of the International Colloquium on Automata, Languages and Programming (ICALP)*, 2018.
- [BCC⁺16] Jonathan Bootle, Andrea Cerulli, Pyrros Chaidos, Jens Groth, and Christophe Petit. Efficient zero-knowledge arguments for arithmetic circuits in the discrete log setting. In *Proceedings of the International Conference on the Theory and Applications of Cryptographic Techniques (EUROCRYPT)*, 2016.
- [BCG⁺18] Jonathan Bootle, Andrea Cerulli, Jens Groth, Sune Jakobsen, and Mary Maller. Arya: Nearly linear-time zero-knowledge proofs for correct program execution. In *Proceedings of the International Conference on the Theory and Application of Cryptology and Information Security (ASIACRYPT)*, 2018.
- [BCHO22] Jonathan Bootle, Alessandro Chiesa, Yuncong Hu, and Michele Orru. Gemini: Elastic snarks for diverse environments. In *Proceedings of the International Conference on the Theory and Applications of Cryptographic Techniques (EUROCRYPT)*, 2022.
- [BDFG20] Dan Boneh, Justin Drake, Ben Fisch, and Ariel Gabizon. Halo Infinite: Recursive zk-SNARKs from any Additive Polynomial Commitment Scheme. Cryptology ePrint Archive, Report 2020/1536, 2020.
- [BEG⁺91] Manuel Blum, Will Evans, Peter Gemmell, Sampath Kannan, and Moni Naor. Checking the correctness of memories. In *Proceedings of the IEEE Symposium on Foundations of Computer Science (FOCS)*, 1991.
- [BFR⁺13] Benjamin Braun, Ariel J. Feldman, Zuocheng Ren, Srinath Setty, Andrew J. Blumberg, and Michael Walfish. Verifying computations with state. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP)*, 2013.
- [BFS20] Benedikt Bünz, Ben Fisch, and Alan Szepieniec. Transparent SNARKs from DARK compilers. In *Proceedings of the International Conference on the Theory and Applications of Cryptographic Techniques (EUROCRYPT)*, 2020.
- [BGH19] Sean Bowe, Jack Grigg, and Daira Hopwood. Recursive proof composition without a trusted setup. Cryptology ePrint Archive, Report 2019/1021, 2019.

- [BGtR23] Jeremy Bruestle, Paul Gafni, and the RISC Zero Team. RISC Zero zkVM: Scalable, transparent arguments of RISC-V integrity, 2023.
- [BSBHR19] Eli Ben-Sasson, Iddo Bentov, Yinon Horesh, and Michael Riabzev. Scalable zero knowledge with no trusted setup. In *Advances in Cryptology—CRYPTO 2019: 39th Annual International Cryptology Conference, Santa Barbara, CA, USA, August 18–22, 2019, Proceedings, Part III* 39, pages 701–732. Springer, 2019.
- [BSCG⁺13a] Eli Ben-Sasson, Alessandro Chiesa, Daniel Genkin, Eran Tromer, and Madars Virza. Snarks for c: Verifying program executions succinctly and in zero knowledge. In Ran Canetti and Juan A. Garay, editors, *Advances in Cryptology – CRYPTO 2013*, pages 90–108, Berlin, Heidelberg, 2013. Springer Berlin Heidelberg.
- [BSCG⁺13b] Eli Ben-Sasson, Alessandro Chiesa, Daniel Genkin, Eran Tromer, and Madars Virza. Tinyram architecture specification, v0. 991. *en. In: (Aug. 2013)*, page 16, 2013.
- [BSCGT13] Eli Ben-Sasson, Alessandro Chiesa, Daniel Genkin, and Eran Tromer. Fast reductions from rams to delegatable succinct constraint satisfaction problems. In *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*, pages 401–414, 2013.
- [BSCTV14] Eli Ben-Sasson, Alessandro Chiesa, Eran Tromer, and Madars Virza. Scalable zero knowledge via cycles of elliptic curves. In *Proceedings of the International Cryptology Conference (CRYPTO)*, 2014.
- [CBBZ23] Binyi Chen, Benedikt Bünz, Dan Boneh, and Zhenfei Zhang. HyperPlonk: Plonk with linear-time prover and high-degree custom gates. In *Proceedings of the International Conference on the Theory and Applications of Cryptographic Techniques (EUROCRYPT)*, 2023.
- [CHM⁺20] Alessandro Chiesa, Yuncong Hu, Mary Maller, Pratyush Mishra, Noah Vesely, and Nicholas Ward. Marlin: Preprocessing zkSNARKs with universal and updatable SRS. In *Proceedings of the International Conference on the Theory and Applications of Cryptographic Techniques (EUROCRYPT)*, 2020.
- [CMT12] Graham Cormode, Michael Mitzenmacher, and Justin Thaler. Practical verified computation with streaming interactive proofs. In *Proceedings of the Innovations in Theoretical Computer Science (ITCS)*, 2012.
- [ET18] Jacob Eberhardt and Stefan Tai. Zokrates - scalable privacy-preserving off-chain computations. In *2018 IEEE International Conference on Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom) and IEEE Smart Data (SmartData)*, pages 1084–1091, 2018.
- [FS86] Amos Fiat and Adi Shamir. How to prove yourself: Practical solutions to identification and signature problems. In *Proceedings of the International Cryptology Conference (CRYPTO)*, pages 186–194, 1986.
- [GKR08] Shafi Goldwasser, Yael Tauman Kalai, and Guy N. Rothblum. Delegating computation: Interactive proofs for muggles. In *Proceedings of the ACM Symposium on Theory of Computing (STOC)*, 2008.
- [GPR21] Lior Goldberg, Shahar Papini, and Michael Riabzev. Cairo—a turing-complete stark-friendly cpu architecture. *Cryptology ePrint Archive*, 2021.
- [Gro16] Jens Groth. On the size of pairing-based non-interactive arguments. In *Proceedings of the International Conference on the Theory and Applications of Cryptographic Techniques (EUROCRYPT)*, 2016.
- [GWC19] Ariel Gabizon, Zachary J. Williamson, and Oana Ciobotaru. PLONK: Permutations over Lagrange-bases for oecumenical noninteractive arguments of knowledge. ePrint Report 2019/953, 2019.

- [KST22] Abhiram Kothapalli, Srinath Setty, and Ioanna Tzialla. Nova: Recursive Zero-Knowledge Arguments from Folding Schemes. In *Proceedings of the International Cryptology Conference (CRYPTO)*, 2022.
- [KT23] Tohru Kohrita and Patrick Towa. Zeromorph: Zero-knowledge multilinear-evaluation proofs from homomorphic univariate commitments. *Cryptology ePrint Archive*, 2023.
- [KZG10] Aniket Kate, Gregory M. Zaverucha, and Ian Goldberg. Constant-size commitments to polynomials and their applications. In *Proceedings of the International Conference on the Theory and Application of Cryptology and Information Security (ASIACRYPT)*, pages 177–194, 2010.
- [Lee21] Jonathan Lee. Dory: Efficient, transparent arguments for generalised inner products and polynomial commitments. In *Theory of Cryptography Conference*, pages 1–34. Springer, 2021.
- [LFKN90] Carsten Lund, Lance Fortnow, Howard Karloff, and Noam Nisan. Algebraic methods for interactive proof systems. In *Proceedings of the IEEE Symposium on Foundations of Computer Science (FOCS)*, October 1990.
- [Lip89] Richard J Lipton. *Fingerprinting sets*. Princeton University, Department of Computer Science, 1989.
- [PST13] Charalampos Papamanthou, Elaine Shi, and Roberto Tamassia. Signatures of correct computation. In *Theory of Cryptography Conference*, pages 222–242. Springer, 2013.
- [SAGL18] Srinath Setty, Sebastian Angel, Trinabh Gupta, and Jonathan Lee. Proving the correct execution of concurrent services in zero-knowledge. In *Proceedings of the USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, October 2018.
- [Set20] Srinath Setty. Spartan: Efficient and general-purpose zkSNARKs without trusted setup. In *Proceedings of the International Cryptology Conference (CRYPTO)*, 2020.
- [SL20] Srinath Setty and Jonathan Lee. Quarks: Quadruple-efficient transparent zkSNARKs. *Cryptology ePrint Archive*, Report 2020/1275, 2020.
- [STW23] Srinath Setty, Justin Thaler, and Riad S. Wahby. Lasso: Unlocking the lookup singularity, 2023.
- [Tha13] Justin Thaler. Time-optimal interactive proofs for circuit evaluation. In *Proceedings of the International Cryptology Conference (CRYPTO)*, 2013.
- [Tha22] Justin Thaler. Proofs, arguments, and zero-knowledge. *Foundations and Trends in Privacy and Security*, 4(2–4):117–660, 2022.
- [Tom23] Tomer Solberg. RISC Zero prover protocol & analysis. https://github.com/ingonyama-zk/papers/blob/main/risc0_protocol_analysis.pdf, 2023.
- [Whi] Barry Whitehat. Lookup singularity. <https://zkresearch.ch/t/lookup-singularity/65/7>.
- [WSR⁺15] Riad S. Wahby, Srinath Setty, Zuocheng (Andy) Ren, Andrew J. Blumberg, and Michael Walfish. Efficient RAM and control flow in verifiable outsourced computation. In *NDSS*, 2015.
- [WTS⁺18] Riad S. Wahby, Ioanna Tzialla, Abhi Shelat, Justin Thaler, and Michael Walfish. Doubly-efficient zkSNARKs without trusted setup. In *Proceedings of the IEEE Symposium on Security and Privacy (S&P)*, 2018.
- [ZGK⁺18] Yupeng Zhang, Daniel Genkin, Jonathan Katz, Dimitrios Papadopoulos, and Charalampos Papamanthou. vRAM: Faster verifiable RAM with program-independent preprocessing. In *Proceedings of the IEEE Symposium on Security and Privacy (S&P)*, 2018.
- [ZXZS20] Jiaheng Zhang, Tiancheng Xie, Yupeng Zhang, and Dawn Song. Transparent polynomial delegation and its applications to zero knowledge proof. In *Proceedings of the IEEE Symposium on Security and Privacy (S&P)*, 2020.

A The Jolt Elements and Constraints in More Detail

The table below lists the elements committed by the prover per CPU step for instructions that are not loads or stores (see Table 2 for the elements involved there).

ELEMENT	PURPOSE	#SIG. BITS
opflags[14]	These are 1-bit elements used to guide the constraint system on if-else branches. (List given below.)	1 bit \times 14
opcode[8]	These bits constitute the 8-bit opcode for the instruction.	1 bit \times 8
rs1, rs2, rd	Indices of the step instruction’s source and destination registers.	5 bits \times 3
PC	The program counter (possible two for certain programs with virtual instructions).	$\log(\text{code}) \times 2$
step_counter	The global timestamp incremented every step.	$\log(\#\text{steps})$ bits
read_ts_code	The timestamp passed as advice for memory-checking when reading the program code.	$\log(\#\text{steps})$ bits
read_ts_rs1, read_ts_rs2	The timestamps passed as advice for memory-checking when reading the values at registers rs1 , rs2 .	$\log(\#\text{steps})$ bits \times 2
imm	The step instruction’s immediate, appropriately sign or zero extended.	W bits
Values read at rs1, rs2	The actual values read from registers rs1 , rs2 .	W bits \times 2
Lookup result	The instruction’s output passed in as advice.	W bits
Extra advice	The non-deterministic advice element used only in division and remainder operations.	W bits
Elements involved in Lasso lookups.	Commitment to subtable outputs	10 bits \times 2c
	Commitment to the chunks $C[c]$ of the lookup query.	22-bits \times c
	Commitment to the step counter	$\log(\#\text{steps})$ bits \times 2c

Table 1: The basic elements involved in most instructions.

A.1 List of Operation Flags employed:

1. This flag is 0 if the first operand is **rs1** and 1 if it is **rs2**.
2. This flag is 0 if the second operand is **rs2** and 1 if it is **imm**.
3. Is this a load instruction?
4. Is this a store instruction?
5. Is this a jump instruction?
6. Is this a branch instruction?
7. Does this instruction update **rd**?
8. Does this instruction involve adding the operands?

9. Does this instruction involve subtracting the operands?
10. Does this instruction involve multiplying the operands?
11. Does this instruction involve non-deterministic advice?
12. Does this instruction assert the lookup output to be false?
13. Does this instruction assert the lookup output to be true?
14. This flag is the sign bit of the immediate.

The memory flags are as follows: if the instruction is a load/store operation that reads/writes k bytes, then the memory flags will be of the form $1^k \parallel 0^{W/8-k}$.

A.2 Supporting byte-addressable loads and stores

One challenge with implementing a zkVM for RISC-V is supporting byte-addressable memory. This requires performing up to $W/8$ memory operations per load/store, one for each byte written or read. This section describes the memory-checking steps involved. Large lookup tables are not involved in these instructions.

Stores. Each store instruction reads the lower k -byte-suffix from **rs2** ($k = 8, 4, 2, 1$ for instructions SD, SW, SH, SB, respectively) and writes the result into memory locations starting from $\text{loc} = \text{rs1} + \text{imm}$ (see Section 5.8 for how this is calculated). There are two steps involved in stores:

1. The prover provides as advice the bytes-decomposition of the value in **rs2**. These values are range-checked and verified to be the correct decomposition.
2. Memory-checking can then write these bytes one by one to their memory locations (which starts at $\text{rs1} + \text{imm}$). Offline memory-checking requires that the prover provide a timestamp of the latest write that occurred to that memory location. This timestamp must be range-checked and verified to be less than the current timestamp.

Loads. The load operations take k bytes of memory ($k = 8, 4, 2, 1$ for instructions LD, LW, LH, LB, respectively) from location $\text{loc} = \text{rs1} + \text{imm}$, sign-extends it to W bits and then stores the result in **rd**.

1. The k bytes are first read using memory-checking. As with stores, range checks and less-than comparisons are performed on the timestamps provided. not required here, as they were enforced during the stores.
2. Jolt employs a small lookup table to get the sign-bit of the highest order byte. Sign-extension and concatenation can then be performed using constraints.

Range-check costs. The CPU circuit requires performing range-checks on inputs of the following bit-sizes: 8 bits (for the bytes in stores) and $\log(\#\text{steps})$ bits (for the timestamps). Both these range checks can be performed using **Lasso** with parameter $c = 1$. This is a special case that requires the prover to commit to only one element of value bounded by the step counter.

Table 2 shows the overheads (on top of the basic non-lookup elements involved in all operations) for load and store operations per byte involved in the operation (up to 4 for RV32 and 8 for RV64). These elements are always 0 in other operations and hence only count towards the prover's cost when performing loads and stores.

Element	Purpose	#Bits per byte
<code>memflags</code>	Unary vector indicating the byte positions that are active in the load or store.	1 bit
<code>memory_v.byte</code>	The actual byte value read/stored	8 bits
<code>memory_timestamp</code>	The timestamps passed as advice for memory-checking.	$\log(\#steps)$ bits
<code>ts_range_check</code>	Element involved in the range checks of <code>memory_ts[k]</code> performed with Lasso	$\log(\#steps)$ bits
<code>max_checks</code>	Elements involved in verifying that the read timestamp is less than the step counter.	3 bits $\times 2$
		5 bits $\times 2$
		$(\log \#steps)$ bits $\times 4$
<code>byte_range_checks</code>	Stores only: Element involved in the 8-bit-range checks of <code>memory_v_bytes[8]</code> performed with Lasso	$(\log \#steps)$ bits

Table 2: The extra elements involved per-byte in loads and stores.

A.3 Summary of CPU Step Constraints

Here, we go through the CPU steps outlined in Figure 1 and add more context in terms of the committed elements used and constraints involved.

1. Read the program code using PC to get the instructions details: `opflags[14]`, `opcode`, `rs1`, `rs2`, `imm`.
 - As discussed in Section 3.3, this involves six reads from program memory: one for each element of the tuple.
2. Read the source registers `rs1`, `rs2`. And then set operands x and y
 - Reading the source registers involves two memory-checking updates. The locations are the registers themselves.
 - The values and timestamps involved are committed advice elements: `(read_val_rs1, read_ts_rs1)` and `(read_val_rs2, read_ts_rs2)`, respectively.
 - The setting of operands x and y involves using `opflags[0]`, `opflags[1]` described earlier.
3. Perform loads and stores using memory-checking.
 - For both loads and stores, the memory location involved is `rs1 + imm`. This sum is calculated using constraints and also involves `opflag[14]` which holds the sign of `imm`.
 - This involves up to two range-checks and a less-than comparison per byte of memory read/written (see Appendix A.2).
 - As RISC-V memory is byte-addressable, the lookup argument involves $W/8$ memory operations. The required bytes-decomposition are the elements in `memory_v_bytes[W/8]`.
4. Construct the lookup element.
 - This is structured as `opcode || z` where z could be $x || y$, $x + y$, $x - y$ or $x * y$.
 - The exact format chosen is guided by many `opflags`.
5. Updating the destination register.

- If the corresponding `opflag` is 1, the lookup result is stored in `rd`.
- The memory-checking write operation here involves the location `rd`, value `result` and timestamp being the current global step count.

6. Updating the PC.

- For Jump instructions, $PC \leftarrow (\text{lookup } \text{result})$. For Branch instructions, $PC \leftarrow PC + \text{imm}$ (sum computed similar to Step 3 above) if and only if the lookup `result` was true. For other instructions, $PC \leftarrow PC + 4$. The right choice is guided by the corresponding `opcodes`.

B Overview of Memory-Checking Arguments

B.1 Merkle trees

One way to implement a memory-checking argument is via Merkle trees. This approach was first introduced to the literature on SNARKs in Pantry [BFR⁺13] and is still used today in so-called *incrementally verifiable computing* (IVC) schemes (see [BSCTV14] for an early example). Conceptually, the VM whose execution is being proved is modified as follows, so as to authenticate its own memory. The VM at all times tracks the root hash of a Merkle tree that has the contents of each of its memory cells at the leaves (let us say that the number of memory cells is N). Every read operation returns not only the value stored in the appropriate memory cell, but also a Merkle authentication path that the machine checks for consistency with the stored root hash. Write operations are preceded by a read operation, which allows the VM to update the root hash in an authenticated manner.

The downside of the Merkle tree approach is that each read and write requires proving knowledge of a Merkle authentication path, which involves $\log N$ many hash evaluations. This can be an effective approach if there are not many reads and writes, but is very expensive for the prover for executions involving many reads and writes.

B.2 Re-ordering the execution trace

The rough idea of this approach to designing a memory-checking argument is to have the prover “re-order” the execution trace so that, rather than entries appearing in increasing order of timestep, they instead are grouped by the memory location read or written at that step, and within each group, entries are sorted by time. This is called a “memory-ordered” execution trace (see [Tha22, Section] for an exposition). It is straightforward to design a quasilinear-size circuit that takes as input a memory-ordered execution trace, and confirms that the value returned by every read operation is indeed the value last written to the appropriate memory cell. zkVM projects taking this, or related, approaches to memory-checking include RISC Zero [BGtR23] and Cairo [GPR21].

The core of the memory-checking argument amounts to confirming that the memory-ordered execution trace is indeed a re-ordering (i.e., permutation) of the time-ordered execution trace. Some early works [BSCGT13] proposed to accomplish this using so-called *routing networks*, a complicated and expensive technique from the PCP literature. Later works [BCG⁺18, ZGK⁺18] proposed to instead use lightweight *permutation-invariant fingerprinting* techniques dating to the memory-checking literature from the late 1980s and early 1990s [Lip89, BEG⁺91]. The key idea in these techniques is to check whether two vectors $a, b \in \mathbb{F}^N$ are permutations of each other (i.e., whether there exists a permutation $\pi: \{1, \dots, N\} \rightarrow \{1, \dots, N\}$ such that $a_i = b_{\pi(i)}$ for all i) by having the verifier pick a random $r \in \mathbb{F}$, and checking that

$$\prod_{i=1}^N (r - a_i) = \prod_{i=1}^N (r - b_i). \quad (11)$$

If a and b are permutations, this check will pass with probability 1, while if they are not permutations, the check will pass with probability only $m/|\mathbb{F}|$. This technique is now pervasive in SNARK design, appearing in works such as Plonk [GWC19]. We refer to this as *permutation-invariant fingerprinting*.

1. Input: vectors $a, b \in \mathbb{F}^N$.
2. Goal: determine whether a and b are permutations of each other.
3. Pick a random $r \in \mathbb{F}^N$, and checks that

$$\prod_{i=1}^N (r - a_i) = \prod_{i=1}^N (r - b_i).$$

Figure 5: Permutation-invariant fingerprinting.

1. Input: vectors $a, b \in \mathbb{F}^N$.
2. Goal: determine whether $a = b$.
3. Pick a random $r \in \mathbb{F}^N$, and checks that

$$\sum_{i=1}^N a_i r^{i-1} = \sum_{i=1}^N b_i r^{i-1}.$$

Figure 6: Reed-Solomon fingerprinting. If $a = b$ then the check passes with probability 1. If $a \neq b$, then the check passes with probability at most $(N - 1)/|\mathbb{F}|$.

We have described the checking procedure as interactive, since the verifier picks the random field element r , but the procedure can be rendered non-interactive via the Fiat-Shamir transformation. In any SNARK that uses permutation-checking, the vectors a and b will be committed by the prover, and r will be chosen by hashing those commitments (and any other messages sent by the prover earlier in the interactive protocol).

B.3 Memory-checking via permutation-checking, without re-ordering

Spartan [Set20], Spice [SAGL18], and descendants including our companion paper Lasso [STW23] build on the memory-checking work of Blum et al. [BEG⁺91] to obtain memory-checking arguments that do not reorder the execution trace.

Overview of the Lasso lookup argument. To illustrate the approach, we begin with a brief overview of (the simplest variant of) the Lasso lookup argument from our companion paper [STW23], which has prover commitment costs linear in the table size. Here, let us consider a VM that is given read-only access to lookup table T . To render the VM's reads more easily checkable, let us modify the VM's reading procedure as follows.

- For each memory cell $j = 1, \dots, N$ maintain a counter c_j , which is supposed to track how many times cell j has been read.
- Every time a read operation to cell b_i returns a (value, count) pair (a, c) , have the VM write the tuple $(a, c + 1)$ to cell b_i . That is, the VM follows every read operation by writing the returned value back to the cell that was just read, incrementing the returned counter value by 1.
- When all the reads are done, the VM makes one final pass over memory. This final set of N reads are not paired with write operations.

Let RS (short for *read-set*) be the vector of all (cell, value, count) tuples returned by the memory across all read operations. Let WS (short for *write-set*) be the vector of all (cell, value, count) tuples across all write operations. WS includes the N writes to initialize memory, i.e., the tuples $(j, T[j], 0)$: $j = 1, \dots, N$. Prior

works [BEG⁺91, Set20, STW23] establish the following lemma, whose proof we omit for brevity.

Lemma 2. *RS and WS are permutations of each other if and only if the result of every read operation to each cell j indeed returns the (value, count) pair last written to that cell. Here, RS and WS are permutations of each other if they are equal as multisets of (cell, value, count) tuples.*

With this lemma in hand, one can obtain an (indexed) lookup argument as follows. Recall that in an indexed lookup argument, the prover has committed to two vectors $a, b \in \mathbb{F}^m$, and wishes to prove that $a_i = T[b_i]$ for all i . View a as the vector of all values returned by the read operations of the modified VM above, and b as the vector specifying the cell targeted by each read operation. The prover next commits to the vector $c \in \mathbb{F}^{m+N}$, whose i 'th entry is purported to be the count returned by the i 'th read operation.

If the prover chooses c honestly, then by Lemma 2, RS and WS are permutations of each other. Conversely, by the same lemma, if RS and WS are permutations of each other, then each read operation returned the correct value, and hence $a_i = T[b_i]$ for all $i = 1, \dots, m$. So it suffices for the lookup argument prover to prove that RS and WS are permutations of other.

To this end, first Reed-Solomon fingerprint each (cell, value, count) tuple. That is, the verifier picks a random $\gamma \in \mathbb{F}$ and sends γ to the prover. The prover then replaces all tuples (b, a, c) in RS or WS with $b + \gamma a + \gamma^2 c$.

This reduces RS and WS from vectors of $N + m$ triples of field elements, to vectors RS' and WS' in \mathbb{F}^{N+m} . If there are no “collisions” (two distinct tuples with matching fingerprints) then RS' and WS' are permutations of each other if and only if RS and WS are. Clearly, the probability of a collision is at most $(N + m)^2 / |\mathbb{F}|$ (a more careful analysis can bound the soundness error by at most $2(N + m) / |\mathbb{F}|$).

Hence, up to the above soundness error, checking whether RS and WS are permutations of each other is equivalent to checking whether RS' and WS' are permutations.

This latter check is done via permutation-invariant fingerprint (Figure 5).

In summary, in the indexed lookup argument, the prover first commits to the vector $c \in \mathbb{F}^{m+N}$ purported to be the counts returned by the read operations specified by the vectors $a, b \in \mathbb{F}^m$. (More precisely, the prover commits to the multilinear extension \tilde{c} of c , using any multilinear polynomial commitment scheme).

The verifier then confirms that RS is a permutation of WS by picking two random field elements γ, α , and ensuring that the following two quantities are equal:

$$\prod_{i=0}^{m-1} (\alpha - (b_i + \gamma \cdot a_i + \gamma^2 \cdot c_i)) \prod_{i=0}^{N-1} (\alpha - (i + \gamma \cdot T[i] + \gamma^2 \cdot c_{m+i})) \quad (12)$$

and

$$\left(\prod_{i=0}^{N-1} (\alpha - (i + \gamma \cdot T[i])) \right) \cdot \left(\prod_{i=0}^{m-1} (\alpha - (b_i + \gamma \cdot a_i + \gamma^2 \cdot (c_i + 1))) \right). \quad (13)$$

Here, Expressions (12) and (13) are the permutation-invariant fingerprints of RS' and WS' respectively.

Lasso forces the prover to correctly compute these two expressions using any grand product argument. Specifically, it Lasso suggests to use either a highly optimized variant of the GKR-protocol [GKR08] due to Thaler [Tha13], which avoids any additional commitment costs for the prover, or a variant with shorter proofs but slightly higher commitment costs [SL20]. At the end of these grand product arguments, the verifier needs to evaluate each of \tilde{T} , \tilde{a} , \tilde{b} , and \tilde{c} at a randomly chosen point. The evaluations of \tilde{a} , \tilde{b} , and \tilde{c} can be obtained from the polynomial commitment scheme used to commit to each of these polynomials. For MLE-structured tables, the evaluation of \tilde{T} can be computed by the verifier with only $O(\log N)$ field operations.

The above argument protocol is implicit in Spark, the sparse polynomial commitment scheme given in Spartan [Set20]. However, Spark's security analysis assumed that (the commitment to the) vector c of purported counts is computed by an honest party. This sufficed for Spartan's application, but not for giving a lookup argument. Our companion paper Lasso shows that the lookup argument is secure even if c is committed by a malicious party.

How Lasso handles large tables. If the lookup table is too large to justifying paying a commitment cost linear in the table size, but the table is decomposable (Definition 2.6), **Lasso** will automatically decompose any lookup into the large table into c lookups into smaller tables, and collate the results into the result of the lookup into the large table. The simple variant of **Lasso** described above (with prover costs linear in table size) is applied to each subtable. See our companion paper [STW23] for details.

Overview of Spice’s memory-checking argument. In (the variant of) the **Lasso** lookup argument described above, every read operation is followed by a write that increments the count returned by the read by one. When supporting read/write memory as required in a memory-checking argument, incrementing counts by one is not sufficient to ensure security. This is because in the memory-checking setting, the machine that is reading from, and writing to, memory may write a value to memory that differs from the value returned by the most recent read operation. This gives an attacker more flexibility than a **Lasso** attacker has. In particular, if one tries to apply **Lasso** in the memory-checking setting, an attacker can potentially answer reads “out of order”, meaning answering a read at time i with a value that won’t be written until time $j > i$.

To prevent this type of attack, Spice updates counts in a different way (which is, unfortunately, more expensive to implement in the context of memory-checking arguments). Specifically, in Spice’s memory-checking argument, the machine is modified to maintain a timestamp ts . As in **Lasso**, every write is preceded by a read. If the read operation returns a count c , the count that is written is not $c + 1$ as in **Lasso**, but rather $\max\{c, \text{ts}\} + 1$. The machine then updates the timestamp to $\max\{c, \text{ts}\} + 1$.

In Spice itself, $\max\{c, \text{ts}\}$ is computed by having the prover provide the bit decomposition of c and ts as untrusted advice. In **Jolt**, we instead compute $\max\{c, \text{ts}\}$ with a single lookup, into the evaluation table of the max function. Note that if the prover is honest, c and ts are always between 0 and the number of steps m that the machine is running for. This ensures that we can use a lookup table of size roughly m^2 .

C A brief overview of two’s complement representation

An *unsigned L -bit data type* refers to a value $z \in \{0, 1, \dots, 2^L - 1\}$. A *signed L -bit data type* (in two’s-complement format) refers to a value $z \in \{-2^{L-1}, \dots, 2^{L-1} - 1\}$. The two’s-complement representation $[z_{L-1}, \dots, z_0] \in \{0, 1\}^L$ of z is the unique vector such that

$$z = -z_{L-1} \cdot 2^{L-1} + \sum_{i=0}^{L-2} 2^i z_i. \quad (14)$$

For clarity, when discussing instructions interpreting their inputs as signed data types represented in two’s-complement format (e.g., Section 5.3), we refer to z_{L-1} as the sign bit of z , and denote this by z_s . We use $z_{<s}$ to refer to $[z_{L-2}, \dots, z_0] \in \{0, 1\}^{L-1}$.

As discussed in Section 3, the use of two’s complement allows instructions to operate identically regardless of whether or not the inputs are interpreted as signed or unsigned. For example, consider the **ADD** instruction when $L = 3$.

When adding three-bit unsigned integers 3 and 4, the addition operation proceeds as follows:

$$3 \text{ (i.e., } 011) + 4 \text{ (i.e., } 100) = 7 \text{ (i.e., } 111).$$

Here, in parenthesis we have provided the binary representations of 3, 4, and 7 when interpreted as unsigned data types in two’s-complement format.

When adding three-bit signed integers 3 and -4 , the addition operation proceeds as follows:

$$3 \text{ (i.e., } 011) + -4 \text{ (i.e., } 100) = -1 \text{ (i.e., } 111).$$

Again, in parentheses we have provided the binary representations of 3 and -4 when interpreted as signed data types in two's complement format.

The above example demonstrates that, when using two's complement binary representations, the input/output behavior of the addition operation is independent of whether the inputs are interpreted as signed or unsigned.