

Mechanised operational semantics of Rowhammer

MARTIN BERGER, University of Sussex, UK and Montanarius Ltd, UK
AMIR NASEREDINI, Huawei R&D UK Ltd, UK

Rowhammer is a hardware vulnerability in dynamic random-access memory (DRAM) in which repeated accesses to one or more aggressor rows can induce bit-flips in nearby victim rows. This phenomenon violates a core assumption of conventional programming language semantics: that reading from or writing to one memory location does not modify others. Despite the security importance of this phenomenon, there is no widely established formal framework connecting Rowhammer faults with program behaviour. This makes it difficult to reason rigorously about the efficacy of proposed Rowhammer defences and the program-level guarantees they provide. To address this gap, we present a probabilistic small-step operational semantics for an idealised imperative language subject to Rowhammer-style faults. The semantics is high-level in that it abstracts from DRAM internals and semiconductor physics. A general probabilistic fault model parameterises the semantics, representing Rowhammer-style faults by assigning probabilities to bit-flips during read or write operations within a specified victim region. The resulting distributions are propagated through programs using the standard monadic structure of probabilistic computation. As a case study, we formalise physical separation, a well-known defence that places program variables sufficiently far apart in physical memory that an access to one variable cannot disturb another. We prove a distribution-independent semantic collapse theorem: for every finite execution, including prefixes of terminating and non-terminating executions, the protected projection of the probabilistic Rowhammer semantics is the Dirac distribution of the corresponding Rowhammer-free execution. Furthermore, we develop an observation-parametric account of secure information flow. Non-interference is expressed as a hyperproperty comparing the distributions of low observations from low-equivalent initial memories. Under physical separation, ordinary non-interference and probability-sensitive Rowhammer non-interference coincide for every observer of protected behaviour. Consequently, physical separation preserves non-interference for every admissible fault model, while every Rowhammer non-interference violation reflects a violation already present in the Rowhammer-free semantics. The development is fully mechanised in Lean using mathlib, relying on no unfinished proofs or problem-specific axioms.

CCS Concepts: • **Theory of computation** → **Program semantics**; • **Software and its engineering** → **Formal methods**; • **Security and privacy** → *Information flow control*; *Systems security*; **Logic and verification**.

Additional Key Words and Phrases: Probability theory, Probabilistic fault models, Rowhammer, Programming language semantics, Operational semantics, Relational safety, Mechanised verification, Lean 4, Secure information flow, Non-interference.

1 Introduction

Most reasoning about software relies on a locality assumption about memory access effects:

Reading or writing one memory location does not modify other locations.

Rowhammer violates this assumption. Repeated activation of one or more aggressor rows in dynamic random-access memory (DRAM) can induce bit-flips in physically nearby victim rows [16, 35]. The resulting mismatch between the memory model assumed by software semantics and the behaviour of its semiconductor implementation creates a semantic gap: hardware research describes disturbance mechanisms, address mappings, and mitigations, whereas software verification reasons about commands, stores, traces, observations, and security properties. This gap is particularly important for security. Rowhammer is usually presented as an integrity failure, but corruption of a high-security value can alter control flow or a later low-security output, while corruption of a low-security location can directly change an attacker-visible result. Existing theories of secure

Authors' Contact Information: [Martin Berger](mailto:Martin.Berger@martinfriedrichberger.net), contact@martinfriedrichberger.net, University of Sussex, Brighton, UK and Montanarius Ltd, London, UK; Amir Naseredini, sahnaseredini@gmail.com, Huawei R&D UK Ltd, Cambridge, UK.

information flow and non-interference [10, 34] therefore cannot simply be applied unchanged: they require a semantics that makes the probabilistic, non-local effect of memory accesses explicit. This paper asks closely connected research questions.

- **RQ1.** Can Rowhammer be modelled compositionally in the familiar language of operational semantics, without modelling device physics?
- **RQ2.** Can the Rowhammer fault mechanism be captured by a small, conceptually non-intrusive semantic interface?
- **RQ3.** Can such an interface be made substantially independent of the target programming language?
- **RQ4.** Can such a model express and justify a concrete Rowhammer defence?
- **RQ5.** Can the abstraction make Rowhammer amenable to existing semantic security frameworks, such as information flow and non-interference reasoning?

We answer all questions affirmatively.

1.1 The semantic view: abstracting memory access

A standard way to give semantics to an imperative language is to treat programs as state transformers. In its most concrete form, assignment and variable lookup are explained by clauses such as

$$[[x := e]](m) = m[x \mapsto v] \quad \text{where} \quad [[e]](m) = v,$$

and

$$[[x]](m) = m(x).$$

where m ranges over memories and e over expressions. These equations present memory as a concrete map and use two primitive operations on that map: read the value stored at a location, and update the value stored at a location. A recurring lesson of programming language semantics is that these two operations should be made explicit. Once lookup and update are abstracted from the concrete function space $\text{Memory} = \text{Loc} \rightarrow \text{Val}$, where Loc are locations in memory, and Val the values being stored, the syntax of the language no longer has to know how memory is represented. The same assignment rule can then be interpreted over different models of store, aliasing, heaps, concurrency, failure, or other effects.

This idea has appeared repeatedly, under different names. McCarthy’s select/store axioms for arrays isolate reading and updating as primitive operations on memories [23]. Monadic and algebraic-effect semantics make the separation systematic: computations with state are interpreted as effectful maps, and state itself may be presented by operations such as lookup and update together with their laws [24, 32]. Modular operational semantics pursues a related goal from the small-step side, factoring stores and other auxiliary components out of individual transition rules so that language constructs can be specified more independently [25]. The common pattern is that memory access becomes an interface rather than an implementation detail.

Our Rowhammer semantics uses this idea. Ordinary memory operations are deterministic:

$$\text{read} : (\text{Loc} \times \text{Memory}) \rightarrow \text{Val} \quad \text{write} : (\text{Loc} \times \text{Val} \times \text{Memory}) \rightarrow \text{Memory}$$

Rowhammer-affected read and write are instead probabilistic operations¹:

$$\text{read} : (\text{Loc} \times \text{Memory}) \rightarrow \text{Dist}(\text{Val} \times \text{Memory}) \quad \text{write} : (\text{Loc} \times \text{Val} \times \text{Memory}) \rightarrow \text{Dist}(\text{Memory})$$

The operations still have the shape of a local memory access, but their results are distributions that include the post-access memory. The distributions model Rowhammer faults.

¹We discuss in Section 5 why the codomain of read needs to be $\text{Dist}(\text{Val} \times \text{Memory})$ and not $\text{Val} \times \text{Dist}(\text{Memory})$.

99 The rest of the semantics is then obtained monadically: deterministic steps are embedded as
100 Dirac distributions, and compound computations are sequenced by monadic bind for the proba-
101 bility monad. Thus the Rowhammer model is not a separate hardware semantics bolted onto the
102 language. It is a small probabilistic reinterpretation of the read/write interface already present in
103 ordinary state-transformer semantics. This is the key abstraction boundary used throughout the
104 paper: Rowhammer enters only at memory accesses, the surrounding operational semantics lifts
105 compositionally around that local probabilistic effect.

106 This move from deterministic to probabilistic state transformers is not unique to Rowhammer:
107 deterministic Turing machines and programming languages naturally become probabilistic by
108 allowing programs to branch according to a probability distribution. In those settings, probabilistic
109 choice is an intentional feature of the syntax or model. Rowhammer is different. Its probabilistic
110 behaviour is an unwanted, observable consequence of the physical substrate, making it conceptually
111 closer to Shannon’s noisy-channel model from information theory: an intended memory operation
112 is transmitted through an unreliable medium and may be received as corrupted state. In contrast,
113 this noise is not merely ambient: it is actively shaped by an attacker’s chosen access patterns and
114 strictly constrained by physical locality in DRAM. Our semantic task is therefore to expose this
115 hardware-level unreliability without modifying the language syntax, placing a probabilistic fault
116 model directly behind the ordinary read/write interface. Ultimately, the novelty lies not in using
117 probabilistic state transformers per se, but in the observation that Rowhammer can be elegantly
118 isolated as a probabilistic reinterpretation of memory access.

120 1.2 Contributions

121 To address the research questions, we make the following contributions:

- 122 • A fine-grained labelled operational semantics in which Rowhammer has a small, explicit
123 interface at reads and writes.
- 124 • A general probabilistic fault model separating spatial confinement from numerical assump-
125 tions about bit-flips.
- 126 • A distribution-independent soundness theorem for physical separation, a well-known
127 Rowhammer defence, covering terminating and non-terminating executions.
- 128 • An observation-parametric theorem that physical separation preserves and reflects probability-
129 sensitive non-interference.

130
131 The complete development is mechanised in Lean [6] using mathlib [37]. The development con-
132 tains no unfinished proofs or problem-specific axioms. An axiom audit (`scripts/axiomAudit.sh`
133 in the repo) confirms that all results rest only on Lean’s three foundational axioms (`propext`,
134 `Classical.choice`, `Quot.sound`): no probability axioms are added. The Lean excerpts shown are
135 automatically extracted verbatim from the sources so they cannot drift from the development. All
136 Lean code is open-sourced and anonymised at [1]. Full operational rules, auxiliary metatheory, and
137 extended mechanisation details are provided in the supplementary material. The appendices cited
138 in the text are provided as anonymised supplementary material.

140 2 Introduction to Rowhammer

141 To appreciate what the present work does and does not promise, it helps to understand the physical
142 mechanism it abstracts. The rest of the paper relies on only a few concepts from DRAM architecture,
143 rows, activation, refresh, and the locality of disturbance, and we introduce them here from first
144 principles. Readers already familiar with Rowhammer may skip to the examples at the end of this
145 section. For the underlying semiconductor physics we refer to [39]. For a broad survey of the attack
146 landscape see [26].

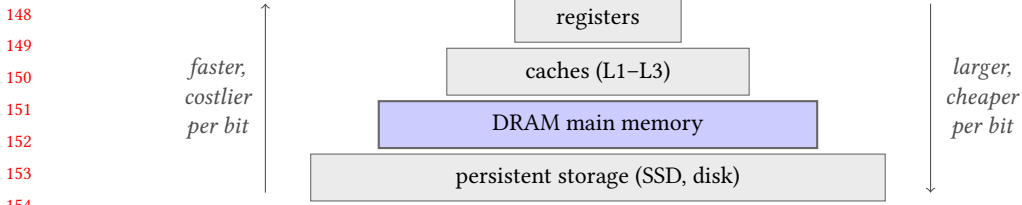


Fig. 1. The memory hierarchy. Each level trades speed for capacity. DRAM (highlighted) wins its place through cell minimalism, one transistor and one capacitor per bit, and the resulting drive towards ever-denser packing is the root cause of the disturbance effects this paper models.

The memory hierarchy. DRAM occupies an important point in a trade-off (Figure 1). Above it sit registers and caches: fast, but expensive. Below sits persistent storage: vast, but orders of magnitude slower. DRAM’s role is to be as large and as cheap as possible while remaining just fast enough. It achieves this with only one transistor and one capacitor per bit, against the six transistors of a cache cell. This economic pressure towards density is Rowhammer’s root cause. The more tightly cells are packed, the stronger the electrical coupling between them, which is why the problem has worsened with every DRAM generation [16].

How DRAM stores data. At the hardware level, each DRAM bit is stored as electrical charge in a capacitor, guarded by an access transistor. Cells are organised into large two-dimensional grids called *banks*, which operate independently of one another. Within a bank, a horizontal line of cells, a *row*, typically a few kilobytes wide, shares a common control wire known as the *wordline*, while cells in the same column share a *bitline*. Individual cells cannot be addressed directly. To access any cell, the memory controller must *activate* the row containing it: raising the row’s wordline voltage connects every cell of the row to its bitline and copies the row’s contents into the bank’s *row buffer*, where reads and writes then take place. Accessing a different row of the same bank requires closing the currently open row, meaning lowering its wordline voltage before activating the new one. Capacitors leak charge. Left alone, a DRAM cell loses its stored charge within a fraction of a second, so the memory controller periodically rewrites every row—a *refresh*. In commodity DRAM, every row is typically guaranteed a refresh once every 64 ms [16]. The reliability of DRAM therefore rests on a quantitative margin: between two refreshes, a cell’s charge must stay on the correct side of a threshold, despite whatever electrical disturbance its neighbourhood produces.

The Rowhammer effect. Kim et al. [16] showed in 2014 that this margin can be exhausted *deliberately*. Every activation of a row causes slight voltage fluctuations on its wordline, and these fluctuations accelerate charge leakage in the cells of *physically adjacent* rows, an electromagnetic coupling effect that grows stronger as cells are packed more densely. A single activation is harmless. But a program that activates the same row hundreds of thousands of times within one refresh interval, *hammering* it, can drain a neighbouring cell past its threshold before the next refresh arrives. The value read thereafter is not the value last written: a bit has flipped in a row that was never written to, or even addressed, by the program. The repeatedly activated rows are called *aggressor* rows and the corrupted ones *victim* rows (Figure 2). The minimum number of activations needed is called the *Rowhammer threshold*: it has dropped sharply as DRAM density has increased, from over a hundred thousand activations on DDR3 to tens of thousands on modern chips [7, 16], which is why the problem is getting worse, not better, with each DRAM generation.

Rowhammer is not an exotic laboratory effect. It is exploitable from ordinary unprivileged code, including JavaScript in a web browser [11], and flipping a single page-table or capability bit can

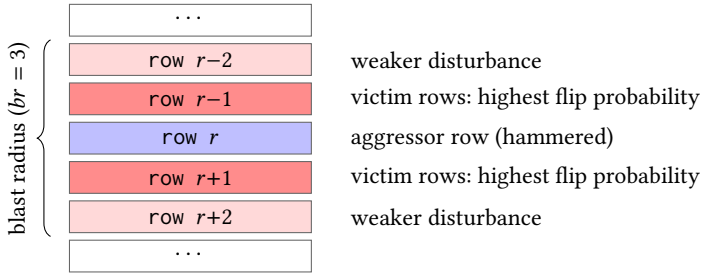


Fig. 2. Disturbance is local. Repeatedly activating (*hammering*) row r accelerates charge leakage in nearby rows. The flip probability is highest for the immediate neighbours $r \pm 1$ and decays with distance. Beyond the blast radius br the effect is negligible. The victim rows are corrupted although the program never writes to or even addresses them.

escalate one hardware fault into arbitrary memory access [35]. The effect has also grown richer since its discovery. Half-Double [17] induces flips beyond immediate neighbours by combining accesses at different distances, TRRespass [7] and Blacksmith [15] bypass in-DRAM mitigations with many-sided and non-uniform hammering patterns, and RAMBleed [20] turns the data-dependence of flip probabilities into a *read* primitive, making Rowhammer a confidentiality problem as well as an integrity problem.

Locality and the blast radius. The one saving grace of DRAM physics is its strong spatial locality: the disturbance an activation causes decays rapidly with physical distance, and beyond a few rows it is negligible. In the rest of the text we call the maximum distance at which an access can disturb a cell the *blast radius*, abbreviated br . The blast radius of a given DRAM device can be approximated empirically, and it is the key intuition behind an entire class of defences: if everything worth protecting is kept at least br rows away from everything an attacker can hammer, the disturbance never reaches it.

Physical separation defences. Defences based on *placement* exploit exactly this locality. CATT [2] partitions kernel from user memory with unused *guard rows*, so that user-space hammering cannot reach kernel rows. ZebRAM [18] interleaves guard rows between all data rows (Figure 3). The correctness claim of such defences is physical: every row an attacker can reach is separated from every row worth protecting by more than the disturbance radius. In this paper the discipline is captured by a predicate $\text{safe}(br, locs)$: the victim region generated by an access to any protected location is disjoint from the protected set $locs$. In the distance-based instance this means that distinct protected rows lie at least br apart. But separation is required *within* the protected set too, because the program’s own accesses hammer their own neighbourhoods. The guard-row layouts above are the same discipline applied to the whole address space, with the attacker-reachable region kept at distance br from everything protected.

The semantic gap. Program verification speaks a different language: commands, stores, traces, observables. A proof about a program is a proof about its semantics, and standard semantics assume that memory changes only when written. Rowhammer breaks exactly that assumption, and placement defences restore it, but no theorem connects the two levels. What is missing is a semantics in which the physical claim (“flips are confined to a blast radius”) and the placement claim (“protected rows are outside every blast radius”) can both be stated, and a theorem that together they restore the ordinary meaning of programs. This paper supplies both. The examples below,

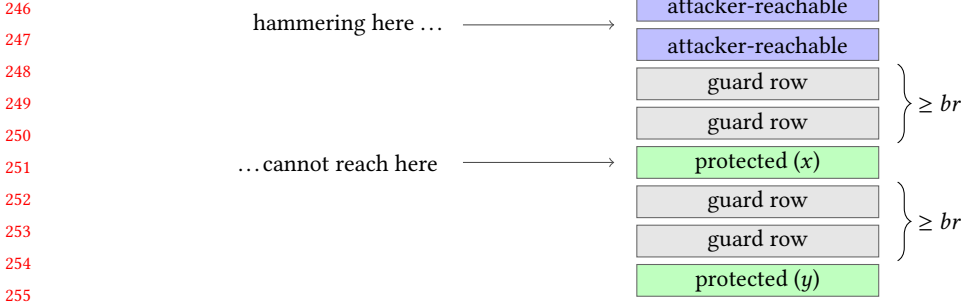


Fig. 3. Physical separation with a blast radius of $br = 3$. Guard rows keep every protected row at least br rows away from every row the program (or an attacker) accesses, including the *other* protected rows, since accesses to x hammer x 's own neighbourhood. This is the layout discipline that the predicate $\text{safe}(br, locs)$ of Section 5.1 abstracts.

adapted from [27], illustrate what goes wrong without separation, and, read against Theorem 6.1, identify exactly which hypothesis each failure violates.

Example 2.1 (Corruption without a write). Let x and y occupy adjacent rows, and let the program hammer y while never assigning to x :

$$x := 1; \text{ while } (0 = 0) \text{ do } y := y + 1.$$

Every access to y disturbs the blast radius of y 's row, which contains x . In every ordinary semantics x holds 1 forever. Under Rowhammer, the probability that x still holds 1 decays with every iteration. No analysis that equates “unwritten” with “unchanged” can see this.

Example 2.2 (A dead branch awakens). Corruption is not limited to data, it can also redirect control. Let x and y be adjacent and consider

$$x := y; x := x; \text{ if } x = y \text{ then skip else } P.$$

Ordinarily the guard is always true and P is dead code. But the reads and writes of x hammer x 's neighbourhood, which contains y : a flip in y between the first assignment and the guard makes the guard false, and the machine executes P . Which code runs is now probabilistic. This is why the protected view of Section 6.1 tracks the *residual program* and the *access trace*, not just memory: a defence that preserved memory but not control flow would not restore the meaning of programs. Note also that both x and y here are the program's *own* variables: if the protected set $locs = \{x, y\}$ places them in adjacent rows, $\text{safe}(br, locs)$ fails, and the theorem correctly does not apply: separation is required *within* the protected set, not only between the program and the outside world.

Example 2.3 (An unowned access). Separation of the program's own rows is not enough either. Suppose $locs = \{x\}$ is safely laid out, but the program also touches a row $z \notin locs$ that happens to be adjacent to x :

$$z := z + 1; \text{ if } x = 0 \text{ then skip else } P.$$

The access to z triggers a fault in z 's blast radius, about which $\text{safe}(br, locs)$ says nothing, and that radius contains x . This is what the well-formedness hypothesis $wf_{locs}(P)$ of Theorem 6.1 rules out: every location the program touches must be part of the protected layout, so that every fault the program can trigger is one the separation condition constrains.

These examples give us the shape of the theorem we need: quantified over fault behaviour (flip probabilities are the attacker’s, not ours), conclusion covering memory, control, and accesses alike, and hypotheses (frame-locality, separation, well-formedness) each of which is violated by some concrete program above when dropped.

3 Mathematical preliminaries

We write $\mathcal{P}(S)$ for the *powerset* of a set S and $\mathcal{P}_{fin}(S)$ for the set of all finite subsets of S . The Rowhammer semantics assigns to each transition a *distribution* over successor configurations. This section fixes the small amount of probability theory involved. Readers familiar with the probability monad (also known as Giry monad) [8, 21] or with monadic treatments of probabilistic programs [5, 19, 33] can skim it for notation. All probability spaces in this paper will be discrete, *i.e.*, finite or countably infinite, so all occurring measure spaces are trivial and we shall not mention them. For a type α , we write $\text{Dist}(\alpha)$ for the type of discrete probability distributions on α , *i.e.*, functions $\mu : \alpha \rightarrow [0, 1]$ with $\sum_x \mu(x) = 1$. The *support* of μ , written $\text{supp}(\mu)$, is the set of x with $\mu(x) > 0$. The *Dirac distribution*, also called *point mass*, at $x \in \alpha$ is the distribution $\delta_x \in \text{Dist}(\alpha)$ given by $\delta_x(y) = 1$ if $y = x$, and $\delta_x(y) = 0$ otherwise, for all $y \in \alpha$. A *probability kernel* from α to β is an input-indexed family of distributions, *i.e.*, a function $\kappa : \alpha \rightarrow \text{Dist}(\beta)$. For each input $a \in \alpha$, the value $\kappa(a)$ is therefore a distribution on β . Any distribution $\mu \in \text{Dist}(\beta)$ determines the constant kernel $a \mapsto \mu$. Likewise, a deterministic function $f : \alpha \rightarrow \beta$ embeds into the probabilistic setting as the Dirac kernel $a \mapsto \delta(f(a))$. In this paper, “kernel” always means a discrete probability kernel of the form $\alpha \rightarrow \text{Dist}(\beta)$. In the Lean development this is represented as a function returning a PMF (probability mass function), rather than by mathlib’s more general measure-theoretic type `ProbabilityTheory.Kernel`. Sequential composition of probabilistic computations is monadic. For our purposes, a monad consists of a type constructor M together with operations *unit* $\text{return}_\alpha : \alpha \rightarrow M(\alpha)$ and *bind* $\gg = : M(\alpha) \rightarrow (\alpha \rightarrow M(\beta)) \rightarrow M(\beta)$, satisfying the left-identity, right-identity, and associativity laws [22, 24]. Intuitively, $M(\alpha)$ is a type of computations producing values of type α , *return* embeds a pure value, and *bind* sequences. The classical probability monad \mathcal{G} instantiates this structure with probability measures on measurable spaces. Since all spaces in this paper are discrete, we work with its discrete analogue, spelled out next. On objects it is given by $\mathcal{G}(\alpha) \triangleq \text{Dist}(\alpha)$. On morphisms, \mathcal{G} acts by pushforward of distributions: for $f : \alpha \rightarrow \beta$, we have $\mathcal{G}(f) : \text{Dist}(\alpha) \rightarrow \text{Dist}(\beta)$, and for $\mu \in \text{Dist}(\alpha)$, we set $\mathcal{G}(f)(\mu)(y) = \sum_{x \in \alpha, f(x)=y} \mu(x)$. The unit $\text{return}_\alpha : \alpha \rightarrow \mathcal{G}(\alpha)$ is the Dirac distribution δ_x , in other words, the point mass concentrated at x . \mathcal{G} ’s bind operation, given a distribution $\mu \in \text{Dist}(\alpha)$ and a probability kernel $\kappa : \alpha \rightarrow \text{Dist}(\beta)$, is written

$$\text{let } x \leftarrow \mu \text{ in } \kappa(x),$$

and defined pointwise by $(\text{let } x \leftarrow \mu \text{ in } \kappa(x))(y) = \sum_{x \in \alpha} \mu(x) \kappa(x)(y)$. In the Lean development, \mathcal{G} is represented by PMF. The Lean development mirrors this notation in the file `DiscreteDist.lean`. There, `Dist` is an abbreviation for mathlib’s PMF, the type of discrete probability mass functions. Thus the formalisation does not assume finite support. The file provides the small interface used in the proofs: Dirac distributions, bind, map, support lemmas, and injectivity of Dirac.

4 A deterministic small-step operational semantics for a WHILE language

The aim of this paper is to express Rowhammer in terms of programming language semantics, so as to be able to reason about program correctness and the efficacy of Rowhammer defences using the extensive tool-sets developed over the previous decades by the programming language and verification communities, but without having to go through semiconductor physics, or even knowledge of DRAM organisation. For this purpose we need a paradigmatic programming language that is simple

enough to allow us to explain the key ideas without being drowned out by language complexity, but at the same time general enough to cover interesting Rowhammer-related phenomena. We choose the well-known WHILE language [29, 38]. Its syntax is given by the following rules.

$$\begin{aligned}
 e &::= c \mid x \mid e + e \\
 b &::= \text{true} \mid \text{false} \mid \neg b \mid b \wedge b \mid e = e \\
 P &::= \text{skip} \mid x := e \mid P; P \mid \text{if } b \text{ then } P \text{ else } P \mid \text{while } b \text{ do } P
 \end{aligned}$$

Here c ranges over integers, and x over memory locations, collectively denoted Loc . It is natural to assume that $\text{Loc} = \mathbb{Z}$, but nothing in our development depends on that. Since the language lacks higher-order functions and dynamic memory allocation, we can associate with each e , b and P the set of locations syntactically mentioned. This gives rise to the `WellFormed` predicates in our Lean formalisation: e.g., `p.WellFormed locs` means that the locations occurring syntactically in the program p are contained in the set locs . Clearly well-formedness is preserved by reduction.

We give a conventional small-step operational semantics with two minor twists. The operational semantics is a deliberately *fine-grained* small-step semantics: arithmetic expressions, Boolean expressions, and programs all reduce one step at a time, so every memory read and every write is an individual transition. Transitions carry a *label* recording the access they perform:

$$\ell ::= \text{read}(x) \mid \text{write}(x) \mid \tau.$$

A read of variable x is labelled $\text{read}(x)$, the write performed by an assignment to x is labelled $\text{write}(x)$, and every other transition, operator reduction, branch selection, sequencing, loop unfolding, is silent (τ). We write Acc for the set of all such labels. The labels play no role in the deterministic semantics itself (only assignment writes modify memory) but help injecting the Rowhammer faults (Section 5). Congruence rules propagate the label of their premise. We use three labelled judgements on configurations

$$\langle e, m \rangle \xrightarrow{a} \langle e', m' \rangle \quad \langle b, m \rangle \xrightarrow{b} \langle b', m' \rangle \quad \langle P, m \rangle \xrightarrow{p} \langle P', m' \rangle$$

In all cases, m ranges over memories, *i.e.*, functions of type $\text{Loc} \rightarrow \mathbb{Z}$. We use `Memory` to refer to the set of all memories. Evaluation is left-to-right, and conjunction is deliberately non-short-circuiting. Values and `skip` have no outgoing transition. The relations are deterministic, and every non-terminal configuration can progress. The definitions are standard, and we give the most relevant rules next.

$$\begin{array}{c}
 \frac{}{\langle x, m \rangle \xrightarrow{\text{read}(x)} \langle m(x), m \rangle} \quad \frac{\langle e, m \rangle \xrightarrow{a} \langle e', m' \rangle}{\langle x := e, m \rangle \xrightarrow{p} \langle x := e', m' \rangle} \quad \frac{}{\langle x := v, m \rangle \xrightarrow{\text{write}(x)} \langle \text{skip}, m[x \mapsto v] \rangle} \\
 \frac{}{\langle \text{while } b \text{ do } P, m \rangle \xrightarrow{p} \langle \text{if } b \text{ then } (P; \text{while } b \text{ do } P) \text{ else skip}, m \rangle}
 \end{array}$$

The full rule sets and this basic metatheory are standard and given in Appendix A. We write $\xrightarrow[p]{s}$ for the reflexive and transitive closure of \xrightarrow{p} . The length $|s|$ counts every transition, expression steps included; the *access trace* $\text{acc}(s)$ is the subsequence of non- τ labels, whose length is the number of memory accesses the run performs. A program *terminates* with final memory m' , written $\langle P, m \rangle \Downarrow m'$, when $\langle P, m \rangle \xrightarrow[p]{s} \langle \text{skip}, m' \rangle$ for some s ; it *diverges*, written $\langle P, m \rangle \Uparrow$, when for every n some s with $|s| = n$ extends its run. The transitions give a partial function on configurations $\langle P, m \rangle$. It will later be convenient to extend this to total functions. We do this by defining a deterministic

393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441

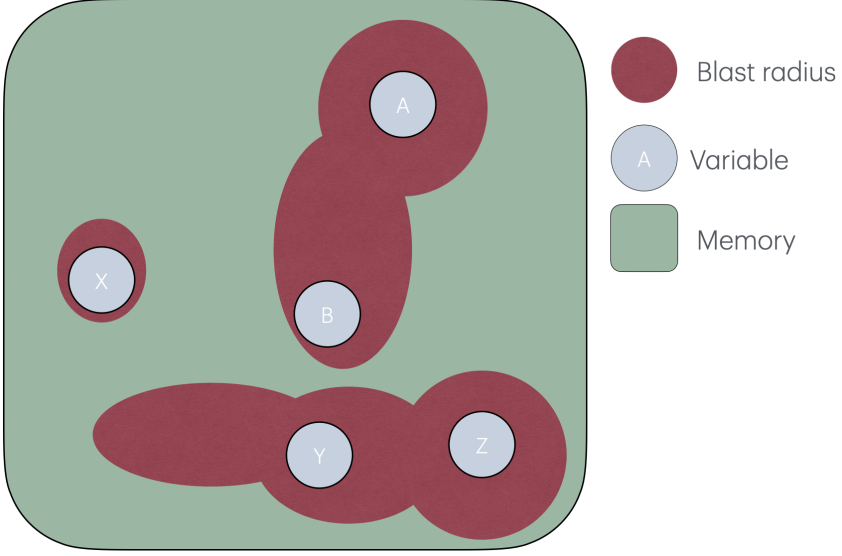


Fig. 4. The idea behind the fault model is that each variable x has a 'blast radius', a set of other memory locations that may get flipped when x is accessed. Locations outside the blast radius will not be flipped. The fault model formalises the shape of each blast radius, and the flip probabilities. Unlike Figures 2 and 3, this image does not show row shapes, because our formal model does not assume any specific geometry of blast radii.

absorbing run denoted $\text{run}_{\text{det, tr}}^n(C)$: iterate the step function n times, letting terminal configurations step to themselves with label τ , and accumulate the labels. Its result is the unique configuration-with-trace reached after n steps.

5 A probabilistic model of Rowhammer

This section addresses RQ1, RQ2 and RQ3: we now refine the ordinary semantics from the previous section with probabilistic Rowhammer fault models. We aim for the interface between the programming language semantics and the fault model to be minimal and canonical. The syntax, evaluation order, and reduction contexts are unchanged, and the refinement touches exactly one place: the $\text{read}(x)$ - and $\text{write}(x)$ -labelled transitions of Section 4 additionally sample from the Rowhammer fault model and produce a probability distribution (e.g., bit-flips in successor memories), while τ -labelled transitions remain deterministic.

With the changes to the read and write functions already explained in the Introduction, this amounts to changing the state transformer semantics of programs from type

$$\text{Memory} \rightarrow \text{Memory}$$

from Section 4 into probabilistic state-transformers

$$\text{Memory} \rightarrow \text{Dist}(\text{Memory})$$

(and likewise for expressions). The fault probabilities are lifted in a canonical way to general programs using the probability monad from Section 3.

442 *Why reads return a joint distribution.* There is a small but important typing choice in the Rowham-
 443 mer semantics. The ordinary read operation has type

$$444 \text{read} : (\text{Loc} \times \text{Memory}) \rightarrow \text{Val}$$

445
 446 Since a Rowhammer access to a location x is intended to disturb only other locations, one might
 447 expect the probabilistic read operation to have type

$$448 (\text{Loc} \times \text{Memory}) \rightarrow (\text{Val} \times \text{Dist}(\text{Memory}))$$

449
 450 returning the value stored at x together with a distribution over successor memories. In the model
 451 used in this paper, this would indeed be sufficient for a *single* read step: the value returned by
 452 the read is determined by the input memory, and the probabilistic fault only affects the successor
 453 memory. We instead use the slightly more general type

$$454 \text{read} : (\text{Loc} \times \text{Memory}) \rightarrow \text{Dist}(\text{Val} \times \text{Memory})$$

455
 456 This is the standard form of a probabilistic state transformer with a return value: it produces
 457 a joint distribution over the observed result and the successor state. The more specialised type
 458 $\text{Val} \times \text{Dist}(\text{Memory})$ embeds into this one by pairing every successor memory with the same
 459 returned value. Using the joint form $\text{Dist}(\text{Val} \times \text{Memory})$ avoids baking this special case into the
 460 semantic interface and ensures that read operations share the same monadic shape as the rest of
 461 the semantics. This leaves the abstraction general enough to model bit-flip errors, where a physical
 462 fault correlates the specific value returned by a noisy read with the resulting successor memory,
 463 without changing the semantic type of reads.
 464

465 5.1 The fault-kernel contract

466
 467 Fault kernels is our name for the probability kernels that model Rowhammer. The idea behind fault
 468 kernels is simple: they describe *what* memory cells may be affected by Rowhammer effects when a
 469 location x is accessed and *how*. Figure 4 visualises this intuition. We formalise this intuition with
 470 the following predicates.

- 471 • $\text{adj}(br, locs, x)$ gives the set of all memory locations within the blast radius br of location x .
- 472 • $\text{pflip}(S, m) \in \text{Dist}(\text{Memory})$ is a probability distribution that is allowed to flip bits in the
 473 finite set S of locations in memory m . Locations outside of S must remain unchanged.

474
 475 We can now state the modelling assumptions we will use in Section 6 for proving the soundness of
 476 physical separation as a Rowhammer defence.

- 477 • **(F1) Access-triggered.** Faults are sampled exactly at the $\text{read}(x)$ - and $\text{write}(x)$ -labelled
 478 transitions of Section 4, τ -transitions are fault-free, modelled as the Dirac distribution. (F1)
 479 is enforced by construction of the rules below, not assumed.
- 480 • **(F2) Spatial confinement (frame-locality).** Every memory in the support of $\text{pflip}(S, m)$
 481 agrees with m outside the victim set S . Within S the corruption may assign arbitrary
 482 probabilities to arbitrary, correlated corruptions: the kernel may behave adversarially,
 483 subject only to this clause.
- 484 • **(F3) History-freeness.** The kernel depends only on the current victim set and the current
 485 memory: pflip is a function $\mathcal{P}_{fin}(\text{Loc}) \rightarrow \text{Memory} \rightarrow \text{Dist}(\text{Memory})$, so a fault cannot
 486 depend on the access history, on time, or on hidden hardware state. This is a genuine mod-
 487 elling restriction. It excludes, for example, charge accumulation across repeated activations
 488 and refresh-cycle effects, generalising to stateful kernels is future work, as discussed in
 489 Section 9.

- **(F4) Stationarity.** The same kernel governs every access of a run, so fault behaviour does not drift over time. Like (F3), this is encoded by `pflip` being a single function rather than assumed as a separate axiom.

Clauses (F1-F4), considered in isolation, do not require the location whose access triggers a fault to lie outside its own victim set: the abstract contract permits $x \in \text{adj}(br, locs, x)$. This generality is not needed for Rowhammer which does not allow “self-hammering”: the aggressor location cannot itself be a victim. Nevertheless, no additional axiom is required for the executions covered by our main theorem. Program well-formedness gives $x \in locs$ for every program access, and safe physical separation then gives $x \notin \text{adj}(br, locs, x)$ by instantiating protected–victim disjointness with the protected location x itself. Clause (F2) consequently ensures that the fault preserves the accessed location. Thus, the abstract contract is slightly more general than the intended DRAM interpretation, while the non-self-disturbance property required for protected executions is a derived consequence of the theorem hypotheses rather than a separate kernel axiom.

5.2 Physical adequacy and scope of our fault model

Our semantics is not a physically faithful device-level model of the electrical mechanisms underlying Rowhammer, nor is it intended to predict the probability of a bit-flip on a particular DRAM module. Our abstraction can be read as a conservative attacker model for Rowhammer-style disturbance: the bit-flips admissible by our model can be more extreme than those observed in physical DRAM. Therefore, a universal theorem proved for every admissible abstract kernel applies a fortiori to any more constrained concrete device model that refines it.

5.3 The fault model, mechanised

The Lean development renders the contract in three declarations. Separation is an abstract structure: a blast-radius function, a safety predicate, and the single fact the soundness proof uses about them: with physical separation, the victim set of an access to a protected location is disjoint from the protected set.

```

519 structure SeparationModel where
520   adj : Nat → LocSet → Loc → LocSet
521   isSafe : Nat → LocSet → Prop
522   /-- Safe layouts: the victim set of an access to a protected row is
523     disjoint from the protected set. -/
524   safe_disjoint : ∀ {br : Nat} {locs : LocSet} {x : Loc},
525     isSafe br locs → x ∈ locs → ∀ y ∈ locs, y ∉ adj br locs x

```

Our abstract fault contract does not build the row structure of DRAM directly into the programming language semantics. Instead, a separation model supplies an abstract victim-set function $\text{adj}(br, locs, x)$, identifying the locations that may be corrupted by an access to x . The theory does not require this function to be derived from a metric, or even to be symmetric. This is both simpler and more general than fixing a particular DRAM organisation.

The mechanisation includes a concrete distance-based instance in which locations are natural numbers and the victim set is determined by numerical distance. This instance is primarily a simple witness that the abstract contract is non-vacuous: numerical address distance should not itself be identified with physical DRAM distance unless locations have already been interpreted as physical row identifiers. A row-based model is a special case of the abstract interface: locations may be taken to denote rows directly, or the victim-set function may be defined through a mapping from program-visible locations to DRAM banks and rows.

A fault kernel is a function from victim sets and memories to distributions over memories. The *type* already encodes (F3) and (F4): the kernel sees neither history nor time.

```
abbrev ProbRHFlip (_sep : SeparationModel) :=
  LocSet → Memory → Dist Memory
```

Clause (F2) is the sole behavioural assumption, stated as a typeclass:

```
class HasProbFlipFrame (sep : SeparationModel) (f : ProbRHFlip sep) : Prop where
  outside_unchanged : ∀ (S : LocSet) (m m' : Memory) (x : Loc),
    m' ∈ (f S m).support → x ∉ S → m' x = m x
```

Clause (F1) is enforced by construction of the transition rules below.

The contract is not vacuous, and not satisfied only by fault-free kernels. The identity kernel `noFlip` (every fault is a no-op) satisfies it and embeds the ordinary semantics (Section 6.2). At the other extreme, so does the worst-case kernel that deterministically corrupts the *entire* victim set at every access:

```
noncomputable def corruptAll (sep : SeparationModel) : ProbRHFlip sep :=
  fun S m => dirac (fun x => if x ∈ S then m x + 1 else m x)
```

The mechanisation proves `corruptAll` frame-local and genuinely faulting (every outcome differs from the pre-fault memory at every victim location), and exhibits concrete physical separation of the distance model together with a well-formed program over it, so every hypothesis of the goal theorem is satisfiable simultaneously with a kernel that really flips bits.

5.4 Probabilistic transitions

The transition judgements

$$\langle e, m \rangle \rightsquigarrow_a^{\text{RH}} \mu, \quad \langle b, m \rangle \rightsquigarrow_b^{\text{RH}} \mu, \quad \langle P, m \rangle \rightsquigarrow_p^{\text{RH}} \mu$$

return a *distribution* μ over successor configurations rather than a single successor, and are otherwise the rules of Section 4 with the same left-to-right, non-short-circuit evaluation order. Exactly two rules sample the kernel, matching (F1). A variable read first obtains the value $m(x)$ and then samples the fault triggered by that access. The sampled fault therefore cannot retroactively change the value returned by the current read, but it may change the memory seen by later reads:

$$\frac{}{\langle x, m \rangle \rightsquigarrow_a^{\text{RH}} (\text{let } m' \leftarrow \text{pflip}(\text{adj}(br, locs, x), m) \text{ in } \delta(\langle m(x), m' \rangle))} \quad (\text{RH-A-READ})$$

Dually, once the right-hand side of an assignment is a value, the assignment writes it and samples the fault triggered by the write, on the *updated* memory:

$$\frac{}{\langle x := v, m \rangle \rightsquigarrow_p^{\text{RH}} (\text{let } m' \leftarrow \text{pflip}(\text{adj}(br, locs, x), m[x \mapsto v]) \text{ in } \delta(\langle \text{skip}, m' \rangle))} \quad (\text{RH-P-ASSIGN})$$

Under the well-formedness and physical separation hypotheses of the main theorem, the accessed location lies outside its victim set. Hence frame-locality preserves the stored value at a read location and preserves the value just installed by an assignment. This fact is derived from physical separation. It is not an additional assumption on the fault kernel. Every other rule is the Dirac or functorial lifting of its deterministic counterpart. The full rule sets for all three judgements are in Section B. There is no transition from $\langle \text{skip}, m \rangle$. The rules are syntax-directed, so every non-terminal program configuration determines a unique one-step distribution: the fault randomises the memory contents,

never which access happens next. The semantics is therefore probabilistic but not nondeterministic. We write $\text{step}_{br,locs}^{\text{RH}}(\langle P, m \rangle)$ for the program-level distribution.

5.5 Finite probabilistic executions

For iteration, define an absorbing extension of the one-step kernel by

- $\widehat{\text{step}}_{br,locs}^{\text{RH}}(\langle \text{skip}, m \rangle) = \delta(\langle \text{skip}, m \rangle)$,
- $\widehat{\text{step}}_{br,locs}^{\text{RH}}(\langle P, m \rangle) = \text{step}_{br,locs}^{\text{RH}}(\langle P, m \rangle)$ assuming that $P \neq \text{skip}$,

and the distribution after exactly n applications of the absorbing kernel by

- $\text{run}_{br,locs}^0(C) = \delta(C)$,
- $\text{run}_{br,locs}^{n+1}(C) = \text{let } C' \leftarrow \widehat{\text{step}}_{br,locs}^{\text{RH}}(C) \text{ in } \text{run}_{br,locs}^n(C')$.

The absorbing extension exists only to define finite runs. The operational semantics still has no transition from a terminal configuration. Because expression reductions are lifted through assignments and conditionals one step at a time, the index n counts genuine small-step transitions. A read or write together with its immediately triggered fault sample counts as one probabilistic transition.

In the mechanisation the absorbing kernel is a total *function* (stepRHHat), and the displayed judgement rules are inductive relations, that the two agree on non-terminal configurations (the kernel computes exactly the judgement's distribution, and the judgement's distribution is unique) is proved ($\text{pRHStep_iff_stepRHHat}$), so results stated over the runs below are statements about the displayed rules.

Because the rules are syntax-directed, *which* access a configuration performs next is a function of the configuration alone. The n -step run therefore extends canonically to a trace-carrying run $\text{run}_{br,locs, \text{tr}}^n$ that accumulates the label sequence of Section 4 alongside the configuration, assigning to every outcome its access trace. Absorbing steps at terminal configurations contribute τ , which the access trace erases. The goal theorem's trace agreement (Section 6.2) is stated over this run.

5.6 Mixed termination

Without mitigations, Rowhammer faults may change not only final values but also whether a program terminates. Indeed, a single initial configuration may induce a probabilistic mixture of terminating and non-terminating executions. Consider the following program:

$$\text{while } x > 17 \text{ do}(y := z + 1; x := x - 1).$$

In the ordinary semantics, every iteration decreases x by one, so the program terminates. Now suppose a Rowhammer bit-flip triggered by writing y may corrupt x . As a concrete illustration, consider a kernel that leaves x unchanged with probability $1 - p$, changes x to $x + 2$ with probability $0 < p < 1$ after the write to y , and otherwise leaves the locations relevant to the example unchanged. Hence the program terminates with positive probability and runs forever with positive probability. This example is deliberately schematic: it demonstrates behaviour permitted by the abstract fault contract, not a quantitative claim about a particular DRAM device.

For sufficiently many execution steps, some probability mass may already be on configurations $\langle \text{skip}, m' \rangle$, while the remaining mass is on non-terminal configurations. The exact probabilities of non-termination are available only as limits. We do not mechanise these limiting probabilities in the present development. The present finite step probability-mass-function approach is sufficient for the theorem we establish in later sections. An explicit account of divergence is possible, but would require additional semantic infrastructure.

5.7 Language independence

The language-parametricity of the model comes from the small semantic signature required by fault kernels. To instantiate the interface for a language, one must identify memories, locations, and the semantic clauses that perform memory actions. The ordinary clauses for those actions are then factored into an intended logical access followed by a call to the fault kernel determined by the accessed location. All other language constructs interact with Rowhammer only through the memory actions they eventually perform. Consequently, the interface can be applied directly to explicitly imperative languages and indirectly to non-imperative languages by applying it to an intermediate language, runtime semantics, or machine semantics in which memory actions are explicit.

The language-independence claim is also supported by recent work on reusable architecture semantics. ArchSem [31], for example, factors instruction set architecture semantics through a small algebraic-effect interface between ISA descriptions and memory models. In that interface, instruction semantics does not directly manipulate a monolithic global machine state. Instead, it emits outcomes such as memory reads and writes, barriers, cache operations, TLB operations, and exceptions. The memory-relevant part of the interface is particularly close to what our model requires: a memory request records an access kind, address, address space, size, and tag information, and the outcome interface contains distinguished memory-read and memory-write effects. Our Rowhammer interface can be understood as an analogous factorisation, but at a different semantic boundary. ArchSem factors ISA behaviour from the memory and concurrency model. We factor the physical fault mechanism from the programming language or machine semantics. To instantiate our model for a new language or architecture, one need not rebuild the Rowhammer theory around that language’s syntax. It suffices to identify the semantic events that correspond to memory accesses, compute the set of physical locations exposed by each such access, and interpret the access through the probabilistic fault kernel. In a WHILE language these events are source-level variable reads and writes, but in an architecture semantics such as ArchSem they would correspond instead to memory outcomes such as MemRead and MemWrite. This is the precise sense in which the interface is substantially independent of the target programming language: the language-specific part is the production of memory-access events, while the Rowhammer-specific part is a reusable interpretation of those events as possible probabilistic disturbances of physical memory.

6 Soundness of physical separation

This section addresses our RQ4 and proves the paper’s first main theorem: under the fault kernel contract of Section 5.1, physical separation eliminates Rowhammer faults from being observable in the results of program execution. We first state the goal informally, then define the projection being compared, give the formal statement together with its mechanised form, walk through every hypothesis, and sketch the proof.

The idea is simple: given a program P such that all variables $locs$ occurring in P are physically separated, meaning that they are all outside each other’s blast radius. Then the probabilistic semantics of $\langle P, m \rangle$ when restricted to $locs$ is the Dirac distribution of the deterministic run of $\langle P, m \rangle$. All formal ingredients are already defined, except comparing a deterministic with a probabilistic run of $\langle P, m \rangle$.

6.1 The protected view

We call any element of $(\text{Prog} \times \text{Memory}) \times \text{Acc}^*$ a *configuration with labels*. Write $\text{run}_{\text{det}, \text{tr}}^n(C) \in (\text{Prog} \times \text{Memory}) \times \text{Acc}^*$ for the deterministic n -step absorbing run with its accumulated label sequence (terminal configurations step to themselves, contributing τ , Section 4), and $\text{run}_{br, locs, \text{tr}}^n(C) \in$

Dist((Prog × Memory) × Acc*) for the trace-carrying Rowhammer run of Section 5.5. The *protected view* of a configuration with labels is

$$\pi_{locs}(\langle P, m \rangle, s) \triangleq (P, m \upharpoonright_{locs}, \text{acc}(s)),$$

the residual program, the restriction of the memory to the protected locations, and the access trace. In Lean, the restriction is a total map on the subtype of protected locations, and the access trace erases τ :

```
def protView (locs : LocSet) :
  (Prog × Memory) × List Access →
  Prog × ({x : Loc // x ∈ locs} → Val) × List Access :=
  fun r => (r.1.1, fun x => r.1.2 x.1, accessTrace r.2)
```

These three components are exactly what the program’s execution can depend on. What the projection discards, unprotected DRAM subject to Rowhammer outside *locs*, is exactly what a physical separation defence sacrifices to the fault, and a well-formed program never reads it.

6.2 Formal statement

We are now in a position to state the soundness theorem, both informally and with extracted Lean.

THEOREM 6.1 (SOUNDNESS OF PHYSICAL SEPARATION, FORMAL). *Let pflip satisfy the fault kernel contract (Section 5.1), let safe(br, locs) hold, and let wf_{locs}(P). Then for every initial memory m and every n ∈ ℕ,*

$$\text{Dist.map } \pi_{locs} \left(\text{run}_{br,locs,tr}^n(\langle P, m \rangle) \right) = \delta \left(\pi_{locs}(\text{run}_{\text{det,tr}}^n(\langle P, m \rangle)) \right).$$

Equivalently, by the support characterisation of the point mass (Section 3): every outcome in the support of the *n*-step Rowhammer run has the residual program, the protected memory, and the access trace of the *n*-step ordinary run. The mechanised statement quantifies in addition over every separation geometry:

```
def PhysicalSeparationSound : Prop :=
  ∀ (sep : SeparationModel) (f : ProbRHFlip sep),
  HasProbFlipFrame sep f →
  ∀ (br : Nat) (locs : LocSet),
  sep.isSafe br locs →
  ∀ (p : Prog), p.WellFormed locs →
  ∀ (m : Memory) (n : Nat),
  Dist.mapD (protView locs)
    (runRHTrace sep f br locs n (p, m))
  = Dist.dirac (protView locs (runDTrace n (p, m)))
```

It might be of interest to see that all hypotheses in the mechanised statement are needed, and we walk through them in order next.

- **sep : SeparationModel.** The theorem is universal over abstract separation geometries (Section 5.3): a blast-radius function *adj*, a safety predicate *safe*, and the disjointness fact *safe_disjoint* connecting them. The proof uses nothing else about the geometry, no row arithmetic, no linear layout, so the result applies to any placement scheme whose safety condition implies the disjointness fact, with the distance model of Section 5.3 as the motivating instance.

- 736 • **f : ProbRHFFlip sep.** The fault kernel. Its type, $\text{LocSet} \rightarrow \text{Memory} \rightarrow \text{Dist}(\text{Memory})$,
737 encodes contract clauses (F3) and (F4): the kernel can consult only the current victim set
738 and memory, and the same kernel governs every access. Nothing about the distribution
739 itself is fixed.
- 740 • **HasProbFlipFrame sep f.** Contract clause (F2), spatial confinement is *the sole behavioural*
741 *hypothesis*. Every memory in the support of $\text{pflip}(S, m)$ agrees with m outside S . Dropping
742 it makes the theorem false immediately: a kernel that flips a protected location on some
743 access satisfies everything else, and its Rowhammer run visibly diverges from the ordinary
744 one on protected memory.
- 745 • **br and locs.** The blast radius and the protected set. Both are abstract parameters threaded
746 through adj where locs is a *finite* set of locations (a Finset), the memory locations the
747 program uses for computation.
- 748 • **sep.isSafe br locs.** Physical separation, the condition the defence establishes, e.g., by
749 guard rows. Through safe_disjoint it yields the only fact the proof needs: the victim
750 set of an access to a protected location is disjoint from the protected set. Without it the
751 theorem fails even for kernels obeying the whole contract: accessing one protected location
752 may corrupt another protected location inside its blast radius ([Theorem 2.2](#)).
- 753 • **p.WellFormed locs.** Every location the program reads or writes is in locs ([Section 4](#)). This
754 clause makes separation compose with execution: it guarantees that every access the run
755 performs is an access to a *protected* location, so safe_disjoint applies to every victim set
756 the run generates. Without it, an access to some $z \notin \text{locs}$ triggers a fault on $\text{adj}(\text{br}, \text{locs}, z)$,
757 about which $\text{safe}(\text{br}, \text{locs})$ says nothing, the blast radius of an unprotected location may
758 well intersect locs ([Theorem 2.3](#)).
- 759 • **m and n.** The initial memory and the number of steps, both universally quantified. Quantify-
760 ing over every finite n is the relational-safety reading of [Theorem 6.1](#): it covers terminating
761 and diverging programs alike, with no limit construction.

762 It may also be interesting to see what is *not* assumed: no bound on flip probabilities (faults may
763 be certain), no independence between locations or between fault events, no unbiasedness, no
764 particular distribution. Within the blast radius the kernel may be adversarial. The guarantee is
765 distribution-free: it follows from spatial confinement and separation alone.

766 6.3 Corollaries

768 Two clauses of [Theorem 6.1](#) are worth recording as standalone consequences of the point mass
769 equation.

771 COROLLARY 6.2 (EXACT TERMINATION CORRESPONDENCE). *Under the hypotheses of [Theorem 6.1](#),*
772 *if $\langle P, m \rangle \xRightarrow[p]{s} \langle \text{skip}, m' \rangle$, then every configuration in the support of $\text{run}_{\text{br}, \text{locs}}^{|\text{s}|}(\langle P, m \rangle)$ is terminal*
773 *and agrees with m' on locs : the Rowhammer execution terminates after exactly the same number of*
774 *abstract-machine steps, with the same protected final memory (and, by [Theorem 6.1](#), the same access*
775 *trace).*

```
777 def TerminationCorrespondence : Prop :=
778   ∀ (sep : SeparationModel) (f : ProbRHFFlip sep),
779   HasProbFlipFrame sep f →
780   ∀ (br : Nat) (locs : LocSet),
781     sep.isSafe br locs →
782     ∀ (p : Prog), p.WellFormed locs →
783     ∀ (m m' : Memory) (s : List Access),
```

```

785 PTrace (p, m) s (Prog.skip, m') →
786   ∀ c ∈ (runRH sep f br locs s.length (p, m)).support,
787     c.1 = Prog.skip ∧ ∀ x ∈ locs, c.2 x = m' x
788
789

```

789 **COROLLARY 6.3 (DIVERGENCE PRESERVATION).** *Under the hypotheses of [Theorem 6.1](#), if $\langle P, m \rangle \uparrow\uparrow$, then for every n every configuration in the support of $\text{run}_{br,locs}^n(\langle P, m \rangle)$ is non-terminal (and agrees with the ordinary run on the protected state).*

```

794 def DivergencePreservation : Prop :=
795   ∀ (sep : SeparationModel) (f : ProbRHFFlip sep),
796     HasProbFlipFrame sep f →
797     ∀ (br : Nat) (locs : LocSet),
798       sep.isSafe br locs →
799       ∀ (p : Prog), p.WellFormed locs →
800         ∀ (m : Memory), Diverges p m →
801           ∀ (n : Nat),
802             ∀ c ∈ (runRH sep f br locs n (p, m)).support,
803               c.1 ≠ Prog.skip
804
805

```

805 6.4 Proof sketch

806 The proof is a probabilistic simulation whose invariant is *parametric in the preserved region*: a predicate `keep` on locations, instantiated only at the very end with membership in `locs`. It proceeds in three main steps.

810 *One-step invariant* (`stepRHHat_agree`). Suppose the Rowhammer and ordinary configurations share the same residual program and their memories agree on `keep`, and suppose (i) every location the program uses satisfies `keep`, and (ii) every fault sample triggered by an access the program can perform preserves `keep`-locations. Then every outcome of one application of the Rowhammer kernel again shares the residual program of the ordinary step and agrees with it on `keep`. The residual program is preserved exactly, not just up to `keep`, because reads return the same values on `keep`-agreeing memories, so both machines make the same control decisions.

817 *Iteration* (`runRHTrace_protected`). The one-step invariant is iterated over the step count n . Because the rules are syntax-directed, the label of the next transition is a function of the residual program alone. Preserving the residual program exactly therefore forces both runs to emit the *same label sequence*, which gives the access-trace clause for free. Well-formedness is preserved along the run, so clauses (i) and (ii) remain available at every step.

822 *Instantiation and conversion.* Take $\text{keep}(\ell) \triangleq \ell \in \text{locs}$. Clause (i) is program well-formedness. Clause (ii) is exactly frame-locality (F2) composed with `safe_disjoint`: an access to $x \in \text{locs}$ triggers a fault confined to $\text{adj}(br, \text{locs}, x)$, which is disjoint from `locs`, so every sample fixes every protected location. The invariant then says that every outcome in the support of $\text{run}_{br,locs,tr}^n$ has the protected view of $\text{run}_{det,tr}^n$. The support characterisation of [Section 3](#) converts this into the point mass equation. The corollaries follow from the central equation together with determinism, progress, and the trace realisation lemma of the ordinary semantics ([Section A](#)). The no-fault embedding is by computation.

830 We can now show that the deterministic semantics is a special case of the probabilistic semantics. This gives us confidence in the correctness of the probabilistic semantics.

THEOREM 6.4 (NO-FAULT EMBEDDING). *Let sep be any separation model, let br be any blast-radius parameter, let $locs$ be any protected set, and let c be any initial configuration. If the Rowhammer fault kernel is the identity kernel $noFlip(sep)$, then n steps with the Rowhammer semantics is exactly the Dirac distribution concentrated on the n -step deterministic trace run: This requires neither safe physical separation nor program well-formedness.*

This theorem is mechanised in Lean as `runRHTrace_noFlip`.

7 Information flow security under Rowhammer

We now address our RQ5 whether the well-developed and widely used theory of information flow [10, 34] can be lifted to our Rowhammer model. We do this by defining *relative non-interference* and prove that physical separation guarantees relative non-interference. Relative non-interference means that any information leak that happens with Rowhammer also happens without Rowhammer². This is of interest, as Rowhammer is usually presented as an integrity failure: an access to one physical row may corrupt data stored in another. But the same mechanism can also create a confidentiality failure: a Rowhammer induced flip in a high location may alter a branch condition, loop condition, or value that is later copied to low memory. Moreover, a fault in a low location may directly change an attacker-visible result. The theory of information flow is tailor-made for studying such phenomena, and provides a natural test of whether physical separation removes the semantic consequences of Rowhammer rather than merely making faults less likely.

Low equivalence and non-interference. In order to explain relative non-interference in detail, we need to sketch the key ideas behind information flow security, see [10, 34] for details. Fix a policy that classifies every memory location as either L or H. Two memories are *low equivalent* when they agree at every L location (they may differ arbitrarily at H locations). A program is *non-interfering* when executions from low-equivalent initial memories end with memories that are also low-equivalent (we discuss termination later). Non-interference is therefore a *hyperproperty* [3]: it relates pairs of executions rather than inspecting one execution in isolation. The ordinary semantics is deterministic, so indistinguishability means equality of the two observations. The Rowhammer semantics maps an initial memory to a probability distribution over executions. Probability-sensitive non-interference consequently requires equality of the induced distributions over low observations.

What physical separation can and cannot guarantee. The physical separation discussed in previous sections must protect every location used by the program, not only its low locations. Protecting only low variables is insufficient: corrupting a high variable may change control flow and thereby alter a later low write, low access, or termination behaviour. The separation condition used in this paper therefore protects the complete set of locations over which the program is well formed. That means we can only give a *no-new-leaks* guarantee: physical separation does not make a program secure that is insecure in the deterministic semantics! For example,

$$L := H$$

explicitly leaks a high value into a low variable and is therefore interfering in both semantics. Physical separation only ensures that the Rowhammer semantics has exactly the same protected observations as the ordinary semantics. Hence, with physical separation, a Rowhammer information flow counterexample is only a counterexample that was already present without faults. Without

²Cf. relative consistency results like Gödel's proof that classical logic does not introduce unsoundness not already present in intuitionistic logic [12], or the Gödel / Cohen results that ZFC set theory is relatively consistent over ZF [4, 9].

883 separation this reasoning fails. If an access to a program's location can disturb another location, a
 884 frame-local fault kernel may change a low variable or a high value controlling future low behaviour.
 885 The fault kernel may assign arbitrary probabilities to such corruptions inside the victim set. Safe
 886 separation rules out this channel structurally: the victim set generated by an access to a protected
 887 location is disjoint from every protected location.

888 8 Mechanised observation-parametric non-interference

889 This section formalises and proves our relative non-interference results. The Rowhammer semantics
 890 of previous sections remains unchanged. Our central design choice is that every observer consumes
 891 the protected view already covered by the physical separation theorem. Thus the security results
 892 are consequences of semantic collapse proven already, rather than a second operational simulation
 893 argument. This gives us additional confidence in our mathematical rendering of Rowhammer as a
 894 programming language feature.
 895

896 8.1 Formal development

897 *Observation-parametric security.* There is no single universally appropriate notion of low obser-
 898 vation (e.g., is termination low observable). The mechanised development therefore parameterises
 899 non-interference by a function on the existing protected view, which contains the residual program,
 900 protected memory, and access trace. We define four concrete observers:

- 901 • Protected low memory only
- 902 • The address trace of reads and writes to low locations
- 903 • Terminal status, protected low memory, and the low access trace
- 904 • Protected low memory and the low access trace, but not terminal status

905 The third observer is progress-sensitive because observations are compared at every finite number
 906 of steps: it can reveal the exact abstract-machine step at which the residual program first becomes
 907 skip. Omitting the terminal-status bit does not by itself yield classical termination-insensitive
 908 non-interference, since low memory or low accesses observed at successive steps may still reveal
 909 progress. The mechanisation therefore makes no claim that its fourth observer is a complete
 910 formalisation of classical termination-insensitive non-interference.
 911

912 *Security policies and low equivalence, protected views and observations.* Let $\text{Sec} \triangleq \{\text{L}, \text{H}\}$ and let
 913 a security policy be a function $\gamma : \text{Loc} \rightarrow \text{Sec}$. The low and high location subtypes are $L_\gamma \triangleq \{x :$
 914 $\text{Loc} \mid \gamma(x) = \text{L}\}$, and $H_\gamma \triangleq \{x : \text{Loc} \mid \gamma(x) = \text{H}\}$. For $m \in \text{Memory}$, the two restricted views are

- 915 • $\text{lowView}_\gamma(m)(x) \triangleq m(x)$ assuming $(x \in L_\gamma)$
- 916 • $\text{highView}_\gamma(m)(x) \triangleq m(x)$ assuming $x \in H_\gamma$.

917 Two memories are low equivalent when

$$918 m_1 \approx_\gamma m_2 \iff \forall x : \text{Loc}. \gamma(x) = \text{L} \implies m_1(x) = m_2(x).$$

919 The mechanisation proves that this pointwise definition is equivalent to

$$920 \text{lowView}_\gamma(m_1) = \text{lowView}_\gamma(m_2),$$

921 and proves reflexivity, symmetry, and transitivity. Let Access contain labels $\text{read}(x)$, $\text{write}(x)$, and
 922 τ . The function accessTrace filters out τ . For the finite protected set locs , define

$$923 \text{PView}(\text{locs}) \triangleq \text{Prog} \times (\{x : \text{Loc} \mid x \in \text{locs}\} \rightarrow \text{Val}) \times \text{List}(\text{Access}).$$

924 The protected projection is

$$925 \text{protView}_{\text{locs}}((Q, m), s) \triangleq (Q, m|_{\text{locs}}, \text{accessTrace}(s)).$$

It retains exactly the three components covered by physical semantic collapse: residual syntax, protected memory, and the access trace. The deterministic protected view at step count n is

$$\text{DView}_{locs}^n(P, m) \triangleq \text{protView}_{locs}(\text{runDTrace}^n(P, m)).$$

The Rowhammer protected-view distribution is

$$\text{RHView}_{sep,f,br,locs}^n(P, m) \triangleq \text{Dist.map protView}_{locs}(\text{runRHTrace}_{sep,f,br,locs}^n(P, m)).$$

Non-interference definitions. Recall that we do not hard-code just one specific notion of how the program execution is observed. Instead, let Ω be an observation type. A protected observer is a function

$$O : \text{PView}(locs) \rightarrow \Omega.$$

The two observations used in the definitions below are

$$\text{ObsOrd}_{locs,O}^n(P, m) \triangleq O(\text{DView}_{locs}^n(P, m))$$

and

$$\text{ObsRH}_{sep,f,br,locs,O}^n(P, m) \triangleq \text{Dist.map } O(\text{RHView}_{sep,f,br,locs}^n(P, m)).$$

See the Lean definitions `rowhammerProtectedView` and `rowhammerObservation` that match this two-stage definition. We can now define our notion of observer parametric non-interference.

Definition 8.1 (Ordinary step bound non-interference). The predicate `OrdinaryNI`($\gamma, locs, O, P$) holds when

$$\forall n \in \mathbb{N}. \forall m_1, m_2 \in \text{Memory}. \quad m_1 \approx_\gamma m_2 \implies \text{ObsOrd}_{locs,O}^n(P, m_1) = \text{ObsOrd}_{locs,O}^n(P, m_2).$$

This is the Lean definition `OrdinaryNI`.

Definition 8.2 (Rowhammer step bound non-interference). For a fixed separation model sep , fault kernel f , blast radius br , and protected set $locs$, the predicate `RowhammerNI`($\gamma, sep, f, br, locs, O, P$) holds when

$$\forall n \in \mathbb{N}. \forall m_1, m_2 \in \text{Memory}. \quad m_1 \approx_\gamma m_2 \implies \text{ObsRH}_{sep,f,br,locs,O}^n(P, m_1) = \text{ObsRH}_{sep,f,br,locs,O}^n(P, m_2).$$

The equality is equality of complete finite distributions. The same fixed kernel f is used in both runs. This is the Lean definition `RowhammerNI`.

Definition 8.3 (Robust Rowhammer non-interference). The predicate `RobustRowhammerNI`($\gamma, sep, br, locs, O, P$) holds when

$$\forall f. \text{HasProbFlipFrame}(sep, f) \implies \text{RowhammerNI}(\gamma, sep, f, br, locs, O, P).$$

Safety of the layout and well-formedness of P are deliberately not fields of this definition, they are hypotheses of the preservation theorem. This matches the Lean definition `RobustRowhammerNI`.

For convenience, the mechanisation also defines

$$\begin{aligned} \text{Admissible}(sep, f, br, locs, P) &\iff \text{HasProbFlipFrame}(sep, f) \\ &\quad \wedge \text{sep.isSafe}(br, locs) \\ &\quad \wedge P.\text{WellFormed}(locs). \end{aligned}$$

The checked theorem statements below take these three assumptions as separate arguments, exactly as their Lean counterparts do. We can now define the four concrete checked observers, promised above. Towards that aim, let

$$L_{\gamma,locs} \triangleq \{x : \text{Loc} \mid x \in locs \wedge \gamma(x) = \text{L}\}.$$

The development defines the following protected observers.

- **Protected-low-memory observer.** It hides residual syntax, terminal status, and the access trace:

$$O_{\gamma,locs}^{\text{mem}}(Q, \bar{m}, s) \triangleq \bar{m}|_{L_{\gamma,locs}}.$$

- **Low-access observer.** Since labels contain addresses but not transferred values, this observer captures address-trace leakage rather than value-labelled I/O leakage.

$$O_{\gamma}^{\text{access}}(Q, \bar{m}, s) \triangleq \text{lowAccessTrace}_{\gamma}(s),$$

where the filter retains $\text{read}(x)$ and $\text{write}(x)$ exactly when $\gamma(x) = L$.

- **Progress-sensitive observer.** Because the definition compares every finite run, this observer reveals exact abstract-machine progress, including the first step-count at which termination is visible.

$$O_{\gamma,locs}^{\text{progress}}(Q, \bar{m}, s) \triangleq (\text{terminal}(Q), \bar{m}|_{L_{\gamma,locs}}, \text{lowAccessTrace}_{\gamma}(s)),$$

where $\text{terminal}(Q)$ is the Boolean that is true exactly for $Q = \text{skip}$.

- **Terminal-status-hidden observer.** This omits the terminal-status bit but is not claimed to formalise classical termination-insensitive non-interference: low memory or low accesses at fixed step count may still reveal progress.

$$O_{\gamma,locs}^{\text{hidden}}(Q, \bar{m}, s) \triangleq (\bar{m}|_{L_{\gamma,locs}}, \text{lowAccessTrace}_{\gamma}(s)).$$

Checked theorem statements. The following statements are proved in Lean without sorry and without new problem-specific axioms.

THEOREM 8.4 (PROTECTED SEMANTIC COLLAPSE). *Let f be frame-local, let $\text{sep.isSafe}(br, locs)$ hold, and let $P.\text{WellFormed}(locs)$ hold. Then, for every memory m and n ,*

$$\text{RHView}_{\text{sep},f,br,locs}^n(P, m) = \delta(\text{DView}_{locs}^n(P, m)).$$

This is theorem `protected_semantic_collapse`, a direct restatement of the existing theorem `physical_separation_sound`.

THEOREM 8.5 (OBSERVATION COLLAPSE). *Under the same frame-locality, safety, and well-formedness assumptions, for every protected observer O , memory m , and step count n ,*

$$\text{ObsRH}_{\text{sep},f,br,locs,O}^n(P, m) = \delta(\text{ObsOrd}_{locs,O}^n(P, m)).$$

This is theorem `observation_collapse`.

THEOREM 8.6 (NON-INTERFERENCE EQUIVALENCE). *Let f be frame-local, let $\text{sep.isSafe}(br, locs)$ hold, and let $P.\text{WellFormed}(locs)$ hold. Then, for every security policy γ and protected observer O ,*

$$\text{OrdinaryNI}(\gamma, locs, O, P) \iff \text{RowhammerNI}(\gamma, \text{sep}, f, br, locs, O, P).$$

This is theorem `noninterference_equivalence`.

COROLLARY 8.7 (ROBUST PRESERVATION). *If $sep.isSafe(br, locs)$, $P.WellFormed(locs)$, and OrdinaryNI($\gamma, locs, O, P$) hold, then*

$$\text{RobustRowhammerNI}(\gamma, sep, br, locs, O, P).$$

Equivalently, every frame-local kernel f satisfies

$$\text{RowhammerNI}(\gamma, sep, f, br, locs, O, P).$$

This is theorem `robust_noninterference_preservation`.

COROLLARY 8.8 (REFLECTION FOR A FIXED KERNEL). *Under frame-locality, safety, and well-formedness,*

$$\neg \text{RowhammerNI}(\gamma, sep, f, br, locs, O, P) \implies \neg \text{OrdinaryNI}(\gamma, locs, O, P).$$

This is theorem `rowhammer_interference_reflection`.

COROLLARY 8.9 (EXISTENTIAL REFLECTION). *Assume $sep.isSafe(br, locs)$ and $P.WellFormed(locs)$.*

Then

$$\begin{aligned} & (\exists f. \text{HasProbFlipFrame}(sep, f) \wedge \neg \text{RowhammerNI}(\gamma, sep, f, br, locs, O, P)) \\ & \implies \neg \text{OrdinaryNI}(\gamma, locs, O, P). \end{aligned}$$

In the Lean code, this is theorem `exists_rowhammer_interference_reflection`.

8.2 Proof sketches

Structural preservation. The existing lemma `stepDHat_wf` proves that one absorbing deterministic step preserves well-formedness over $locs$. Its proof follows the syntax of the residual program. Expression reduction preserves expression well-formedness. Sequencing preserves the well-formedness of the untouched continuation. A conditional either reduces its guard or selects an already well-formed branch. Unfolding a loop constructs a conditional and sequence from already well-formed components. This lemma is what allows the step bound simulation invariant to be applied recursively to the successor program.

One-step agreement. The generic relation $\text{Agree}(keep, m_f, m_d)$ states that the faulty and deterministic memories agree at every location satisfying $keep$. Lemma `stepRHHat_agree` assumes that every location used by the program is kept and that every supported fault outcome preserves all kept locations. If the input memories agree on $keep$, then every supported Rowhammer successor has exactly the same residual program as the deterministic successor and its memory still agrees with the deterministic memory on $keep$. The proof is by structural induction on the program, using corresponding agreement lemmas for arithmetic and Boolean expressions. In a read or write case, frame preservation ensures that the sampled fault cannot alter a kept location. Since values read from kept locations agree, both semantics reduce the same expression, write the same value, select the same branch, and unfold the same control construct.

Finite step bound invariant. Lemma `runRHTrace_protected` iterates one-step agreement by induction on n . For every outcome r in the support of the n -step Rowhammer trace run, it proves three facts:

- The residual program of r equals the deterministic residual program.
- The memories agree on every kept location.
- The complete accumulated label lists are equal.

In the successor case, `stepRHHat_agree` supplies the next related configuration, `stepDHat_wf` supplies well-formedness of its residual program, and the induction hypothesis supplies the remaining step bound.

1079 *Instantiation by physical separation.* For protected semantic collapse, the proof instantiates
 1080 $keep(x)$ with $x \in locs$. The premise that faults preserve kept locations is derived by combin-
 1081 ing frame-locality with safe separation: for an access to $x \in locs$, `safe_disjoint` shows that
 1082 every $\ell \in locs$ lies outside the victim set, and `outside_unchanged` therefore preserves ℓ in ev-
 1083 ery supported fault outcome. The finite step bound invariant then shows that every supported
 1084 Rowhammer outcome has the same protected view as the deterministic run. The finite-distribution
 1085 lemma `mapD_eq_dirac_of_support` turns this support-wise fact into the point mass equation of
 1086 Theorem 8.4.

1087 *Observation collapse.* Unfolding the two observation definitions exposes a map of O over the
 1088 Rowhammer protected-view distribution. Rewriting that distribution with Theorem 8.4 leaves O
 1089 mapped over a Dirac distribution. The finite-distribution law

$$1091 \text{Dist.map } O (\delta(v)) = \delta(O(v))$$

1092 gives Theorem 8.5.

1093 *Non-interference equivalence and its corollaries.* For low-equivalent m_1 and m_2 , observation
 1094 collapse rewrites the two Rowhammer observation distributions to

$$1095 \delta(\text{ObsOrd}_{locs,O}^n(P, m_1)) \quad \text{and} \quad \delta(\text{ObsOrd}_{locs,O}^n(P, m_2)).$$

1096 Ordinary non-interference implies equality of the two points and hence congruence of their Dirac
 1097 distributions. Conversely, equality of the two Dirac distributions implies equality of their points
 1098 by injectivity of δ . This proves Theorem 8.6. Robust preservation introduces an arbitrary frame-
 1099 local kernel and applies the forward direction. Fixed-kernel and existential reflection are direct
 1100 contrapositives of that same direction.

1103 8.3 No-fault sanity results

1104 The checked development also proves, without any safety or well-formedness assumption, that the
 1105 no-fault kernel embeds deterministic execution exactly:

$$1106 \text{runRHTrace}_{sep, \text{noFlip}(sep), br, locs}^n(c) = \delta(\text{runDTrace}^n(c)).$$

1107 The proof is an induction on n using the existing one-step no-fault theorem, Dirac bind, and
 1108 mapping over a Dirac distribution. It yields

$$1109 \text{ObsRH}_{sep, \text{noFlip}(sep), br, locs, O}^n(P, m) = \delta(\text{ObsOrd}_{locs, O}^n(P, m))$$

1110 and therefore

$$1111 \text{OrdinaryNI}(\gamma, locs, O, P) \iff \text{RowhammerNI}(\gamma, sep, \text{noFlip}(sep), br, locs, O, P).$$

1112 In Lean they are `runRHTrace_noFlip`, `noFault_observation`, and `ordinaryNI_iff_noFaultNI`.

1117 9 Conclusion

1118 This paper has given an abstract, compositional operational semantics for a representative program-
 1119 ming language in the presence of Rowhammer. The semantics treats Rowhammer as a probabilistic
 1120 effect on memory accesses, rather than as a direct model of DRAM organisation or semiconductor
 1121 physics. The no-fault embedding theorem, Theorem 6.4, recovers the ordinary deterministic se-
 1122 mantics as the point mass special case of the Rowhammer semantics, answering RQ1. The model
 1123 isolates the Rowhammer effect as a small generalisation from state transformers to probabilistic
 1124 state transformers: memory reads and writes are intercepted by a fault kernel, while the rest of the
 1125 operational semantics is lifted compositionally using the probability monad. This gives the small,
 1126 conceptually non-intrusive interface sought in RQ2. The interface is also largely independent of

1128 the ambient programming language. The development uses only that programs perform memory
 1129 reads and writes, and that the ordinary operational semantics can expose those accesses. This is
 1130 not by itself a formal semantics for every language or machine model, but it gives a clear route
 1131 for instantiating the same Rowhammer interface in other languages, intermediate representations,
 1132 or instruction-set semantics. In this sense, the paper affirms RQ3. The semantic-collapse theorem,
 1133 **Theorem 6.1**, proves the soundness of physical separation: if every protected program access has a
 1134 victim set disjoint from the protected locations, then every admissible Rowhammer kernel induces
 1135 exactly the deterministic protected behaviour. This justifies physical separation as a Rowhammer
 1136 defence within the abstract fault contract, answering RQ4. Finally, the information flow results show
 1137 that the abstraction supports standard semantic security reasoning. In particular, Theorems 8.4, 8.5
 1138 and 8.6 transport non-interference reasoning from ordinary executions to Rowhammer-affected
 1139 executions, and back again, under the same physical separation hypotheses. Thus the model does
 1140 not merely describe faulty executions, it makes Rowhammer amenable to established semantic
 1141 security frameworks, addressing RQ5.

1142

1143 10 Future work

1144

1145 *Refining the Rowhammer model.* The fault kernel contract history-freeness (F3) and stationarity
 1146 (F4) exclude genuinely history-dependent faults such as charge accumulation across repeated
 1147 activations and refresh-cycle effects. Generalising requires extending the Rowhammer state with
 1148 an activation history, refresh state, or other hardware state. Our theorems would then quantify
 1149 over arbitrary stateful kernels satisfying the same protected-region preservation condition, and
 1150 we expect the parametric-invariant proof architecture to carry over. Moreover, DRAM vendors
 1151 typically deploy proprietary, undocumented mitigations, most notably Target Row Refresh (TRR),
 1152 whose opaque nature complicates formal security guarantees. Bridging the gap between our formal
 1153 models and these black-box defences requires rigorous experimental reverse engineering. Pioneering
 1154 methodologies such as the U-TRR framework [13], uncover the internal mechanics of TRR by taking
 1155 data-retention failures as a side channel. Alongside this, complementary characterisation of DRAM's
 1156 disturbance sensitivities [30] establishes clear empirical boundaries on when such defences succeed
 1157 or fail, and in turn informs principled mitigations such as BlockHammer [40]. By experimentally
 1158 exposing the true behaviour and limitations of in-DRAM mitigations, such efforts provide the
 1159 precise physical parameters and fault models that a formal framework like ours needs in order to
 1160 deliver end-to-end software security guarantees.

1160

1161 *Deploying physical separation.* Our soundness theorem is conditional on the data placement
 1162 satisfying physical separation. Establishing that on modern hardware is nontrivial: DRAM vendors
 1163 do not publish row geometry, in-DRAM address remapping obscures physical adjacency, and the
 1164 effective blast radius must be measured per device. Empirical work on learning the Rowhammer
 1165 parameters of a given device [28] is complementary to our results: it supplies the physical facts that
 1166 instantiate *br* and *safe*, while the theorem turns those facts into a semantic guarantee. The separation
 1167 between the abstract geometry (`SeparationModel`) and the semantics in our mechanisation is
 1168 designed exactly for this division of labour.

1169

1170 *Probabilistic non-termination.* A natural extension is to complement the finite step bound seman-
 1171 tics with an explicit account of eventual termination and divergence. This could be done either by
 1172 taking limits of the finite step bound termination probabilities, by developing a sub-probability
 1173 semantics over final memories, or by constructing a probability measure on infinite execution paths.
 1174 A fully event-based treatment, in which divergence is identified with a measurable set of infinite
 1175 executions, may additionally require the construction of an appropriate probability measure on
 1176

1176

1177 infinite paths and an adequacy theorem relating that measure to the n -step operational semantics.
1178 Each route introduces nontrivial order-theoretic or measure-theoretic mechanisation.

1179 *Connections with Shannon-Scott information.* A further direction is to relate the present proba-
1180 bilistic account of Rowhammer to semantic accounts of information based on both Shannon entropy
1181 and Scott domains. Shannon's theory of communication provides a natural language for noisy
1182 channels and probabilistic leakage [36], while domain-theoretic and topological methods are com-
1183 monly used in programming language semantics to model partial information and non-terminating
1184 computation. Recent work has connected these two perspectives in the setting of information
1185 flow and non-interference [14]. Our results compare the distributions of low observations after a
1186 fixed number of execution steps. It would be interesting to investigate whether this account can
1187 be extended to an infinite-horizon semantics in which Rowhammer-induced probabilistic leakage,
1188 partial observations, and non-termination are treated in a single Shannon-Scott framework.

1189 **Acknowledgements**

1190 Chunyan Mu drew our attention to Rowhammer being an open problem from a non-interference angle. We
1191 thank David Clark for numerous discussions about information flow and non-interference, Davide Bartolini
1192 and Xubin Tan for helping us to understand memory controllers, and Lev Mukhanov for alerting us to the
1193 difficulty of reproducing the Rowhammer attacks. Fredrik Dahlqvist helped us better to understand Kozen-style
1194 probabilistic semantics of probabilistic programming languages. We thank Yusuke Izawa and Mohammad
1195 M. Ahmadpanah for helpful comments, which greatly improved the presentation of this paper.

1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225

References

- [1] Anonymous. 2026. Lean proofs for "Mechanised operational semantics of Rowhammer". https://anonymous.4open.science/r/rh_separation_paper_anon_lean-48C7/.
- [2] Ferdinand Brasser, Lucas Davi, David Gens, Christopher Liebchen, and Ahmad-Reza Sadeghi. 2017. CAN't Touch This: Software-only Mitigation against Rowhammer Attacks Targeting Kernel Memory. In *26th USENIX Security Symposium (USENIX Security 17)*. USENIX Association, Vancouver, BC, Canada, 117–130. <https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/brasser>
- [3] Michael R. Clarkson and Fred B. Schneider. 2010. Hyperproperties. *Journal of Computer Security* 18, 6 (2010), 1157–1210. doi:10.3233/JCS-2009-0393
- [4] Paul J. Cohen. 1963. The Independence of the Continuum Hypothesis. *Proceedings of the National Academy of Sciences* 50, 6 (1963), 1143–1148. doi:10.1073/pnas.50.6.1143
- [5] Fredrik Dahlqvist, Alexandra Silva, and Dexter Kozen. 2020. *Semantics of Probabilistic Programming: A Gentle Introduction*. Cambridge University Press, 1–42. doi:10.1017/9781108770750.002
- [6] Leonardo de Moura and Sebastian Ullrich. 2021. The Lean 4 Theorem Prover and Programming Language. In *Automated Deduction – CADE 28*. Springer, 625–635. doi:10.1007/978-3-030-79876-5_37
- [7] Pietro Frigo, Emanuele Vannacci, Hasan Hassan, Victor van der Veen, Onur Mutlu, Cristiano Giuffrida, Herbert Bos, and Kaveh Razavi. 2020. TRRespass: Exploiting the Many Sides of Target Row Refresh. In *2020 IEEE Symposium on Security and Privacy (SP)*. 747–762. doi:10.1109/SP40000.2020.00090
- [8] Michèle Giry. 1982. A categorical approach to probability theory. In *Categorical Aspects of Topology and Analysis*. Lecture Notes in Mathematics, Vol. 915. Springer, 68–85. doi:10.1007/BFb0092872
- [9] Kurt Gödel. 1940. *The Consistency of the Continuum-Hypothesis*. Number 3 in Annals of Mathematics Studies. Princeton University Press, Princeton, NJ. doi:10.1515/9781400881635
- [10] Joseph A. Goguen and José Meseguer. 1982. Security Policies and Security Models. In *Proceedings of the 1982 IEEE Symposium on Security and Privacy*. 11–20. doi:10.1109/SP.1982.10014
- [11] Daniel Gruss, Clémentine Maurice, and Stefan Mangard. 2016. Rowhammer.js: A Remote Software-Induced Fault Attack in JavaScript. In *Detection of Intrusions and Malware, and Vulnerability Assessment (DIMVA)*. Springer, 300–321. doi:10.1007/978-3-319-40667-1_15
- [12] Kurt Gödel. 1933. On Intuitionistic Arithmetic and Number Theory. *Ergebnisse eines mathematischen Kolloquiums 4* (1933), 34–38. doi:10.1093/oso/9780195147209.003.0063 English translation by Stefan Bauer-Mengelberg and Jean van Heijenoort. Reprinted in: *Kurt Gödel: Collected Works, Volume I*, Oxford University Press (1986).
- [13] Hasan Hassan, Yahya Can Tuğrul, Jeremie S. Kim, Victor van der Veen, Kaveh Razavi, and Onur Mutlu. 2021. Uncovering In-DRAM RowHammer Protection Mechanisms: A New Methodology, Custom RowHammer Patterns, and Implications. In *54th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO-54)*. doi:10.1145/3466752.3480110
- [14] Sebastian Hunt, David Sands, and Sandro Stucki. 2023. Reconciling Shannon and Scott with a Lattice of Computable Information. *Proc. ACM Program. Lang.* 7, POPL, Article 68 (jan 2023), 30 pages. doi:10.1145/3571740
- [15] Patrick Jattke, Victor van der Veen, Pietro Frigo, Stijn Gunter, and Kaveh Razavi. 2022. BLACKSMITH: Scalable Rowhammering in the Frequency Domain. In *2022 IEEE Symposium on Security and Privacy (SP)*. 716–734. doi:10.1109/SP46214.2022.9833772
- [16] Yoongu Kim, Ross Daly, Jeremie Kim, Chris Fallin, Ji Hye Lee, Donghyuk Lee, Chris Wilkerson, Konrad Lai, and Onur Mutlu. 2014. Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors. In *Proceedings of the 41st Annual International Symposium on Computer Architecture (ISCA '14)*. 361–372. doi:10.1109/ISCA.2014.6853210
- [17] Andreas Kogler, Jonas Juffinger, Sulaiman Qazi, Yoongu Kim, Moritz Lipp, Nicolas Boichat, Eric Shiu, Matt Nissler, and Daniel Gruss. 2022. Half-Double: Hammering From the Next Row Over. In *31st USENIX Security Symposium (USENIX Security 22)*. USENIX Association, 3807–3824. <https://www.usenix.org/conference/usenixsecurity22/presentation/kogler-half-double>
- [18] Radhesh Krishnan Konoth, Marco Oliverio, Andrei Tatar, Dennis Andriesse, Herbert Bos, Cristiano Giuffrida, and Kaveh Razavi. 2018. ZebRAM: Comprehensive and Compatible Software Protection Against Rowhammer Attacks. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*. USENIX Association. <https://www.usenix.org/conference/osdi18/presentation/konoth>
- [19] Dexter Kozen. 1981. Semantics of Probabilistic Programs. *J. Comput. System Sci.* 22, 3 (1981), 328–350. doi:10.1016/0022-0000(81)90036-2
- [20] Andrew Kwong, Daniel Genkin, Daniel Gruss, and Yuval Yarom. 2020. RAMBleed: Reading Bits in Memory Without Accessing Them. In *IEEE Symposium on Security and Privacy (S&P)*. IEEE, 695–711. doi:10.1109/SP40000.2020.00020
- [21] Francis William Lawvere. 1962. The Category of Probabilistic Mappings: With Applications to Stochastic Processes, Statistics, and Pattern Recognition. (1962). <https://lawverearchives.com/wp-content/uploads/2025/07/1962.probmap.pdf> Manuscript, 12 pages.

- 1275 [22] Saunders Mac Lane. 1998. *Categories for the Working Mathematician* (2 ed.). Graduate Texts in Mathematics, Vol. 5.
1276 Springer, New York. doi:10.1007/978-1-4757-4721-8
- 1277 [23] John McCarthy. 1962. Towards a Mathematical Science of Computation. In *Information Processing 1962: Proceedings of*
1278 *IFIP Congress 62*, Cicely M. Popplewell (Ed.). North-Holland, Amsterdam, 21–28. [https://www-formal.stanford.edu/](https://www-formal.stanford.edu/jmc/towards.pdf)
1279 [jmc/towards.pdf](https://www-formal.stanford.edu/jmc/towards.pdf)
- 1280 [24] Eugenio Moggi. 1991. Notions of Computation and Monads. *Information and Computation* 93, 1 (1991), 55–92.
1281 doi:10.1016/0890-5401(91)90052-4
- 1282 [25] Peter D. Mosses. 2004. Modular Structural Operational Semantics. *Journal of Logical and Algebraic Methods in*
1283 *Programming* 60–61 (2004), 195–228. doi:10.1016/j.jlap.2004.03.008
- 1284 [26] Onur Mutlu and Jeremie S. Kim. 2020. RowHammer: A Retrospective. *IEEE Transactions on Computer-Aided Design of*
1285 *Integrated Circuits and Systems* 39, 8 (2020), 1555–1571. doi:10.1109/TCAD.2019.2915318
- 1286 [27] Amir Naseredini. 2024. *Towards Automatic Analysis of Microarchitectural Attacks*. Ph. D. Dissertation. University
1287 of Sussex. [https://sussex.figshare.com/articles/thesis/Towards_automatic_analysis_of_microarchitectural_attacks/](https://sussex.figshare.com/articles/thesis/Towards_automatic_analysis_of_microarchitectural_attacks/24610929)
1288 [24610929](https://sussex.figshare.com/articles/thesis/Towards_automatic_analysis_of_microarchitectural_attacks/24610929)
- 1289 [28] Amir Naseredini, Martin Berger, Matteo Sammartino, and Shale Xiong. 2023. ALARM: Active LeArning of Rowhammer
1290 Mitigations. In *Proceedings of the 11th International Workshop on Hardware and Architectural Support for Security and*
1291 *Privacy* (Chicago, IL, USA) (HASP '22). Association for Computing Machinery, New York, NY, USA, 1–9. doi:10.1145/
1292 [3569562.3569563](https://doi.org/10.1145/3569562.3569563)
- 1293 [29] Hanne Riis Nielson and Flemming Nielson. 2007. *Semantics with Applications: An Appetizer (Undergraduate Topics in*
1294 *Computer Science)*. Springer-Verlag, Berlin, Heidelberg. doi:10.1007/978-1-84628-692-6
- 1295 [30] Lois Orosa, A. Giray Yağlıkcı, Haocong Luo, Ataberk Olgun, Jisung Park, Hasan Hassan, Minesh Patel, Jeremie S.
1296 Kim, and Onur Mutlu. 2021. A Deeper Look into RowHammer’s Sensitivities: Experimental Analysis of Real DRAM
1297 Chips and Implications on Future Attacks and Defenses. In *54th Annual IEEE/ACM International Symposium on*
1298 *Microarchitecture (MICRO-54)*. doi:10.1145/3466752.3480069
- 1299 [31] Thibaut Pérami, Thomas Bauereiss, Brian Campbell, Zongyuan Liu, Nils Laueremann, Alasdair Armstrong, and Peter
1300 Sewell. 2026. ArchSem: Reusable Rigorous Semantics of Relaxed Architectures. *Proc. ACM Program. Lang.* 10, POPL
1301 (2026), 204–234. doi:10.1145/3776650
- 1302 [32] Gordon Plotkin and John Power. 2002. Notions of Computation Determine Monads. In *Foundations of Software Science and*
1303 *Computation Structures (Lecture Notes in Computer Science, Vol. 2303)*. Springer, 373–393. doi:10.1007/3-540-45931-
1304 [6_24](https://doi.org/10.1007/3-540-45931-6_24)
- 1305 [33] Norman Ramsey and Avi Pfeffer. 2002. Stochastic lambda calculus and monads of probability distributions. In
1306 *Proceedings of the 29th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*. ACM,
1307 154–165. doi:10.1145/503272.503288
- 1308 [34] Andrei Sabelfeld and Andrew C. Myers. 2003. Language-Based Information-Flow Security. *IEEE Journal on Selected*
1309 *Areas in Communications* 21, 1 (2003), 5–19. doi:10.1109/JSAC.2002.806121
- 1310 [35] Mark Seaborn and Thomas Dullien. 2015. Exploiting the DRAM rowhammer bug to gain kernel privileges. <https://googleprojectzero.blogspot.com/2015/03/exploiting-dram-rowhammer-bug-to-gain.html> Code at <https://github.com/google/rowhammer-test>.
- 1311 [36] Claude E. Shannon. 1948. A Mathematical Theory of Communication. *The Bell System Technical Journal* 27, 3 (1948),
1312 379–423. doi:10.1002/j.1538-7305.1948.tb01338.x
- 1313 [37] The mathlib Community. 2020. The Lean Mathematical Library. In *Proceedings of the 9th ACM SIGPLAN International*
1314 *Conference on Certified Programs and Proofs* (New Orleans, LA, USA) (CPP 2020). Association for Computing Machinery,
1315 New York, NY, USA, 367–381. doi:10.1145/3372885.3373824
- 1316 [38] Cláudio Vasconcelos and António Ravara. 2016. The While language. arXiv:1603.08949 [cs.PL] doi:10.48550/arXiv.
1317 [1603.08949](https://arxiv.org/abs/1603.08949)
- 1318 [39] Andrew Walker, Sungkwon Lee, and Dafna Beery. 2021. On DRAM Rowhammer and the Physics of Insecurity. *IEEE*
1319 *Transactions on Electron Devices* PP (03 2021), 1–11. doi:10.1109/TED.2021.3060362
- 1320 [40] A. Giray Yağlıkcı, Minesh Patel, Jeremie S. Kim, Roknoddin Azizi, Ataberk Olgun, Lois Orosa, Hasan Hassan, Jisung
1321 Park, Konstantinos Kanellopoulos, Taha Shahroodi, Saugata Ghose, and Onur Mutlu. 2021. BlockHammer: Preventing
1322 RowHammer at Low Cost by Blacklisting Rapidly-Accessed DRAM Rows. In *IEEE International Symposium on High-*
1323 *Performance Computer Architecture (HPCA)*. doi:10.1109/HPCA51647.2021.00037

Supplementary Appendix: Mechanised high-level operational semantics of Rowhammer

A Material for Section 4: A deterministic small-step operational semantics for a WHILE language

This appendix presents the full small-step operational semantics that were given only partially in Section 4. For more on small-step operational semantics see any textbook on programming language semantics like e.g., [4–6].

Recall that the syntax (arithmetic expressions, Boolean expressions, programs) are given by

$$e ::= c \mid x \mid e + e$$

$$b ::= \text{true} \mid \text{false} \mid \neg b \mid b \wedge b \mid e = e$$

$$P ::= \text{skip} \mid x := e \mid P; P \mid \text{if } b \text{ then } P \text{ else } P \mid \text{while } b \text{ do } P$$

We identify arithmetic constants with their denoted values. Thus $v, v_1, v_2 \in \text{Val}$ range over arithmetic values, represented by constants in the expression syntax. Boolean values are $t, t_1, t_2 \in \{\text{true}, \text{false}\}$.

Evaluation is left-to-right. Conjunction is deliberately non-short-circuiting, matching the expression semantics used elsewhere in the paper. Although deterministic expression reduction does not modify memory, memories are included in expression configurations so that the ordinary and faulty transition systems have the same shape.

Transitions carry a *label* recording the memory access they perform:

$$\ell ::= \text{read}(x) \mid \text{write}(x) \mid \tau.$$

A read of variable x is labelled $\text{read}(x)$, the write performed by an assignment to x is labelled $\text{write}(x)$, and every other transition is silent (τ); congruence rules propagate the label of their premise. The labels play no role in the deterministic semantics itself, but they localise the fault points of the probabilistic Rowhammer semantics of Section 5 and make access counting definitional. Single steps are always written with their label; every multi-step notion below is indexed by the *sequence* of labels performed (written with a double arrow), so the labelled judgement is the only transition relation in play.

A.0.1 Arithmetic expressions. An arithmetic-expression configuration has the form $\langle e, m \rangle$. The labelled relation

$$\langle e, m \rangle \xrightarrow{\ell}_a \langle e', m' \rangle$$

is the least relation generated by the following rules. Arithmetic values have no outgoing transition.

$$\frac{}{\langle x, m \rangle \xrightarrow{\text{read}(x)}_a \langle m(x), m \rangle} \quad (\text{A-READ})$$

$$\frac{\langle e_1, m \rangle \xrightarrow{\ell}_a \langle e'_1, m' \rangle}{\langle e_1 + e_2, m \rangle \xrightarrow{\ell}_a \langle e'_1 + e_2, m' \rangle} \quad (\text{A-PLUS-L})$$

$$\frac{\langle e_2, m \rangle \xrightarrow{\ell}_a \langle e'_2, m' \rangle}{\langle v_1 + e_2, m \rangle \xrightarrow{\ell}_a \langle v_1 + e'_2, m' \rangle} \quad (\text{A-PLUS-R})$$

$$\frac{v = v_1 + v_2}{\langle v_1 + v_2, m \rangle \xrightarrow{\tau}_a \langle v, m \rangle} \quad (\text{A-PLUS})$$

1373 *A.0.2 Boolean expressions.* A Boolean-expression configuration has the form $\langle b, m \rangle$. The labelled
1374 relation

$$\langle b, m \rangle \xrightarrow{b} \langle b', m' \rangle$$

1376 is generated by the following rules. The Boolean values true and false have no outgoing transition.
1377

$$\frac{\langle b, m \rangle \xrightarrow{b} \langle b', m' \rangle}{\langle \neg b, m \rangle \xrightarrow{b} \langle \neg b', m' \rangle} \quad (\text{B-NOT-STEP})$$

$$\frac{}{\langle \neg \text{true}, m \rangle \xrightarrow{\tau} \langle \text{false}, m \rangle} \quad (\text{B-NOT-TRUE})$$

$$\frac{}{\langle \neg \text{false}, m \rangle \xrightarrow{\tau} \langle \text{true}, m \rangle} \quad (\text{B-NOT-FALSE})$$

$$\frac{\langle b_1, m \rangle \xrightarrow{b} \langle b'_1, m' \rangle}{\langle b_1 \wedge b_2, m \rangle \xrightarrow{b} \langle b'_1 \wedge b_2, m' \rangle} \quad (\text{B-AND-L})$$

$$\frac{\langle b_2, m \rangle \xrightarrow{b} \langle b'_2, m' \rangle}{\langle t_1 \wedge b_2, m \rangle \xrightarrow{b} \langle t_1 \wedge b'_2, m' \rangle} \quad (\text{B-AND-R})$$

$$\frac{t = t_1 \wedge t_2}{\langle t_1 \wedge t_2, m \rangle \xrightarrow{\tau} \langle t, m \rangle} \quad (\text{B-AND})$$

$$\frac{\langle e_1, m \rangle \xrightarrow{a} \langle e'_1, m' \rangle}{\langle e_1 = e_2, m \rangle \xrightarrow{b} \langle e'_1 = e_2, m' \rangle} \quad (\text{B-EQ-L})$$

$$\frac{\langle e_2, m \rangle \xrightarrow{a} \langle e'_2, m' \rangle}{\langle v_1 = e_2, m \rangle \xrightarrow{b} \langle v_1 = e'_2, m' \rangle} \quad (\text{B-EQ-R})$$

$$\frac{t = (v_1 = v_2)}{\langle v_1 = v_2, m \rangle \xrightarrow{\tau} \langle t, m \rangle} \quad (\text{B-EQ})$$

1406 *A.0.3 Programs.* A program configuration has the form $\langle P, m \rangle$. The labelled program transition
1407 relation

$$\langle P, m \rangle \xrightarrow{p} \langle P', m' \rangle$$

1410 is the least relation generated by the following rules.

1411 An assignment first reduces its right-hand side one expression step at a time. Once the expression
1412 is a value, the assignment performs the write.

$$\frac{\langle e, m \rangle \xrightarrow{a} \langle e', m' \rangle}{\langle x := e, m \rangle \xrightarrow{p} \langle x := e', m' \rangle} \quad (\text{P-ASSIGN-STEP})$$

$$\frac{}{\langle x := v, m \rangle \xrightarrow{\text{write}(x)} \langle \text{skip}, m[x \mapsto v] \rangle} \quad (\text{P-ASSIGN})$$

1420 Sequential composition reduces its left-hand program until that program has terminated.
1421

$$\frac{\langle P, m \rangle \xrightarrow{p} \langle P', m' \rangle}{\langle P; Q, m \rangle \xrightarrow{p} \langle P'; Q, m' \rangle} \quad (\text{P-SEQ-STEP})$$

$$\frac{}{\langle \text{skip}; Q, m \rangle \xrightarrow{p} \langle Q, m \rangle} \quad (\text{P-SEQ-SKIP})$$

A conditional reduces its guard one Boolean-expression step at a time and then selects the corresponding branch.

$$\frac{\langle b, m \rangle \xrightarrow{b} \langle b', m' \rangle}{\langle \text{if } b \text{ then } P \text{ else } Q, m \rangle \xrightarrow{p} \langle \text{if } b' \text{ then } P \text{ else } Q, m' \rangle} \quad (\text{P-IF-STEP})$$

$$\frac{}{\langle \text{if true then } P \text{ else } Q, m \rangle \xrightarrow{p} \langle P, m \rangle} \quad (\text{P-IF-TRUE})$$

$$\frac{}{\langle \text{if false then } P \text{ else } Q, m \rangle \xrightarrow{p} \langle Q, m \rangle} \quad (\text{P-IF-FALSE})$$

A while loop unfolds by one standard small step. Its guard is subsequently reduced by the conditional rules above.

$$\frac{}{\langle \text{while } b \text{ do } P, m \rangle \xrightarrow{p} \langle \text{if } b \text{ then } (P; \text{while } b \text{ do } P) \text{ else skip}, m \rangle} \quad (\text{P-WHILE})$$

There is no transition from $\langle \text{skip}, m \rangle$; these are the terminal program configurations.

A.0.4 Finite and terminating executions. A finite execution is indexed by the sequence of labels it performs. Write

$$\langle P, m \rangle \xRightarrow[p]{s} \langle P', m' \rangle$$

for $s = \ell_1 \cdots \ell_k$ when there are configurations

$$\langle P, m \rangle = C_0 \xrightarrow{p} \ell_1 C_1 \xrightarrow{p} \ell_2 \cdots \xrightarrow{p} \ell_k C_k = \langle P', m' \rangle.$$

The length $|s|$ is the number of transitions performed. Because expression reductions are propagated through assignments and conditionals one step at a time, it counts both expression reductions and command reductions. The *access trace* $\text{acc}(s)$ is the subsequence of non- τ labels of s , in order; its length is the number of memory accesses (reads and writes) the run performs. Both are definitions over the label sequence, not counts of rule occurrences inside derivation trees. The analogous trace-indexed relations $\xRightarrow[a]{s}$ and $\xRightarrow[b]{s}$ at the expression levels are defined in the same way.

A program terminates with final memory m' , written

$$\langle P, m \rangle \Downarrow m',$$

when $\langle P, m \rangle \xRightarrow[p]{s} \langle \text{skip}, m' \rangle$ for some label sequence s . It terminates within at most n small steps, written

$$\langle P, m \rangle \Downarrow_{\leq n} m',$$

when moreover some such s has $|s| \leq n$. A program diverges from m , written

$$\langle P, m \rangle \Uparrow,$$

1471 when for every $n \in \mathbb{N}$ there are a label sequence s with $|s| = n$ and a configuration $\langle P_n, m_n \rangle$ such
 1472 that $\langle P, m \rangle \xRightarrow{s}_p \langle P_n, m_n \rangle$. The index n here is a literal bound on the number of small-step transitions.
 1473 Sequential composition cannot duplicate or reuse a budget already consumed by its left-hand
 1474 program.
 1475

1476 **A.0.5 Basic metatheory.** The lemmas of this subsection are standard and proved on paper. The
 1477 mechanisation contains those the soundness theorem rests on (determinism and progress in func-
 1478 tional form, trace realisation, divergence, and well-formedness preservation). The remaining lemmas
 1479 below are stated for completeness of the presentation and are not part of the Lean development.
 1480

1481 **LEMMA A.1 (VALUES ARE TERMINAL).** *There are no ℓ and C with $\langle v, m \rangle \xrightarrow{\ell}_a C$; none with $\langle t, m \rangle \xrightarrow{\ell}_b$
 1482 C ; and none with $\langle \text{skip}, m \rangle \xrightarrow{\ell}_p C$.*

1483 **LEMMA A.2 (DETERMINISM).** *Each of the labelled relations is a partial function on configurations,
 1484 including its label: if $C \xrightarrow{\ell_1}_p C_1$ and $C \xrightarrow{\ell_2}_p C_2$ then $\ell_1 = \ell_2$ and $C_1 = C_2$, and likewise for $\xrightarrow{\cdot}_a$ and $\xrightarrow{\cdot}_b$.*

1486 **PROOF.** Rule induction. For each syntactic form at most one rule applies: the congruence rules
 1487 require a non-value in the position being reduced (values are terminal by [Theorem A.1](#), so their
 1488 premises are otherwise unsatisfiable), and the redex rules require values in all positions. \square
 1489

1490 **LEMMA A.3 (PROGRESS).** *If e is not a value then $\langle e, m \rangle \xrightarrow{\ell}_a \langle e', m' \rangle$ for some ℓ, e', m' ; if b is not a
 1491 value then $\langle b, m \rangle \xrightarrow{\ell}_b \langle b', m' \rangle$ for some ℓ, b', m' ; and if $P \neq \text{skip}$ then $\langle P, m \rangle \xrightarrow{\ell}_p \langle P', m' \rangle$ for some $\ell,$
 1492 P', m' .*

1493 **PROOF.** Induction on the syntax, using totality of memories for A-READ. \square
 1494

1495 **COROLLARY A.4 (TERMINATION-DIVERGENCE DICHOTOMY).** *For every configuration $\langle P, m \rangle$, exactly
 1496 one of the following holds: $\langle P, m \rangle \Downarrow m'$ for a unique m' , or $\langle P, m \rangle \Uparrow$.*

1497 **PROOF.** By [Theorem A.2](#) the maximal transition sequence from $\langle P, m \rangle$ is unique, and by [Theo-](#)
 1498 [rem A.3](#) it can stop only at a terminal configuration $\langle \text{skip}, m' \rangle$. If that sequence is finite it yields
 1499 $\Downarrow m'$, with m' unique by determinism iterated along the run; if it is infinite, every n is realised and
 1500 \Uparrow holds. The two cases exclude each other because a terminated run has no continuation. \square
 1501

1502 **LEMMA A.5 (MEMORY INVARIANCE OF DETERMINISTIC EXPRESSION STEPS).** *If $\langle e, m \rangle \xrightarrow{\ell}_a \langle e', m' \rangle$ or
 1503 $\langle b, m \rangle \xrightarrow{\ell}_b \langle b', m' \rangle$, then $m' = m$.*

1504 Deterministically, expression-level small steps are thus redundant: only P-ASSIGN modifies the
 1505 memory. This lemma is the precise sense in which the fine-grained presentation is overkill for the
 1506 ordinary semantics; the Rowhammer refinement of [Section B](#) breaks exactly this lemma, a read
 1507 may disturb the blast radius of the accessed location, and nothing else in the design.
 1508

1509 **LEMMA A.6 (EXPRESSION EVALUATION TERMINATES, WITH EXPLICIT COST).** *Let $r(e)$ be the number
 1510 of variable occurrences plus the number of operators in e . Then $\langle e, m \rangle \xRightarrow{s}_a \langle v, m \rangle$ for a unique value v
 1511 and some s with $|s| = r(e)$, and similarly for Boolean expressions. Each transition fires exactly one
 1512 A-READ or one redex rule under congruences, so r strictly decreases; $\text{acc}(s)$ is the sequence of reads of
 1513 e , in left-to-right order.
 1514*

1515 **LEMMA A.7 (ADEQUACY OF EXPLICIT READS).** *$\langle e, m \rangle \xRightarrow{s}_a \langle v, m \rangle$ for some s iff $\llbracket e \rrbracket^{\text{ord}}(m) = v$, and
 1516 $\langle b, m \rangle \xRightarrow{s}_b \langle t, m \rangle$ for some s iff $\llbracket b \rrbracket^{\text{ord}}(m) = t$, where $\llbracket \cdot \rrbracket^{\text{ord}}$ is the atomic expression evaluation defined
 1517 above.
 1518
 1519*

Each dynamic variable read is exactly one $\text{read}(x)$ -labelled transition and each dynamic assignment write exactly one $\text{write}(x)$ -labelled transition, so the number of potential access-triggered Rowhammer events in a run is the length of its access trace. The total transition count is an exact cost for this abstract machine, but it is *not* the number of memory accesses: it also counts operator reductions, branch selection, sequencing elimination, and loop unfolding, silent (τ) steps that trigger no fault in the semantics of Section 5.

B Material for Section 5: A probabilistic model of Rowhammer

This appendix presents the full small-step probabilistic operational semantics that was given only partially in Section 5. For more on probabilistic operational semantics see e.g., [1–3].

The syntax and evaluation order are those of Section A. Throughout, the blast radius br and the protected set $locs$ are ambient parameters, and pflip is a fault kernel subject to the contract (F1–F4) of Section 5.1. The judgements

$$\langle e, m \rangle \rightsquigarrow_a^{\text{RH}} \mu, \quad \langle b, m \rangle \rightsquigarrow_b^{\text{RH}} \mu, \quad \langle P, m \rangle \rightsquigarrow_p^{\text{RH}} \mu$$

relate a configuration to a *distribution* over successor configurations, with $\mu \in \text{Dist}(\text{AExp} \times \text{Memory})$, $\mu \in \text{Dist}(\text{BExp} \times \text{Memory})$, and $\mu \in \text{Dist}(\text{Prog} \times \text{Memory})$ respectively. Exactly two rules sample the kernel, RH-A-READ and RH-P-ASSIGN: every redex rule is the Dirac lifting of its deterministic counterpart, and every congruence rule pushes the reduction context through the premise distribution with Dist.map .

B.0.1 Arithmetic expressions. Arithmetic values have no outgoing transition. A variable read first obtains the value $m(x)$ and then samples the fault triggered by that access.

$$\frac{}{\langle x, m \rangle \rightsquigarrow_a^{\text{RH}} (\text{let } m' \leftarrow \text{pflip}(\text{adj}(br, locs, x), m) \text{ in } \delta(\langle m(x), m' \rangle))} \quad \text{(RH-A-READ)}$$

$$\frac{\langle e_1, m \rangle \rightsquigarrow_a^{\text{RH}} \mu}{\langle e_1 + e_2, m \rangle \rightsquigarrow_a^{\text{RH}} \text{Dist.map}(\lambda\langle e'_1, m' \rangle. \langle e'_1 + e_2, m' \rangle)(\mu)} \quad \text{(RH-A-PLUS-L)}$$

$$\frac{\langle e_2, m \rangle \rightsquigarrow_a^{\text{RH}} \mu}{\langle v_1 + e_2, m \rangle \rightsquigarrow_a^{\text{RH}} \text{Dist.map}(\lambda\langle e'_2, m' \rangle. \langle v_1 + e'_2, m' \rangle)(\mu)} \quad \text{(RH-A-PLUS-R)}$$

$$\frac{v = v_1 + v_2}{\langle v_1 + v_2, m \rangle \rightsquigarrow_a^{\text{RH}} \delta(\langle v, m \rangle)} \quad \text{(RH-A-PLUS)}$$

B.0.2 Boolean expressions. The Boolean values true and false have no outgoing transition.

$$\frac{\langle b, m \rangle \rightsquigarrow_b^{\text{RH}} \mu}{\langle \neg b, m \rangle \rightsquigarrow_b^{\text{RH}} \text{Dist.map}(\lambda\langle b', m' \rangle. \langle \neg b', m' \rangle)(\mu)} \quad \text{(RH-B-NOT-STEP)}$$

$$\frac{}{\langle \neg \text{true}, m \rangle \rightsquigarrow_b^{\text{RH}} \delta(\langle \text{false}, m \rangle)} \quad \text{(RH-B-NOT-TRUE)}$$

$$\frac{}{\langle \neg \text{false}, m \rangle \rightsquigarrow_b^{\text{RH}} \delta(\langle \text{true}, m \rangle)} \quad \text{(RH-B-NOT-FALSE)}$$

$$\frac{\langle b_1, m \rangle \rightsquigarrow_b^{\text{RH}} \mu}{\langle b_1 \wedge b_2, m \rangle \rightsquigarrow_b^{\text{RH}} \text{Dist.map}(\lambda\langle b'_1, m' \rangle. \langle b'_1 \wedge b_2, m' \rangle)(\mu)} \quad \text{(RH-B-AND-L)}$$

$$\frac{\langle b_2, m \rangle \rightsquigarrow_b^{\text{RH}} \mu}{\langle b_1 \wedge b_2, m \rangle \rightsquigarrow_b^{\text{RH}} \text{Dist.map}(\lambda\langle b'_2, m' \rangle. \langle b_1 \wedge b'_2, m' \rangle)(\mu)} \quad \text{(RH-B-AND-R)}$$

$$\frac{t = t_1 \wedge t_2}{\langle t_1 \wedge t_2, m \rangle \rightsquigarrow_b^{\text{RH}} \delta(\langle t, m \rangle)} \quad (\text{RH-B-AND})$$

$$\frac{\langle e_1, m \rangle \rightsquigarrow_a^{\text{RH}} \mu}{\langle e_1 = e_2, m \rangle \rightsquigarrow_b^{\text{RH}} \text{Dist.map}(\lambda \langle e'_1, m' \rangle. \langle e'_1 = e_2, m' \rangle)(\mu)} \quad (\text{RH-B-EQ-L})$$

$$\frac{\langle e_2, m \rangle \rightsquigarrow_a^{\text{RH}} \mu}{\langle v_1 = e_2, m \rangle \rightsquigarrow_b^{\text{RH}} \text{Dist.map}(\lambda \langle e'_2, m' \rangle. \langle v_1 = e'_2, m' \rangle)(\mu)} \quad (\text{RH-B-EQ-R})$$

$$\frac{t = (v_1 = v_2)}{\langle v_1 = v_2, m \rangle \rightsquigarrow_b^{\text{RH}} \delta(\langle t, m \rangle)} \quad (\text{RH-B-EQ})$$

1580 **B.0.3 Programs.** An assignment first reduces its right-hand side one expression step at a time.
 1581 Once the expression is a value, the assignment writes that value and samples the fault triggered by
 1582 the write, on the *updated* memory.
 1583

$$\frac{\langle e, m \rangle \rightsquigarrow_a^{\text{RH}} \mu}{\langle x := e, m \rangle \rightsquigarrow_p^{\text{RH}} \text{Dist.map}(\lambda \langle e', m' \rangle. \langle x := e', m' \rangle)(\mu)} \quad (\text{RH-P-ASSIGN-STEP})$$

$$\frac{}{\langle x := v, m \rangle \rightsquigarrow_p^{\text{RH}} (\text{let } m' \leftarrow \text{pflip}(\text{adj}(br, locs, x), m[x \mapsto v]) \text{ in } \delta(\langle \text{skip}, m' \rangle))} \quad (\text{RH-P-ASSIGN})$$

1588 Sequential composition reduces its left-hand program until that program has terminated.
 1589

$$\frac{\langle P, m \rangle \rightsquigarrow_p^{\text{RH}} \mu}{\langle P; Q, m \rangle \rightsquigarrow_p^{\text{RH}} \text{Dist.map}(\lambda \langle P', m' \rangle. \langle P'; Q, m' \rangle)(\mu)} \quad (\text{RH-P-SEQ-STEP})$$

$$\frac{}{\langle \text{skip}; Q, m \rangle \rightsquigarrow_p^{\text{RH}} \delta(\langle Q, m \rangle)} \quad (\text{RH-P-SEQ-SKIP})$$

1596 A conditional reduces its guard one Boolean-expression step at a time and then selects the
 1597 corresponding branch.
 1598

$$\frac{\langle b, m \rangle \rightsquigarrow_b^{\text{RH}} \mu}{\langle \text{if } b \text{ then } P \text{ else } Q, m \rangle \rightsquigarrow_p^{\text{RH}} \text{Dist.map}(\lambda \langle b', m' \rangle. \langle \text{if } b' \text{ then } P \text{ else } Q, m' \rangle)(\mu)} \quad (\text{RH-P-IF-STEP})$$

$$\frac{}{\langle \text{if true then } P \text{ else } Q, m \rangle \rightsquigarrow_p^{\text{RH}} \delta(\langle P, m \rangle)} \quad (\text{RH-P-IF-TRUE})$$

$$\frac{}{\langle \text{if false then } P \text{ else } Q, m \rangle \rightsquigarrow_p^{\text{RH}} \delta(\langle Q, m \rangle)} \quad (\text{RH-P-IF-FALSE})$$

1609 A while loop unfolds by one standard small step. Its guard is subsequently reduced by the
 1610 conditional rules above.
 1611

$$\frac{}{\langle \text{while } b \text{ do } P, m \rangle \rightsquigarrow_p^{\text{RH}} \delta(\langle \text{if } b \text{ then } (P; \text{while } b \text{ do } P) \text{ else skip}, m \rangle)} \quad (\text{RH-P-WHILE})$$

1612 There is no transition from $\langle \text{skip}, m \rangle$. These are the terminal program configurations. The rules
 1613 are syntax-directed, so every non-terminal configuration determines a unique one-step distribution,
 1614
 1615
 1616
 1617

1618 written $\text{step}_{br,locs}^{\text{RH}}(\langle P, m \rangle)$ at the program level. The absorbing kernel $\widehat{\text{step}}_{br,locs}^{\text{RH}}$ and the finite runs
 1619 $\text{run}_{br,locs}^n$ and $\text{run}_{br,locs,tr}^n$ built from it are defined in Section 5.5.

1620

1621 B.1 The mechanised proof of the central theorem

1622 For reference we reproduce the Lean proof of Theorem 6.1 from `Soundness.lean`. The prose sketch
 1623 is in Section 6.4.

1624

```

1625 theorem physical_separation_sound : PhysicalSeparationSound := by
1626   intro sep f hframe br locs hsafe p hwf m n
1627   apply Dist.mapD_eq_dirac_of_support
1628   intro r hr
1629   have hsub :  $\forall x \in \text{locs}, x \in \text{locs} := \text{fun } x \text{ hx} \Rightarrow \text{hx}$ 
1630   have hpres :  $\forall x \in \text{locs}, \forall (m_0 m' : \text{Memory}),$ 
1631      $m' \in (f (\text{sep.adj } br \text{ locs } x) m_0). \text{support} \rightarrow \forall \ell, \ell \in \text{locs} \rightarrow m' \ell = m_0 \ell :=$ 
1632      $\text{fun } x \text{ hx } m_0 m' \text{ hm}' \ell \text{ h}\ell \Rightarrow$ 
1633      $\text{hframe.outside\_unchanged } \_ m_0 m' \ell \text{ hm}' (\text{sep.safe\_disjoint } \text{hsafe } \text{hx } \ell \text{ h}\ell)$ 
1634   obtain  $\langle h1, h2, h3 \rangle :=$ 
1635      $\text{runRHTrace\_protected } (\text{keep} := (\cdot \in \text{locs})) \text{ hsub hpres } n \text{ hwf}$ 
1636      $(\text{AgreeP.refl } \_ m) r \text{ hr}$ 
1637   simp only [protView, Prod.mk.injEq]
1638   refine  $\langle h1, ?\_ , ?\_ \rangle$ 
1639   · funext x
1640     exact h2 x.1 x.2
1641   · rw [h3]
```

1641 The script follows the three moves of the sketch. The `intro` line introduces the separation geometry,
 1642 the kernel with its frame-locality hypothesis (`hframe`, clause F2), the blast radius and protected set,
 1643 physical separation (`hsafe`), the well-formed program (`hwf`), the initial memory, and the step bound.
 1644 `Dist.mapD_eq_dirac_of_support` is the support characterisation of the point mass (Section 3): it
 1645 reduces the distributional equation to showing that every outcome `r` in the support of the trace-
 1646 carrying Rowhammer run has the protected view of the deterministic run. The two `haves` discharge
 1647 the two clauses of the parametric invariant at `keep` $\triangleq (\cdot \in \text{locs})$: `hsub` is clause (i), used locations
 1648 are kept, trivial at this instantiation, where it is exactly well-formedness, `hpres` is clause (ii), fault
 1649 samples preserve kept locations, composed from `safe_disjoint` (the victim set of a protected
 1650 access avoids `locs`) and `hframe.outside_unchanged` (samples fix everything outside the victim
 1651 set). The iterated invariant `runRHTrace_protected`, started from the reflexive agreement on the
 1652 initial memory, then yields the three components of the protected view: `h1`, equal residual programs,
 1653 `h2`, memories agreeing on `locs`, `h3`, equal label sequences. The remainder repackages them as an
 1654 equality of protected views: `funext` turns pointwise agreement on protected locations into equality
 1655 of the restricted memories as functions on the subtype, and `h3` rewrites the access traces.

1656

1657 Appendix References

- 1658 [1] Fredrik Dahlqvist, Alexandra Silva, and Dexter Kozen. 2020. *Semantics of Probabilistic Programming: A Gentle Introduction*.
 1659 Cambridge University Press, 1–42. doi:10.1017/9781108770750.002
- 1660 [2] Alessandra Di Pierro, Chris Hankin, and Herbert Wiklicky. 2010. Probabilistic Semantics and Program Analysis. In
 1661 *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 1–42. doi:10.1007/978-3-642-13678-8_1
- 1662 [3] Ugo Dal Lago and Margherita Zorzi. 2012. Probabilistic operational semantics for the lambda calculus. *RAIRO -*
 1663 *Theoretical Informatics and Applications* 46 (2012), 413–450. doi:10.1051/ita/2012012
- 1664 [4] Hanne Riis Nielson and Flemming Nielson. 2007. *Semantics with Applications: An Appetizer (Undergraduate Topics in*
 1665 *Computer Science)*. Springer-Verlag, Berlin, Heidelberg. doi:10.1007/978-1-84628-692-6
- 1666 [5] Tobias Nipkow and Gerwin Klein. 2014. *Concrete Semantics: With Isabelle/HOL*. Springer. doi:10.1007/978-3-319-10542-0

- 1667 [6] Glynn Winskel. 1993. *The Formal Semantics of Programming Languages*. The MIT Press. doi:10.7551/mitpress/3054.001.
1668 0001
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715