



This is a repository copy of *Relative security: formally modeling and (dis)proving resilience against semantic optimization vulnerabilities*.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/231110/>

Version: Accepted Version

---

### Proceedings Paper:

Dongol, B., Griffin, M., Popescu, A. [orcid.org/0000-0001-8747-0619](https://orcid.org/0000-0001-8747-0619) et al. (1 more author) (2024) Relative security: formally modeling and (dis)proving resilience against semantic optimization vulnerabilities. In: 2024 IEEE 37th Computer Security Foundations Symposium (CSF). 2024 IEEE 37th Computer Security Foundations Symposium (CSF), 08-12 Jul 2024, Enschede, Netherlands. Institute of Electrical and Electronics Engineers (IEEE), pp. 403-418. ISBN: 9798350362046. ISSN: 1940-1434. EISSN: 2374-8303.

<https://doi.org/10.1109/csf61375.2024.00027>

---

© 2024 The Authors. Except as otherwise noted, this author-accepted version of a paper published in 2024 IEEE 37th Computer Security Foundations Symposium (CSF) is made available via the University of Sheffield Research Publications and Copyright Policy under the terms of the Creative Commons Attribution 4.0 International License (CC-BY 4.0), which permits unrestricted use, distribution and reproduction in any medium, provided the original work is properly cited. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>

### Reuse

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

### Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.



[eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk)  
<https://eprints.whiterose.ac.uk/>

# Relative Security: Formally Modeling and (Dis)Proving Resilience Against Semantic Optimization Vulnerabilities

Brijesh Dongol\*, Matt Griffin\*, Andrei Popescu† and Jamie Wright†

\*Department of Computer Science, University of Surrey, UK Email: {b.dongol, matt.griffin}@surrey.ac.uk

†Department of Computer Science, University of Sheffield, UK Email: {a.popescu, jwright8}@sheffield.ac.uk

**Abstract**—Meltdown and Spectre are vulnerabilities known as transient execution vulnerabilities, where an attacker exploits speculative execution (a semantic optimization present in most modern processors) to break confidentiality. We introduce *relative security*, a general notion of information-flow security that models this type of vulnerability by contrasting the leaks that are possible in a “vanilla” semantics with those possible in a different semantics, often obtained from the vanilla semantics via some optimizations. We describe incremental proof methods, in the style of Goguen and Meseguer’s unwinding, both for proving and for disproving relative security, and deploy these to formally establish the relative (in)security of some standard Spectre examples. Both the abstract results and the case studies have been mechanized in the Isabelle/HOL theorem prover.

## 1. Introduction

Meltdown [20] and Spectre [19] are transient execution vulnerabilities, which exploit timing-based side-channels caused by speculative execution optimizations, as present in almost all modern processors. Mitigating against such vulnerabilities is of significant importance for computer security. This paper proposes a new model-theoretic and proof-theoretic framework for describing and verifying resilience against transient execution attacks and beyond. We introduce *relative security* (§2), a general notion that focuses not on the absolute (information-flow) security of a system, but on the difference in security between a basic, “vanilla” system and a system that is enhanced to allow optimizations in the system’s execution. More precisely, relative security expresses that there is *no* difference in information flow: Any leak occurring in the vanilla system can also occur in the optimization-enhanced one. We instantiate this notion to a programming language with speculative semantics (§3), yielding intuitive results on some standard example programs. While building on a rich literature that addresses this problem from a formal modeling perspective (§4), relative security’s innovation is in natively capturing *fully interactive attackers* and *dynamic creation of secrets*, as required, e.g., by operating system processes. There are two key features enabling this: **1)** a fine-grained attacker model that allows secrets, attacker actions and attacker observations anywhere on the execution trace, and **2)** a view of *leaks as first-class citizens* that takes advantage of this fine granularity.

Our second contribution is general methods for *proving* and *disproving* relative security in an incremental fashion

(§5), generalizing the unwinding method by Goguen and Meseguer [12]. On the proof front (§5.1–5.4), our method allows the incremental construction, from any two optimization-enhanced execution traces that exhibit a leak, of two counterpart vanilla traces that exhibit the same leak. The four traces are constrained by both “secrecy contracts” and “interaction contracts”, which postulate local similarities and disimilarities between pairs of them.

Due to the complex dynamics between the four involved traces, the *disproof* front of relative security (§5.5) is also interesting and benefits from the idea of unwinding: A counterexample requires indicating a concrete leak in the optimization-enhanced system, followed by a proof employing a form of *secret-directed unwinding* that shows how this leak cannot be reproduced in the vanilla system. Both the proof and disproof methods have been validated on our example programs (§6).

The relative security framework, as well as the language semantics and examples, have been mechanized in the Isabelle/HOL prover [26] (§7). More details on our results and proofs can be found in an extended technical report [7].

## 2. Defining Relative Security

This section motivates and introduces our notion of relative security. We start with example programs exhibiting various features we would like to cover (§2.1), after which we develop the abstract framework: On top of system models capturing the program semantics (§2.2), we introduce *leakage models* which allow us to express the essence of relative security (§2.3). Then we instantiate these to more concrete (*state-wise*) *attacker models* (§2.4, §2.5), where leaks are defined from secrets and attacker actions and observations.

### 2.1. Motivating examples

Next we discuss example C programs, including some standard ones from the Spectre literature [5, 18]. In all examples, we assume that the attacker controls the inputs and sees the outputs (unless specified otherwise), and can also infer the locations accessed for reading (as in Spectre attacks). The program in Listing 1 is vulnerable to a Spectre bounds-check-bypass attack. If the value of  $x$  is greater than or equal to  $N$ , under normal execution it is impossible for an attacker to infer the *value* of  $a[x]$  (which is needed to load a value from  $b$  into  $t$ ). However, under the speculative execution

```

1 unsigned fun1(unsigned x) {
2     unsigned t = 0;
3     if (x < N) {
4         t = b[a[x] * 512]; }
5     return t; }

```

Listing 1: Spectre Bounds Check Bypass (BCB)

```

1 unsigned fun2(unsigned x) {
2     unsigned t = 0;
3     if (x < N) {
4         _mm_lfence();
5         t = b[a[x] * 512]; }
6     return t; }

```

Listing 2: Fix 1 for Spectre BCB

```

1 unsigned fun3(unsigned x) {
2     unsigned t = 0;
3     if (x < N) {
4         unsigned v = a[x];
5         _mm_lfence();
6         t = b[v * 512]; }
7     return t; }

```

Listing 3: Fix 2 for Spectre BCB

optimization, a branch misprediction enables this load to occur, potentially allowing the attacker to later perform a side-channel attack on the cache to infer the value of  $a[x]$ .

A suboptimal fix to this particular function is straightforward. As shown in Listing 2, one can introduce a load fence (`_mm_lfence`) instruction prior to loading from  $b$  and thus calculating the value of  $a[x]$ . This effectively resolves speculation by disallowing the processor from continuing execution until it can be certain that the branch will be taken, i.e., that  $x$  is indeed less than  $N$ .

A less obvious fix (given in Listing 3) is also possible. This program seemingly suffers from a transient execution vulnerability because speculative execution loads the value of  $a[x]$  into the cache, whereas non-speculative execution does not. However, since both  $x$  and the base address of  $a$  are known (or inferable) to the attacker, loading  $a[x]$  does not leak any additional information under misprediction.

Now consider the program in Listing 4, discussed by Cheang et al. [5]. The authors refer to it as being *conditionally secure* because its susceptibility to a transient execution attack depends on the value of  $N$ . For  $N = 0$ , the program is insecure because  $b[a[0] * 512]$  can be loaded into the cache (thus leaking  $a[0]$ ) only under misprediction.

Now suppose  $N = 2$  (i.e.,  $N$  is some value greater than 0), and consider a trace in which the attacker calls `fun4(3)`, triggering a misprediction and loading the value of  $a[0]$  into  $v$ . Should such an execution be ruled insecure? Unlike the program in Listing 1, here the value leaked from  $a$  is always from index 0. Thus, the leak through misprediction described above is also possible through a normal execution, e.g., where the attacker chooses  $x = 1$ . Thus, the program does *not* have a transient execution vulnerability.

These examples have certain features in common: First, there are *secrets* to be protected—some parts of the

```

1 unsigned fun4(unsigned x) {
2     unsigned t = 0;
3     if (x < N) {
4         unsigned v = a[0];
5         t = b[v * 512]; }
6     return t; }

```

Listing 4: Conditionally secure BCB

```

1 void fun5() {
2     unsigned t = 0;
3     unsigned x = 1;
4     while (x != 0) {
5         scanf("%u", &x);
6         if (x < N) {
7             unsigned v = a[x];
8             _mm_lfence();
9             t = b[v * 512];
10            printf("%u", t); } } }

```

Listing 5: Secure interactive program

memory (inside or outside the bounds of array  $a$ ). Second, there are *observations* that an attacker can make—via standard channels such as function return or side channels, e.g., the affected cache. Finally, there are *actions* that an (active) attacker can take to interact with, and influence the execution of the program—via inputs to the functions.

While all the examples so far are only *end-to-end interactive*, i.e., take an input at the beginning and return an output at the end, this does not need to be the case. For example, Listing 5 shows a (possibly nonterminating) fully interactive program, which keeps inputting an integer value (via the `scanf` function) and accessing the arrays  $a$  and  $b$  based on it until the inputted value is 0. It also outputs a (possibly infinite) stream of elements from the array  $b$  (via the `printf` function). So here, depending on the particular setting, one may wish to assume that the attacker/observer continuously interacts with the program via both observations (printed values) and actions (scanned inputs)—and the Listing 5 program is still Spectre-secure under these assumptions.

Thus, actions and observations can take place not only at the execution’s beginning and end respectively, but *fully interactively*, and possibly *in(de)finitely*. In some cases, this could also be true about secrets: One may wish to protect not only the initial memory, but also inputs that are under the control of trusted parties. Listing 6 shows a contrived example illustrating this: Now the program reads from both an untrusted and a trusted source, and one may wish to ensure that the input from the trusted source is not leaked. The call to a procedure named `writeOnSecretFile` illustrates the potential processing of some trusted input stored in  $y$  (here, together with that of untrusted input stored in  $x$ ), which could influence a file on the disk; in this example, such an influence is harmless, since that file is assumed unobservable by the attacker.

In conclusion, we want to model fully interactive (both acting and observing) attackers and fully interactive secret

```

1 void fun6() {
2   unsigned t = 0;
3   unsigned x = 1;
4   while (x != 0) {
5     x = getUntrustedInput();
6     unsigned y = getTrustedInput();
7     if (x < N) {
8       unsigned v = a[x];
9       writeOnSecretFile(x, y);
10      _mm_lfence();
11      t = b[v * 512];
12      printf("%u", t); } } }

```

Listing 6: Secure secret-interactive program

uploading, while also allowing interaction to take place in(de)finitely, i.e., factoring in infinite executions.

## 2.2. System models

For a set  $A$ , a *sequence* over  $A$  is an item in  $\text{Seq}(A) = A^* \cup A^\mathbb{N}$ , i.e., a finite or infinite list of elements from  $A$ .

A *system model* is a triple  $SM = (\text{State}, \text{istate}, \Rightarrow)$ , consisting of: **1**) a set  $\text{State}$  of *states*, ranged over by  $s$ , **2**) a predicate  $\text{istate} : \text{State} \rightarrow \text{Bool}$  that describes the *initial states*, and **3**) a binary relation  $\Rightarrow : \text{State} \times \text{State} \rightarrow \text{Bool}$  on states called the *transition relation*.

The next concepts are relative to a system model  $SM = (\text{State}, \text{istate}, \Rightarrow)$ . We say that a state  $s \in \text{State}$  is *final*, written  $\text{final}(s)$ , when there is no transition out of  $s$ , i.e.,  $\neg \exists s'. s \Rightarrow s'$ . An (*execution*) *trace* is a nonempty (finite or infinite) maximal sequence of states  $s_0 s_1 \dots \in \text{Seq}(\text{State})$  such that  $\text{istate}(s_0)$  holds and  $s_i \Rightarrow s_{i+1}$  for all  $i$ . Maximality refers to the suffix relation; it is equivalent to saying that, if the sequence is finite then its last state is final. We let  $\text{Trace} \subseteq \text{Seq}(\text{State})$ , ranged over by  $\pi, \rho$ , be the set of traces,  $\text{Trace}^{\text{fin}}$  be the subset of  $\text{Trace}$  consisting of the finite traces only.

## 2.3. Very abstract relative security: leakage models

We fix two system models  $SM_{\text{van}} = (\text{State}_{\text{van}}, \text{istate}_{\text{van}}, \Rightarrow_{\text{van}})$  and  $SM_{\text{opt}} = (\text{State}_{\text{opt}}, \text{istate}_{\text{opt}}, \Rightarrow_{\text{opt}})$  such as  $\text{State}_{\text{van}} \subseteq \text{State}_{\text{opt}}$ . We will call  $SM_{\text{van}}$  the *vanilla* system model, and refer to its states, transition relation etc. as vanilla states, vanilla transition etc. Its set of traces is denoted by  $\text{Trace}_{\text{van}}$ . The vanilla system model represents the plain system, featuring no optimization—which will act as reference for our (relative) notion of security. Moreover, we will call  $SM_{\text{opt}}$  the *optimization-enhanced* system model, since it will stand for the system that has been optimized in various ways, e.g., with speculative or out-of-order executions. Its set of traces is denoted by  $\text{Trace}_{\text{opt}}$ . We will use the following short names: “vtrace” for “vanilla trace” (i.e., an element of  $\text{Trace}_{\text{van}}$ ), and “otrace” for “optimization-enhanced trace” (i.e., an element of  $\text{Trace}_{\text{opt}}$ ).

We are interested in expressing the information-flow security of  $SM_{\text{opt}}$  relative to that of  $SM_{\text{van}}$ , in order to assess whether the optimizations have introduced further vulnerabilities (as is the case with Spectre). However, we

do not just want to check the implication “if  $SM_{\text{van}}$  does not leak then  $SM_{\text{opt}}$  does not leak either”, because that would be too coarse: We wish to allow  $SM_{\text{van}}$  to have some (presumably acceptable, or at least known) leaks, and to check that  $SM_{\text{opt}}$  has no *additional* leaks.

What we seem to need for this is the ability to talk *leak-wise*, i.e., to express not just the *absence* of any leak (as done by traditional notions such as noninterference), but explicitly the *very notion* of a leak. To this end, we introduce leakage models: Given any system model  $SM = (\text{State}, \text{istate}, \Rightarrow)$ , a *leakage model* for  $SM$  is a pair  $(\text{Leak}, \text{leakVia})$ , where  $\text{Leak}$ , ranged over by  $l$ , is a set of entities called *leaks*, and  $\text{leakVia}$  is a predicate in  $\text{Trace} \times \text{Trace} \times \text{Leak} \rightarrow \text{Bool}$ . We think of  $\text{leakVia}(\pi_1, \pi_2, l)$  as expressing that the traces  $\pi_1$  and  $\pi_2$  exhibit the leak  $l$ . Indeed, it is well known [6] that whatever the leaks are, one requires two traces (i.e., two alternative executions) rather than just one to exhibit a leak.

Note that we work very abstractly, gradually instantiating our concepts. For now, we leave the notion of a leak unspecified. In the next subsection, §2.4, we will get more concrete, taking leaks to be pairs of sequences of secrets produced by alternative execution traces—which, for suitable choices of the notion of secret, recovers what is typically taken to constitute a leak in a language-based setting [30]. In §3, when we have at our disposal some concrete system models given by a programming language semantics, we further instantiate the secrets to be the initial memories and inputs and outputs over certain trusted channels.

We now have all the ingredients for an abstract, leak-centric definition of relative security:

**Def. 1** Let  $\mathcal{LM}_{\text{van}} = (\text{Leak}, \text{leakVia}_{\text{van}})$  and  $\mathcal{LM}_{\text{opt}} = (\text{Leak}, \text{leakVia}_{\text{opt}})$  be leakage models for  $SM_{\text{van}}$  and  $SM_{\text{opt}}$  respectively (having the same set of leaks  $\text{Leak}$ ). We say that  $(SM_{\text{opt}}, \mathcal{LM}_{\text{opt}})$  satisfies *relative security* w.r.t.  $(SM_{\text{van}}, \mathcal{LM}_{\text{van}})$ , written  $(SM_{\text{opt}}, \mathcal{LM}_{\text{opt}}) \geq_{\checkmark} (SM_{\text{van}}, \mathcal{LM}_{\text{van}})$ , when the following holds:

$$\forall l \in \text{Leak}_{\text{opt}}. \forall \pi_1, \pi_2 \in \text{Trace}_{\text{opt}}. \text{leakVia}_{\text{opt}}(\pi_1, \pi_2, l) \rightarrow \exists \hat{\pi}_1, \hat{\pi}_2 \in \text{Trace}_{\text{van}}. \text{leakVia}_{\text{van}}(\hat{\pi}_1, \hat{\pi}_2, l)$$

We say that  $(SM_{\text{opt}}, \mathcal{LM}_{\text{opt}})$  satisfies *finitary relative security* w.r.t.  $(SM_{\text{van}}, \mathcal{LM}_{\text{van}})$ , written  $(SM_{\text{opt}}, \mathcal{LM}_{\text{opt}}) \geq_{\checkmark}^{\text{fin}} (SM_{\text{van}}, \mathcal{LM}_{\text{van}})$ , when the above property holds when restricted to finite traces, i.e., replacing  $\text{Trace}_{\text{opt}}$  with  $\text{Trace}_{\text{opt}}^{\text{fin}}$  and  $\text{Trace}_{\text{van}}$  with  $\text{Trace}_{\text{van}}^{\text{fin}}$ .

This definition expresses that, for the given notion of leak, the optimization-enhanced system  $SM_{\text{opt}}$  does not exhibit any leaks besides those already exhibited by the vanilla system  $SM_{\text{van}}$ . (Thinking of  $\checkmark$  as expressing security, the notation  $\geq_{\checkmark}$  suggests an “at least as secure as” reading.) The notations  $\hat{\pi}_1$  and  $\hat{\pi}_2$  for vanilla traces remind of their dependency on the otraces  $\pi_1$  and  $\pi_2$ —with the caveat that each of the two vtraces may depend on *both*  $\pi_1$  and  $\pi_2$ .

We believe the finitary (i.e., termination-conditioned) variant of relative security  $\geq_{\checkmark}^{\text{fin}}$  is of interest for both historic reasons (since often the discussion in the literature is restricted to finite traces, e.g., [14]) and pragmatic reasons

(since, as we shall see, finitary security is amenable to a simpler unwinding proof method).  $\geq^{\text{fin}}$  is the same as  $\geq_{\checkmark}$  for terminating programs such as the ones in §2.1's Listings 1–4, but  $\geq_{\checkmark}$  is more suitable for possibly nonterminating interactive programs as in our Listings 5 and 6.

## 2.4. Less abstract relative security: attacker models

Our current definition of relative security is parameterized by leaks. But still, what *is* a leak more concretely? To give a plausible answer, remember the key ingredients identified in our examples (§2.1): secret uploading and observer interaction, and our desire to capture fully interactive versions of these. These lead us to the next definition. Given a system model  $(\text{State}, \text{istate}, \Rightarrow)$ , an *attacker model* for it is a tuple  $(\text{Sec}, S, \text{Obs}, O, \text{Act}, A)$  where:

- $\text{Sec}$ ,  $\text{Obs}$  and  $\text{Act}$  are sets of items called *secrets*, *observations* and *actions* respectively;
- $S : \text{Trace} \rightarrow \text{Seq}(\text{Sec})$ ,  $O : \text{Trace} \rightarrow \text{Seq}(\text{Obs})$  and  $A : \text{Trace} \rightarrow \text{Seq}(\text{Act})$  are functions called the *secrecy*, *observation* and *action* functions, respectively.

Thus, attacker models indicate what needs protection (the secrets) and what attacker actions / observations are available.

The functions  $S$ ,  $O$  and  $A$  provide a natural notion of leak. We take  $\text{Leak}$  to be  $\text{Seq}(\text{Sec}) \times \text{Seq}(\text{Sec})$  and  $\text{leakVia}(\pi_1, \pi_2, (\sigma_1, \sigma_2))$  to be  $S(\pi_1) = \sigma_1 \wedge S(\pi_2) = \sigma_2 \wedge A(\pi_1) = A(\pi_2) \wedge O(\pi_1) \neq O(\pi_2)$ , which means:

- $\pi_1$  and  $\pi_2$  have the sequences of secrets  $\sigma_1$  and  $\sigma_2$ ;
- the attacker took the *same* actions during  $\pi_1$  and  $\pi_2$  ...
- which led to the attacker making *different* observations.

$\text{leakVia}(\pi_1, \pi_2, (\sigma_1, \sigma_2))$  therefore says that, via  $\pi_1$  and  $\pi_2$ , the attacker can observationally distinguish between  $\sigma_1$  and  $\sigma_2$ . Note that it was crucial to require that the distinction be made while the attacker takes the *same* actions—otherwise it would not tell us anything about the secrets (as it could simply be a consequence of the different actions).

Thus, any attacker model for  $SM$  induces a leakage model for  $SM$ . It is worth spelling out what the definition of relative security becomes in this more concrete setting:

**Def. 2** Let  $\mathcal{AM}_{\text{van}} = (\text{Sec}, S_{\text{van}}, \text{Obs}_{\text{van}}, O_{\text{van}}, \text{Act}_{\text{van}}, A_{\text{van}})$  and  $\mathcal{AM}_{\text{opt}} = (\text{Sec}, S_{\text{opt}}, \text{Obs}_{\text{opt}}, O_{\text{opt}}, \text{Act}_{\text{opt}}, A_{\text{opt}})$  be attacker models for  $SM_{\text{van}}$  and  $SM_{\text{opt}}$  respectively (with same set of secrets  $\text{Sec}$ ). We say that  $(SM_{\text{opt}}, \mathcal{AM}_{\text{opt}})$  satisfies *relative security* w.r.t.  $(SM_{\text{van}}, \mathcal{AM}_{\text{van}})$ , written  $(SM_{\text{opt}}, \mathcal{AM}_{\text{opt}}) \geq_{\checkmark} (SM_{\text{van}}, \mathcal{AM}_{\text{van}})$ , when:

$$\begin{aligned} & \forall \sigma_1, \sigma_2 \in \text{Seq}(\text{Sec}). \forall \pi_1, \pi_2 \in \text{Trace}_{\text{opt}}. \\ & S_{\text{opt}}(\pi_1) = \sigma_1 \wedge S_{\text{opt}}(\pi_2) = \sigma_2 \wedge \\ & A_{\text{opt}}(\pi_1) = A_{\text{opt}}(\pi_2) \wedge O_{\text{opt}}(\pi_1) \neq O_{\text{opt}}(\pi_2) \\ & \longrightarrow \\ & \exists \hat{\pi}_1, \hat{\pi}_2 \in \text{Trace}_{\text{van}}. S_{\text{van}}(\hat{\pi}_1) = \sigma_1 \wedge S_{\text{van}}(\hat{\pi}_2) = \sigma_2 \wedge \\ & A_{\text{van}}(\hat{\pi}_1) = A_{\text{van}}(\hat{\pi}_2) \wedge O_{\text{van}}(\hat{\pi}_1) \neq O_{\text{van}}(\hat{\pi}_2) \end{aligned}$$

*Finitary relative security*,  $(SM_{\text{opt}}, \mathcal{AM}_{\text{opt}}) \geq^{\text{fin}}_{\checkmark} (SM_{\text{van}}, \mathcal{AM}_{\text{van}})$ , is again defined by restricting to finite traces.

The above just expands the aforementioned construction of leakage models from attacker models, so Def. 2 is a

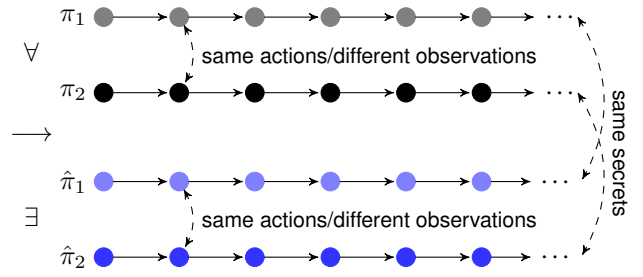


Fig. 1: Visualizing relative security

particular case of Def. 1. Its conclusion can be reformulated without explicitly quantifying over secrets as follows:

$$\begin{aligned} & \forall \pi_1, \pi_2 \in \text{Trace}_{\text{opt}}. \\ & A_{\text{opt}}(\pi_1) = A_{\text{opt}}(\pi_2) \wedge O_{\text{opt}}(\pi_1) \neq O_{\text{opt}}(\pi_2) \\ & \longrightarrow \\ & \exists \hat{\pi}_1, \hat{\pi}_2 \in \text{Trace}_{\text{van}}. \\ & S_{\text{van}}(\hat{\pi}_1) = S_{\text{van}}(\pi_1) \wedge S_{\text{van}}(\hat{\pi}_2) = S_{\text{van}}(\pi_2) \wedge \\ & A_{\text{van}}(\hat{\pi}_1) = A_{\text{van}}(\hat{\pi}_2) \wedge O_{\text{van}}(\hat{\pi}_1) \neq O_{\text{van}}(\hat{\pi}_2) \end{aligned} \quad (\star)$$

The  $(\star)$  formulation, which is the one we will prefer, facilitates an intuitive reading of relative security, as visualized in Fig. 1: Provided the otraces  $\pi_1$  and  $\pi_2$  have the same actions and different observations, the vtraces  $\hat{\pi}_1, \hat{\pi}_2$  must be proved to exist such that they have the same secrets as  $\pi_1$  and  $\pi_2$  respectively, and also have between each other the same actions and different observations.

In line with our aforementioned gradual instantiation approach, for now we have left  $S$ ,  $A$  and  $O$  unspecified. In the next section, §2.5, we make these more concrete by assuming they operate “state-wise”, i.e., the secrets, actions and observations are produced locally from each state of a traces. (Another natural choice would be to assume they operate transition-wise rather than state-wise; these two choices are equivalent, since states can also encode transition information.) For our §2.1 motivating examples, when applied to a trace:  $S$  will give the initial memory, as well as any potential secrets received or sent during the execution (e.g., on the trusted input channel in Listing 6); moreover,  $A$  will give any inputs on untrusted channels; finally,  $O$  will give any returned or printed values via untrusted channels, and the memory locations accessed for reading—see §3.3.

## 2.5. State-wise attacker models

Next, we will consider an assumption on attacker models that makes them even less abstract (thus bringing us closer to considering concrete examples). Given a system model  $SM = (\text{State}, \text{istate}, \Rightarrow)$ , an attacker model  $\mathcal{AM} = (\text{Sec}, S, \text{Obs}, O, \text{Act}, A)$  for it is said to be *state-wise* when there exist the predicates and functions  $\text{isSec} : \text{State} \rightarrow \text{Bool}$ ,  $\text{getSec} : \text{State} \rightarrow \text{Sec}$ ,  $\text{isInt} : \text{State} \rightarrow \text{Bool}$ ,  $\text{getObs} : \text{State} \rightarrow \text{Obs}$ , and  $\text{getAct} : \text{State} \rightarrow \text{Act}$  that define the functions  $S$ ,  $O$  and  $A$  state-wise as follows. For any trace  $\pi = s_0 s_1 \dots$ :



- Let  $s_{i_0} s_{i_1} \dots$  for  $i_0 < i_1 < \dots$  be its subsequence consisting of states where `isSec` holds. We define  $S(\pi)$  as  $\text{getSec}(s_{i_0}) \text{getSec}(s_{i_1}) \dots$
- Let  $s_{j_0} s_{j_1} \dots$  for  $j_0 < j_1 < \dots$  be its subsequence consisting of the states where `isInt` holds. We define  $O(\pi)$  as  $\text{getObs}(s_{j_0}) \text{getObs}(s_{j_1}) \dots$  and  $A(\pi)$  as  $\text{getAct}(s_{j_0}) \text{getAct}(s_{j_1}) \dots$

What is being required above is that the trace functions  $S$ ,  $O$  and  $A$  are defined by “filtermap”, filtering with a predicate and mapping with a getter function.

We think of the predicates `isSec` and `isInt` as determining whether (the next transition from) a state uploads a secret, and is a point of interaction (observation or action), respectively. For example, in our concrete programming language models, the program counter stored in the state will determine the next statement to be executed, and therefore `isSec` and `isInt` will check whether this next statement is secret-uploading or interaction-producing. (We do not require `isSec` and `isInt` to be disjoint, although for some systems this can be a reasonable assumption for *a priori* excluding obvious leaks.) Moreover, we think of the functions `getSec`, `getObs` and `getAct` as actually extracting that particular secret, observation or action.

### 3. Language-Based Instantiation

Next we describe a concrete instance of relative security (§3.3), via a programming language with speculative execution (§3.1) that can express our running examples (§3.2).

#### 3.1. The IMP Language with Speculative Semantics

We introduce IMP, a simple language that exhibits interesting security aspects stemming from speculative execution. Its speculative semantics follows the ideas of Cheang et al. [5], maintaining runtime configurations for nested speculative executions—and only for those that result from *misprediction*, since they are the only security-relevant ones.

**Syntax.** The set  $\text{Lit}$  of *literals*, ranged over by  $i, j$  etc., is taken to be  $\mathbb{N}$ .  $\text{Var}$ , the set of (*scalar*-)variables, ranged over by  $x, y, z$  etc., is a fixed countably infinite set; and so is  $\text{AVar}$ , the set of *array-variables*, ranged over by  $a$ .  $\text{Op}$  is the set of binary arithmetic operators, e.g.,  $+$ ,  $*$ , etc;  $\text{COp}$ , that of binary comparison operators, e.g.,  $<$ ,  $=$ , etc; and  $\text{BOp}$  that of binary boolean operators, e.g.,  $\wedge$ ,  $\vee$ , etc.  $\text{ICh}$ , ranged over by  $ich$ , and  $\text{OCh}$ , ranged over by  $och$ , are sets of input and output channels, respectively.

The sets of (arithmetic) expressions, boolean expressions and commands are defined by the grammar:

$$\begin{aligned}
\text{Exp} &::= \text{Lit} \mid \text{Var} \mid \text{AVar}[\text{Exp}] \mid \text{Exp Op Exp} \\
\text{BExp} &::= \text{true} \mid \text{false} \mid \text{Exp COp Exp} \mid \\
&\quad \text{not BExp} \mid \text{BExp BOp BExp} \\
\text{Com} &::= \text{Start} \mid \text{Input}_{\text{ICh}} \text{Var} \mid \text{Output}_{\text{OCh}} \text{Exp} \mid \\
&\quad \text{Fence} \mid \text{Var} = \text{Exp} \mid \text{AVar}[\text{Exp}] = \text{Exp} \mid \\
&\quad \text{Jump } \mathbb{N} \mid \text{IfJump BExp } \mathbb{N} \mathbb{N}
\end{aligned}$$

We let  $e$  range over  $\text{Exp}$ ,  $b$  over  $\text{BExp}$ , and  $c$  over  $\text{Com}$ .

$$\begin{array}{c}
\text{STARTORFENCEOROUTPUT} \\
\frac{c_{pc} \in \{\text{Start}, \text{Fence}\} \cup \{\text{Output}_{och} \mid e \in \text{Exp}\}}{((pc, \mu), inp) \Rightarrow_B ((pc + 1, \mu), inp)}
\end{array}$$

$$\begin{array}{c}
\text{VARASSIGN} \\
\frac{c_{pc} = (x = e)}{((pc, \mu), inp) \Rightarrow_B ((pc + 1, \mu[x \leftarrow \llbracket e \rrbracket(\mu)]), inp)}
\end{array}$$

$$\begin{array}{c}
\text{AVARASSIGN} \\
\frac{c_{pc} = (a[e] = e')}{((pc, \mu), inp) \Rightarrow_B ((pc + 1, \mu[a, \llbracket e \rrbracket(\mu)] \leftarrow \llbracket e' \rrbracket(\mu)]), inp)}
\end{array}$$

$$\begin{array}{c}
\text{INPUT} \\
\frac{c_{pc} = (\text{Input}_{ich} \ x) \quad inp_{ich} = i \cdot is'}{((pc, \mu), inp) \Rightarrow_B ((pc + 1, \mu[x \leftarrow i]), inp[ich \leftarrow is'])}
\end{array}$$

$$\begin{array}{c}
\text{JUMP} \\
\frac{c_{pc} = (\text{Jump } pc')}{((pc, \mu), inp) \Rightarrow_B ((pc', \mu), inp)}
\end{array}$$

$$\begin{array}{c}
\text{IFJUMP} \\
\frac{c_{pc} = (\text{IfJump } b \ pc_1 \ pc_2) \quad pc' = (\text{if } \llbracket b \rrbracket(\mu) \text{ then } pc_1 \text{ else } pc_2)}{((pc, \mu), inp) \Rightarrow_B ((pc', \mu), inp)}
\end{array}$$

Fig. 2: Basic semantics for program  $P = c_0; \dots; c_n$ . We implicitly assume  $pc \leq n$  as a condition in each rule.

Thus, IMP has the standard basic mechanisms for manipulating scalar and array variables, and (un)conditional jumps, `Jump` and `IfJump`, as control structures. It is also an I/O interactive language, accepting inputs on input channels and producing outputs on output channels. A *program*  $P = c_0; c_1; \dots; c_n$  is a nonempty list of commands where  $c_0 = \text{Start}$ .  $\text{Prog}$  denotes the set of programs.

The set  $\text{Val}$  of *values*, ranged over by  $v, w$  etc., is  $\mathbb{N}$ .  $\text{Loc}$ , the set of *locations*, ranged over by  $l$ , is also  $\mathbb{N}$ .

**Basic semantics.** Fig. 2 shows the *basic (small-step) semantics*, denoted  $\Rightarrow_B$ , parameterized by a fixed program  $P = c_0; \dots; c_n$ . It maintains input streams and memories, which are consumed and respectively updated while the program counter moves through the program’s list of commands.

In detail,  $\Rightarrow_B$  is a relation between pairs  $(cfg, inp)$  where  $inp : \text{ICh} \rightarrow \text{Seq}(\text{Val})$  is an input-channel indexed family of *input streams*, and  $cfg$  is a *configuration*, i.e., a pair  $(pc, \mu)$ :

- $pc \in \mathbb{N}$  is the *program counter (PC)*, pointing to the  $pc$ ’th statement in the program,  $c_{pc}$ .
- $\mu$  is a *memory*, i.e., a triple  $(vs, avs, hp)$  where  $vs : \text{Var} \rightarrow \text{Val}$  is a *variable store* assigning values to the (scalar) variables,  $avs : \text{AVar} \rightarrow \text{Loc} \times \mathbb{N}$  is an *array-variable store* assigning to each array-variable its starting location in the heap and its size, and  $hp : \text{Loc} \rightarrow \text{Val}$  is a *heap* assigning values to locations.

If  $cfg = (pc, (vs, avs, hp))$ , we let  $\text{pcOf}(cfg) = pc$ . We use standard notation for updates, e.g.,  $\mu[x \leftarrow v]$ , and expression evaluation.  $\llbracket e \rrbracket(\mu)$ .  $A_\perp$  denotes the extension of a set  $A$  with an “undefined” element  $\perp$ . The *output of a con-*

$$\text{STANDARD} \quad \frac{(cfg, inp) \Rightarrow_B (cfg', inp')}{(cfg, inp, L) \Rightarrow_N (cfg', inp', L \cup \text{readLocs}(cfg))}$$

Fig. 3: Normal semantics for program  $P = c_0; \dots; c_n$ . We implicitly assume  $pc \leq n$ .

$$\text{IFJUMPISPRED} \quad \frac{\begin{array}{l} c_{pc} = (\text{IfJump } b \text{ } pc_1 \text{ } pc_2) \\ pc' = (\text{if } \llbracket b \rrbracket(\mu) \text{ then } pc_2 \text{ else } pc_1) \end{array}}{((pc, \mu), inp) \Rightarrow_M ((pc', \mu), inp)}$$

$$\text{STANDARD} \quad \frac{\begin{array}{l} \neg \text{isCond}(cfg_k) \vee \neg \text{mispred}(ps, pcs) \\ (k > 0 \rightarrow \neg \text{isOFence}(cfg_k) \wedge \neg \text{resolve}(ps, pcs)) \\ (cfg_k, inp) \Rightarrow_B (cfg', inp') \end{array}}{C' = cfg_0 \dots \dots \cdot cfg_{k-1} \cdot cfg' \quad L' = L \cup \text{readLocs}(cfg_k)} \\ (ps, cfg_0 \dots \dots \cdot cfg_k, inp, L) \Rightarrow_S (ps, C', inp', L')$$

$$\text{MISPRED} \quad \frac{\begin{array}{l} \text{isCond}(cfg_k) \quad \text{mispred}(ps, pcs) \\ (cfg_k, inp) \Rightarrow_B (cfg', inp') \quad (cfg_k, inp) \Rightarrow_M (cfg'', inp'') \end{array}}{C' = cfg_0 \dots \dots \cdot cfg_{k-1} \cdot cfg' \cdot cfg'' \quad L' = L \cup \text{readLocs}(cfg_k)} \\ (ps, cfg_0 \dots \dots \cdot cfg_k, inp, L) \Rightarrow_S (\text{update}(ps, pcs), C', inp', L')$$

$$\text{RESOLVE} \quad \frac{k > 0 \quad \text{resolve}(ps, pcs) \quad C' = cfg_0 \dots \dots \cdot cfg_{k-1}}{(ps, cfg_0 \dots \dots \cdot cfg_k, inp, L) \Rightarrow_S (\text{update}(ps, pcs), C', inp, L)}$$

$$\text{FENCE} \quad \frac{k > 0 \quad \neg \text{resolve}(ps, pcs) \quad \text{isFence}(cfg_k)}{(ps, cfg_0 \dots \dots \cdot cfg_k, inp, L) \Rightarrow_S (pcs, cfg_0, inp, L)}$$

Fig. 4: Speculative semantics for program  $P = c_0; \dots; c_n$  under misprediction oracle (PState, mispred, resolve, update). Each time, we implicitly assume  $pc \leq n$  and  $pcs = pc_0 \dots \dots pc_k$  where  $pc_i = \text{pcOf}(cfg_i)$  for  $i \in \{0, \dots, k\}$ .

figuration,  $\text{outOf}(pc, \mu) \in (\text{OCh} \times \text{Val})_\perp$ , is  $(och, \llbracket e \rrbracket(\mu))$  if  $c_{pc}$  has the form  $\text{Output}_{och} e$ , and  $\perp$  otherwise.

**The read location.** For modeling Spectre vulnerabilities, we will record memory reads (as in [5]). We let  $\text{readLocs}(pc, \mu)$  be the (possibly empty) set of locations that are read by the current command  $c_{pc}$ —computed from all sub-expressions of  $c_{pc}$  of the form  $a[e]$ . For example, if  $c_{pc}$  is the assignment  $x = a[b[3]]$ , then  $\text{readLocs}$  returns two locations: counting from 0, the 3rd location of  $b$  and the  $b[3]$ ’th location of  $a$ .

**Normal semantics.** Fig. 3 shows the “normal” semantics  $\Rightarrow_N$ , which is the basic semantics extended to accumulate the read locations, hence account for cache side-channels.

**Speculative semantics.** Finally, Fig. 4 shows the *speculative semantics*  $\Rightarrow_S$ , which augments normal semantics with speculative steps that go wrong, along a *misprediction* (taking the wrong branch). In addition to the program  $P = c_0; \dots; c_n$ ,  $\Rightarrow_S$  is parameterized by a *misprediction oracle*, i.e., a tuple (PState, mispred, resolve, update) where: PState is a set of *predictor states* ranged over by  $ps$ ,  $\text{mispred} : \text{PState} \times \{0, \dots, n\}^* \rightarrow \text{Bool}$ ,  $\text{resolve} : \text{PState} \times \{0, \dots, n\}^* \rightarrow \text{Bool}$ , and  $\text{update} : \text{PState} \times \{0, \dots, n\}^* \rightarrow \text{PState}$ .

The oracle decides for when misprediction occurs thus triggering a new speculation (the mispred predicate), and for when a speculation is being resolved (the resolve predicate). Both occurrence and resolution depend on the predictor state (which evolves via the function update) and on the list of PCs of the nested speculations, an element of  $\{0, \dots, n\}^*$ . Indeed, the semantics allows for nested speculation: The runtime environment can mispredict thus taking a wrong branch, and while running on that wrong branch it can mispredict again, and so on. Resolution, i.e., the act of the runtime environment figuring out that the prediction was wrong and reverting the corresponding speculative execution, occurs when the oracle dictates that, via resolve.

The speculative semantics  $\Rightarrow_S$  operates on (*multi-speculative*) states, i.e., tuples  $s = (ps, cfs, inp, L)$  where  $ps \in \text{PState}$  is the current predictor state,  $cfs$  is a non-empty list of configurations.  $inp$  is a family of input streams and  $L$  is the set of locations read so far. We think of  $cfs$  as a stack of configurations, one for each *speculation level* in a nested speculative execution. At level 0 we have the configuration  $cfg_0$  for normal, non-speculative execution.

At each moment, only the top of the configuration stack,  $cfg_k$ , is active. One type of transitions that  $cfg_k$  can take is standard transitions, shown in the rule STANDARD. These occur only when (1) the misprediction option is not on the table (either because the command of  $cfg_k$  is not a conditional, or because the oracle does not demand a misprediction at this time) and (2) if in speculation mode, i.e.,  $k > 0$ , no other blocking factors occur.

If the command of  $cfg_k$  is a conditional and the oracle demands misprediction, then a new mispredicting speculation is launched—as shown in the rule MISPRED. In this case,  $cfg_k$  takes a normal execution step on the correct branch (yielding the configuration  $cfg'$ ) while at the same time  $cfg_k$  takes a step on the wrong branch (yielding the configuration  $cfg''$  at speculation level  $k + 1$ , which becomes the active configuration). The step on the wrong branch is achieved using a dual of the  $\Rightarrow_B$  transition for conditional commands, denoted  $\Rightarrow_M$  (where  $M$  stands for “misprediction”)—also shown in Fig. 4, as the rule IFJUMP-MISPRED. The fact that, in the semantics, the new speculation execution moves on the wrong branch at the same time with execution one level below moving on the correct branch is a succinct way of expressing that the new speculation involves misprediction (as opposed to correct prediction).

When in speculative mode ( $k > 0$ ) and the oracle demands it, a resolution occurs: The  $k$ -level speculation is reverted, i.e., the configuration  $cfg_k$  is dropped and  $cfg_{k-1}$  re-becomes the active configuration—as shown by the rule RESOLVE. Finally, when in speculative mode and encountering a Fence command, the entire stack of nested speculations is reverted, leaving only the configuration at level 0 (for normal execution)—as shown by the rule FENCE.

### 3.2. Representing the running examples in IMP

We assume  $\text{ICh} = \text{OCh} = \{T, U\}$ , where  $T$  and  $U$  represent a trusted and an untrusted channel, respectively. More-

over, we assume a binary operation  $F \in \text{Op}$  whose semantics is some (fixed but) unspecified function in  $\text{Val}^2 \rightarrow \text{Val}$ .

Our motivating examples from §2.1 are expressed in IMP as depicted in Fig. 5. Expectedly, the if statements are modeled by IfJump, and the while loops by IfJump in conjunction with unconditional Jump. In fun1–fun5 we use only the untrusted channel  $U$  for inputs and outputs—thus matching our §2.1’s assumption that the attacker controls the inputs and sees the outputs. In fun6 we use both the trusted and untrusted channels, matching our assumption that the attacker controls only specific inputs and outputs. Since our language lacks procedures, we use  $\text{Output}_T(F(x, y))$  to encode the  $\text{writeOnSecretFile}(x, y)$  procedure of fun6.

### 3.3. Instantiating relative security to IMP

We fix an IMP program  $P = c_0; \dots; c_n$ , and some initial memory  $\mu_0$ , inputs  $\text{inp}_0$  and predictor state  $ps_0$ . We first instantiate the system models  $\mathcal{SM}_{\text{van}}$  and  $\mathcal{SM}_{\text{opt}}$ :

- $\text{State}_{\text{van}}$  consists of non-speculative states, i.e., triples  $(\text{cfg}, \text{inp}, L)$ ;  $\Rightarrow_{\text{van}}$  is the normal semantics,  $\Rightarrow_N$ ; and  $\text{istate}_{\text{van}}(\text{cfg}, \text{inp}, L)$  says that  $\text{cfg} = (0, \mu_0)$ ,  $\text{inp} = \text{inp}_0$  and  $L = \emptyset$ .
- $\text{State}_{\text{opt}}$  consist of the (multi-speculative) states;  $\Rightarrow_{\text{opt}}$  is  $\Rightarrow_S$ ; and  $\text{istate}_{\text{opt}}(s)$  for  $s = (ps, \text{cfgs}, \text{inp}, L)$  says that  $ps = ps_0$ ,  $\text{inp} = \text{inp}_0$ ,  $L = \emptyset$  and  $\text{cfgs}$  is a (singleton) list consisting of one configuration,  $(0, \mu_0)$ .

Thus, both systems start with the specified initial memory and input stream, and no read locations or speculation.

For instantiating the state-wise attacker models  $\mathcal{AM}_{\text{opt}}$  and  $\mathcal{AM}_{\text{van}}$ , we note that the non-speculative states in  $\text{State}_{\text{van}}$  can be seen as particular cases of (multi-speculative) states in  $\text{State}_{\text{opt}}$ . Therefore, agreeing to take the operators for  $\mathcal{AM}_{\text{van}}$  as the restrictions of those for  $\mathcal{AM}_{\text{opt}}$ , we will allow ourselves to only define the latter, while also omitting the subscript  $\text{opt}$ .

To determine  $\text{isInt}$  and  $\text{getInt}$ , we ask what the attacker’s capabilities are. Of course, the attacker controls/observes untrusted inputs/outputs. Moreover, we assume a strong attacker observing execution time and control flow, following the leakage model of *constant-time security* [4]. Finally, we assume the attacker observes the read locations (e.g., via probing the cache)—another usual assumption of constant-time security. We thus define  $\text{isInt}(s) = \neg \text{final}(s)$  and  $\text{getInt}(s) = (\text{getAct}(s), \text{getObs}(s))$  where, if  $s = (ps, \text{cfg}_0 \dots \text{cfg}_k, \text{inp}, L)$ , we have:

$$\begin{aligned} \text{getAct}(s) &= \begin{cases} \text{head}(\text{inp}_U), & \text{if } k = 0 \text{ and } \text{isInput}_U(\text{cfg}_0) \\ \perp, & \text{otherwise} \end{cases} \\ \text{getObs}(s) &= \begin{cases} (\text{outOf}(\text{cfg}_0), L), & \text{if } k = 0 \text{ and } \text{isOutput}_U(\text{cfg}_0) \\ \perp, & \text{otherwise} \end{cases} \end{aligned}$$

As shown by the definition of  $\text{isInt}$ , interaction happens everywhere except for the final states (where the system is idle)—this pervasiveness of interaction means that the attacker can observe the execution time. For most interaction points,  $\text{getInt}$  (which pairs actions via  $\text{getAct}$  with observations via  $\text{getObs}$ ) reveals nothing beyond the fact that an

execution step has been taken, i.e., returns  $(\perp, \perp)$ ; exceptions are when the current command is an untrusted input or output command (which we expressed by the  $\text{isInput}_U$  and  $\text{isOutput}_U$  predicates). Note also that inputs and outputs can only take place non-speculatively, i.e., at speculation level 0.

The second question, which will determine  $\text{isSec}$  and  $\text{getSec}$ , is what the sensitive data, i.e., the secrets, are. For the bulk of our examples, fun1–fun5, we take  $\text{isSec}(s)$  to say that  $s$  is initial and  $\text{getSec}(s)$  to return the memory part of this initial state (namely  $\mu_0$ ). In short, here the entire initial memory constitutes the secrets.

For fun6, we change  $\text{isSec}$  and  $\text{getSec}$  to account for the interactive nature of the secrets, entered in the system as trusted inputs and produced as trusted outputs. Thus,  $\text{isSec}$  now says that the state is not final and not speculating and, for a state  $s = (ps, \text{cfgs}, \text{inp}, L)$ ,  $\text{getSec}(s)$  now returns triples  $(A, B, C)$  where:

- $A$  is (as before) the memory part of  $s$  provided  $s$  is initial, and  $\perp$  otherwise;
- $B$  is  $\text{head}(\text{inp}_T)$  if the command of  $\text{cfgs}_0$  is an  $\text{Input}_T$  (trusted input) command, and  $\perp$  otherwise;
- $C$  is either the (trusted) output produced if the command of  $\text{cfgs}_0$  is an  $\text{Output}_T$  one, and  $\perp$  otherwise.

In short, the secrets are now three-fold: memory if in initial state, trusted input if any, and trusted output if any.

With this instantiation, one can check that the intuitive (in)security of these examples discussed in §2.1 corresponds to their (not) satisfying relative security. For example, fun1 is not relatively secure essentially because the attacker can use the  $\text{Input}_U$   $x$  command to pass an input  $x \geq N$ , which leaks  $a[x]$  (part of the initial memory) via the read-locations component of the observation ( $L$ ). Indeed, the value of  $a[x]$  is collected by the semantics among the read locations while executing line 4 in the speculative semantics at speculation level 1; but not under normal semantics, under which line 4 is not executed because  $x < N$ . This violation of relative security is impossible for fun2 and fun3, where the fences stop speculation via the FENCE rule in Fig. 4 before the location  $a[x]$  can be collected. In fact, fun2–fun6 are all relatively secure, like our informal analysis in §2.1 concluded. But actual proofs for these will be based on the unwinding proof method we introduce in §5, and will be sketched in §6.

## 4. Connections with notions from the literature

Of the vast literature on Spectre and Meltdown [3], we will only discuss formal modeling and verification aspects.

**The contract/policy pattern.** Most of the approaches focus on protecting the secrecy of (aspects of) the initial state, expressed as an indistinguishability relation  $\simeq$  on  $\text{State}$ , and on explicitly modeling only the attacker’s observations (our function  $O$ ) and not the attacker’s actions (our function  $A$ ). Notably, this is the case of Guanciale et al.’s *conditional non-interference* [13], an extension of noninterference [11, 29] used to build a detailed formal model of Spectre vulnerabilities in the presence of speculative and out-of-order execution; and of Guarnieri et al.’s *speculative non-*



<pre> 0 : Start ; 1 : Input<sub>U</sub> x ; 2 : t = 0 ; 3 : IfJump (x &lt; N) 4 5 ; 4 :   t = b[a[x] * 512] ; 5 : Output<sub>U</sub> t </pre> <p style="text-align: center;">fun1</p>	<pre> 0 : Start ; 1 : Input<sub>U</sub> x ; 2 : t = 0 ; 3 : IfJump (x &lt; N) 4 6 ; 4 :   Fence ; 5 :   t = b[a[x] * 512] ; 6 : Output<sub>U</sub> t </pre> <p style="text-align: center;">fun2</p>	<pre> 0 : Start ; 1 : Input<sub>U</sub> x ; 2 : t = 0 ; 3 : IfJump (x &lt; N) 4 7 ; 4 :   v = a[x] ; 5 :   Fence ; 6 :   t = b[v * 512] ; 7 : Output<sub>U</sub> t </pre> <p style="text-align: center;">fun3</p>	<pre> 0 : Start ; 1 : Input<sub>U</sub> x ; 2 : t = 0 ; 3 : IfJump (x &lt; N) 4 6 ; 4 :   v = a[0] ; 5 :   t = b[v * 512] ; 6 : Output<sub>U</sub> t </pre> <p style="text-align: center;">fun4</p>	
<pre> 0 : Start ; 1 : t = 0 ; 2 : x = 1 ; 3 : IfJump (not (x == 0)) 4 11 ; 4 :   Input<sub>U</sub> x ; 5 :   IfJump (x &lt; N) 6 10 ; 6 :     v = a[x] ; 7 :     Fence ; 8 :     t = b[v * 512] ; 9 :   Output<sub>U</sub> t ; 10 :  Jump 3 ; 11 : Output<sub>U</sub> 0 ; </pre> <p style="text-align: center;">fun5</p>				<pre> 0 : Start ; 1 : t = 0 ; 2 : x = 1 ; 3 : IfJump (not (x == 0)) 4 13 ; 4 :   Input x ; 5 :   Input<sub>T</sub> y ; 6 :   IfJump (x &lt; N) 7 12 ; 7 :     v = a[x] ; 8 :     Output<sub>T</sub> (F(x, y)) ; 9 :     Fence ; 10 :    t = b[v * 512] ; 11 :   Output<sub>U</sub> t ; 12 :  Jump 3 ; 13 : Output<sub>U</sub> 0 ; </pre> <p style="text-align: center;">fun6</p>

Fig. 5: Our six motivating examples from §2.1 written in IMP

*interference* [14] (which forms the basis of Spectator, an automatic tool for proving security against Spectre).

These and several other approaches are surveyed in a recent SoK paper by Cauligi et al. [4], which, borrowing concepts and terminology from Guarnieri et al. [15], describes these approaches uniformly under the following *contract/policy* pattern, characterized by three parameters: **1**) an *execution model*  $\alpha$ , indicating the states explored during executions, **2**) a *leakage model*  $l$ , indicating the observations that are possible along executions, and **3**) a (*secrecy*) *policy*  $p$ , indicating the secrets stored in the initial state, via an indistinguishability relation on states  $\simeq_p$ .

An execution model  $\alpha$  corresponds to an operational semantics (our system models  $\mathcal{SM} = (\text{State}, \text{istate}, \Rightarrow)$ ). The combination between a leakage and an execution model, called a *contract*, determines a state-observation function  $\llbracket \cdot \rrbracket_l^\alpha : \text{State} \rightarrow \text{Seq}(\text{Obs})$ . Given execution and leakage models  $\alpha$  and  $l$  and a policy  $p$ , direct noninterference states that the observation function cannot detect any secret, in that  $\forall s_1, s_2 \in \text{State}. s_1 \simeq_p s_2 \longrightarrow \llbracket s_1 \rrbracket_l^\alpha = \llbracket s_2 \rrbracket_l^\alpha$ . Relative noninterference generalizes this to two contracts, namely (in our terminology) a vanilla one  $(\alpha, l)$  and an optimization-enhanced one  $(\alpha', l')$ , and a common policy  $p$ , stating that any secrets that are leaked by  $(\alpha', l')$  are also leaked by  $(\alpha, l)$ , in that  $\forall s_1, s_2 \in \text{State}. s_1 \simeq_p s_2 \wedge \llbracket s_1 \rrbracket_{l'}^{\alpha'} = \llbracket s_2 \rrbracket_{l'}^{\alpha'} \longrightarrow \llbracket s_1 \rrbracket_l^\alpha = \llbracket s_2 \rrbracket_l^\alpha$ .

Assuming that the state-observation function stems from a trace observation function (which seems true in all interesting cases), the contract/policy pattern is an instance of our relative security by instantiating our attacker models:  $S(\pi)$  as a singleton sequence, namely the  $\simeq_p$ -equivalence class of the trace's starting state,  $O(\pi)$  as the trace observation function underlying  $\llbracket \cdot \rrbracket_l^\alpha$ , and  $A(\pi)$  as the empty sequence.

The attacker actions  $A$  playing no role in the above models and the secrets being restricted to the initial state means that such models, while easily capturing examples like the ones in our Listings 1–4, are not directly suitable for more interactive examples as in Listings 5 and 6. While interaction features are not completely out of scope for these models, we believe our approach to allow interaction natively in the abstract models can simplify reasoning.

**TPOD.** Cheang et al.'s *TPOD* (*trace-property dependent observational nondeterminism*) [5] is an exception to the above rule, in that it explicitly captures both active attackers and interactive uploading of the secrets. TPOD extends observational determinism [35] from a two-trace to a four-trace property. It is parameterized by a notion of low-equivalence  $\equiv_{\text{low}}$  on states, describing what an attacker cannot distinguish, and by low and high operations taken at each transition between states:  $\text{op}_{\text{low}}(s)$  and  $\text{op}_{\text{high}}(s)$  are the high and low operations applied when transiting from state  $s$ . Low equivalence and the high/low operations are extended from states to traces componentwise. Finally, rather than considering two system models, (in our terminology) a vanilla one and an optimization-enhanced one, TPOD uses a single system for both and instead uses a set  $T$  of (what we call) *vtraces*. With these parameters fixed, TPOD is expressed as follows (where  $\pi_i^0$  denotes the starting state of  $\pi_i$ ):

$$\begin{aligned}
& \forall \pi_1, \pi_2 \in \text{Trace} \setminus T. \quad \forall \hat{\pi}_1, \hat{\pi}_2 \in T. \\
& \quad \text{op}_{\text{low}}(\hat{\pi}_1) = \text{op}_{\text{low}}(\hat{\pi}_2) \equiv \text{op}_{\text{low}}(\pi_2) = \text{op}_{\text{low}}(\pi_1) \wedge \\
& \quad \text{op}_{\text{high}}(\hat{\pi}_1) = \text{op}_{\text{high}}(\pi_1) \wedge \text{op}_{\text{high}}(\hat{\pi}_2) = \text{op}_{\text{high}}(\pi_2) \wedge \\
& \quad \hat{\pi}_1 \equiv_{\text{low}} \hat{\pi}_2 \wedge \pi_1^0 \equiv_{\text{low}} \pi_2^0 \\
& \quad \longrightarrow \pi_1 \equiv_{\text{low}} \pi_2
\end{aligned}$$

We can parse TPOD in terms of our relative security ingredients: **1**) the sequence of low-equivalence classes

$s_0/\equiv_{\text{low}} s_1/\equiv_{\text{low}} \dots$  of a trace  $\pi = s_0 s_1 \dots$  form the observations  $O(\pi)$ ; **2**) the low operations  $\text{op}_{\text{low}}(\pi)$  form the actions  $A(\pi)$ ; **3**) the high operations  $\text{op}_{\text{high}}(\pi)$  form the secrets  $S(\pi)$ .

However, TPOD is *not* a particular case of our relative security because of two reasons, highlighted in the above formula. First, as shown by the highlighted equality, the TPOD formula constrains the vtraces  $\hat{\pi}_1$  and  $\hat{\pi}_2$  to have the same low operations (in our terminology, the same actions) not only with each other, but also with their counterpart otraces  $\pi_1$  and  $\pi_2$ ; this (over)constrains any leak exhibited by the optimization-enhanced system to be reproduced by the vtraces under the same attacker actions, which in our opinion is not justified. Indeed, our design of relative security suggests that the similarities between the vtraces and the otraces should refer to the underlying secrets, not to the attacker-taken actions; this is because the standard assumption is that the attacker is free to take *any* actions to exhibit a leak. Interestingly, this overconstraining problem with TPOD is illustrated by an example in the TPOD paper itself [5], namely the conditionally vulnerable program we showed in Listing 4 (Fig. 3(c) in [5]): Assuming  $N > 0$ , under the interpretation of low actions as the inputs to the function `fun4` (which is natural, and is the one endorsed in [5]), TPOD requires that the vtraces reproduce the leak with the same input—which is impossible because, as we explained in §2.1 (following a discussion from [5]), reproducing the leak in the vanilla semantics requires a specific input smaller than  $N$  which is different from the one in the otraces, e.g., if  $N = 2$  then we need the vanilla input to be 0 or 1. Thus, contrary to the authors’ suggestion, Listing 4’s example does *not* satisfy TPOD; though it satisfies relative security.

The other difference between relative security and TPOD is shown in the highlighted quantifier: Because it quantifies universally rather than existentially over the vtraces  $\hat{\pi}_1$  and  $\hat{\pi}_2$ , the TPOD formula asks that any leak of  $\pi_1$  and  $\pi_2$  is reproduced not just by *some* pair of vtraces  $\hat{\pi}_1$  and  $\hat{\pi}_2$  (which seems natural), but rather, more demandingly, by *all* pairs of vtraces that happen to have the same secrets with  $\pi_1$  and  $\pi_2$ .

In summary, TPOD is a strong property that mischaracterizes intuitively secure programs that our relative security characterizes correctly. However, it is the TPOD strength that enables successful automatic verification which deems secure several interesting programs (as elaborated in [5]).

## 5. (Dis)Proof Methods for Relative Security

This section develops incremental proof and disproof methods for relative security, which can be enabled provided the secrets, observations and actions of attacker models are themselves defined incrementally on traces (§2.5). After discussing design challenges stemming from the four-trace constraint system specific to relative security (§5.1), we converge to a definition of unwinding relation (§5.2); we distinguish between the finite-trace case and the general case, the former allowing for simpler conditions. Our discussion culminates with the formal statement that unwinding indeed ensures relative security (§5.3). We also introduce a compositional variant of unwinding (§5.4). Finally, we look into the

dual problem, of incrementally *disproving* relative security, which leads to a proof method that we call secret-directed unwinding because it is a form of unwinding that shows the impossibility of saturating given sequences of secrets (§5.5).

We fix two system models  $\mathcal{SM}_{\text{van}} = (\text{State}_{\text{van}}, \text{istate}_{\text{van}}, \Rightarrow_{\text{van}})$  and  $\mathcal{SM}_{\text{opt}} = (\text{State}_{\text{opt}}, \text{istate}_{\text{opt}}, \Rightarrow_{\text{opt}})$  such that  $\text{State}_{\text{van}} \subseteq \text{State}_{\text{opt}}$ , and state-wise attacker models for them,  $\mathcal{AM}_{\text{van}} = (\text{Sec}, S_{\text{van}}, \text{Obs}_{\text{van}}, O_{\text{van}}, \text{Act}_{\text{van}}, A_{\text{van}})$  and  $\mathcal{AM}_{\text{opt}} = (\text{Sec}, S_{\text{opt}}, \text{Obs}_{\text{opt}}, O_{\text{opt}}, \text{Act}_{\text{opt}}, A_{\text{opt}})$  (having the same set of secrets  $\text{Sec}$ ). Recall from §2.5 that the attacker models being state-wise means the existence of the predicates  $\text{isSec}_u : \text{State}_u \rightarrow \text{Bool}$ ,  $\text{getSec}_u : \text{State}_u \rightarrow \text{Sec}$ ,  $\text{isInt}_u : \text{State}_u \rightarrow \text{Bool}$ ,  $\text{getObs}_u : \text{State}_u \rightarrow \text{Obs}_u$ , and  $\text{getAct}_u : \text{State}_u \rightarrow \text{Act}_u$  that define the functions  $S_u$ ,  $O_u$  and  $A_u$  state-wise, where  $u \in \{\text{opt}, \text{van}\}$ .

We will describe methods for (dis)proving  $(\mathcal{SM}_{\text{opt}}, \mathcal{AM}_{\text{opt}}) \geq_{\checkmark} (\mathcal{SM}_{\text{van}}, \mathcal{AM}_{\text{van}})$  and  $(\mathcal{SM}_{\text{opt}}, \mathcal{AM}_{\text{opt}}) \geq_{\text{fin}}^{\text{fin}} (\mathcal{SM}_{\text{van}}, \mathcal{AM}_{\text{van}})$ . We will only use the subscripts `van` and `opt` for the state-set components, e.g.,  $\text{State}_{\text{van}}$  and  $\text{State}_{\text{opt}}$ , and omit them for the other components, e.g.,  $\Rightarrow$ .

### 5.1. Design aspects

Unwinding is a (bi)simulation-like [32] method specialized in proving noninterference and related two-trace properties [12, 21]. It exhibits a winning strategy for a two-player game that incrementally follows a trace  $\pi_1$  controlled by an antagonist, while constructing a similar trace  $\pi_2$  controlled by a protagonist while “countering” any possible leak.

For relative security, there are additional challenges when designing an unwinding-like proof method since we must cope with not two but four traces  $\pi_1, \pi_2, \hat{\pi}_1, \hat{\pi}_2$ , as depicted in Fig. 1. Here, the otraces  $\pi_1$  and  $\pi_2$  on the one hand, and also their counterpart vtraces  $\hat{\pi}_1$  and  $\hat{\pi}_2$  on the other, both pairs correspond to the pair of traces from traditional unwinding. However, there are different requirements on  $\pi_1, \pi_2$  (which are *allowed* to exhibit a leak) and  $\hat{\pi}_1, \hat{\pi}_2$  (which must *reproduce* the leak exhibited by  $\pi_1, \pi_2$ ).

We are after a mechanism for building the vtraces  $\hat{\pi}_1$  and  $\hat{\pi}_2$  incrementally from the otraces  $\pi_1$  and  $\pi_2$  in such a way that they create the same leak, i.e., take the same actions *and* generate the same secrets *yet* produce different observations. During the build process, we will have to simultaneously maintain the following relationships between these traces:

- (R1) the otraces  $\pi_1$  and  $\pi_2$  have the same actions but different observations;
- (R2) the vtraces  $\hat{\pi}_1$  and  $\hat{\pi}_2$  have the same actions but different observations;
- (R3) the otrace  $\pi_i$  and its vtrace counterpart  $\hat{\pi}_i$  have the same secrets for  $i \in \{1, 2\}$ .

We must also factor in the polarities of these relationships. First, as the traces evolve during the unwinding game, relationship (R1) will be *assumed*, whereas relationships (R2) and (R3) must be *guaranteed*. Moreover, assuming/guaranteeing that two generated sequences *stay equal* (like for the sequences of actions in relationships (R1)–(R3)) has a different flavour from assuming/guaranteeing that two generated sequences *become different* (like for the sequences

of observations in relationships (R1) and (R2))—roughly, we are talking about a safety versus a liveness property.

To capture these nuances, we will maintain the status of the observations’ divergence at (R1) and (R2), by remembering whether  $\pi_1$  and  $\pi_2$  have (already) differed in their observations (meaning the status is “diff”) or not (yet), meaning the status is “eq” (and similarly for  $\hat{\pi}_1$  and  $\hat{\pi}_2$ ). Thus, we can guarantee (R2) when assuming (R1) by:

- starting with the status of  $\pi_1$  vs.  $\pi_2$ , as well as that of  $\hat{\pi}_1$  vs.  $\hat{\pi}_2$ , set to eq;
- making sure the status of  $\hat{\pi}_1$  vs.  $\hat{\pi}_2$  changes (from eq to diff) as soon as the status for  $\pi_1$  vs.  $\pi_2$  changes.

The vtraces  $\hat{\pi}_1$  and  $\hat{\pi}_2$  will grow as  $\pi_1$  and  $\pi_2$  grow, using one of the following transition-matching mechanisms:

- $\hat{\pi}_i$  takes a step to match a step by  $\pi_i$ ; or
- both  $\hat{\pi}_1$  and  $\hat{\pi}_2$  take a synchronized step to match a synchronized step by  $\pi_1$  and  $\pi_2$ ; or
- $\hat{\pi}_1$  and/or  $\hat{\pi}_2$  ignore the step taken either separately or synchronously by  $\pi_1$  and/or  $\pi_2$ .

These matching patterns are the most natural choices for the reactive growth of  $\hat{\pi}_i$ ’s based on that of the  $\pi_i$ ’s. But for extra flexibility in proofs we will also allow variations of these patterns—e.g.,  $\hat{\pi}_1$  alone taking a step in response to a synchronized step by  $\pi_1$  and  $\pi_2$ , or conversely  $\hat{\pi}_1$  and  $\hat{\pi}_2$  taking a synchronized step in response to a step by  $\pi_1$ . The side-conditions when applying these matching mechanisms ensure that we maintain “(R1) implies (R2) and (R3)”.

In what follows, we will refer to the conditions (R1)–(R3) more casually and intuitively, as:

- an *interaction contract*, corresponding to “(R1) implies (R2)”, about the actions of  $\hat{\pi}_1$  and  $\hat{\pi}_2$  being the same and the observations being (eventually) different (the latter provided this is the case for  $\pi_1$  and  $\pi_2$ );
- a *secrecy contract*, corresponding to “((R1) implies (R2))”, about the produced secrets being the same between  $\hat{\pi}_i$  and  $\pi_i$ .

In addition to the above discussed *reactive* growth of the  $\hat{\pi}_i$ ’s, we will also allow their *proactive* growth. This will ensure further flexibility in proofs by enabling the  $\hat{\pi}_i$ ’s to take “independent” moves, i.e., moves not triggered by moves of the  $\pi_i$ ’s. However, to ensure soundness we will need to restrict proactive growth using *timers* (as we will explain).

## 5.2. Definition of unwinding

Below is the formal definition. We let  $\text{Status} = \{\text{eq}, \text{diff}\}$ , where eq signifies equality and diff difference/divergence, and  $\mathbb{N}_\infty = \mathbb{N} \cup \{\infty\}$ , the set of natural numbers extended with  $\infty$ —these will be our timer parameters. We let  $\text{SState}_u = \text{State}_u \times \text{State}_u \times \text{Status}$ , where  $u \in \{\text{opt}, \text{van}\}$ .

**Def. 3** A relation  $\Delta : \mathbb{N}_\infty \rightarrow (\mathbb{N}_\infty \times \mathbb{N}_\infty) \rightarrow \text{SState}_{\text{opt}} \rightarrow \text{SState}_{\text{van}} \rightarrow \text{Bool}$  is an *unwinding* when, for all  $v \in \mathbb{N}_\infty$ ,  $v_1, v_2 \in \mathbb{N}_\infty$ ,  $s_1, s_2 \in \text{State}_{\text{opt}}$ ,  $\hat{s}_1, \hat{s}_2 \in \text{State}_{\text{van}}$  and  $st, \hat{st} \in \text{Status}$ , if  $\Delta v(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st})$  then:

- $st = \text{eq}$  implies  $\text{isInt}(s_1) \leftrightarrow \text{isInt}(s_2)$ ;

- $\text{final}(s_1) \leftrightarrow \text{final}(s_2) \leftrightarrow \text{final}(\hat{s}_1) \leftrightarrow \text{final}(\hat{s}_2)$ ;
- either  $\text{react}(\Delta)(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st})$ ,  
or  $\exists w < v. \text{proact}(\Delta) w(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st})$ .

The predicates  $\text{react}(\Delta)$  and  $\text{proact}(\Delta)$ , covering the aforementioned reactive and proactive components of unwinding, are defined in Figs. 6 and 7, respectively. The definitions in these figures use the parameters of relative security either directly (e.g.,  $\text{isSec}$  and  $\text{isInt}$ ) or via the following auxiliary functions  $\text{newStat}$ ,  $\text{eqSec}$  and  $\text{eqAct}$ :

$$\begin{aligned} \text{eqSec}(s, s') &\doteq (\text{isSec}(s) \leftrightarrow \text{isSec}(s')) \wedge \\ &\quad (\text{isSec}(s) \rightarrow \text{getSec}(s) = \text{getSec}(s')) \\ \text{eqAct}(s, s') &\doteq (\text{isInt}(s) \leftrightarrow \text{isInt}(s')) \wedge \\ &\quad (\text{isInt}(s) \rightarrow \text{getAct}(s) = \text{getAct}(s')) \\ \text{newStat}(st, s, s') &\doteq \begin{cases} \text{diff} & \text{if } \text{isInt}(s) \wedge \text{isInt}(s') \wedge \\ & \text{getObs}(s) \neq \text{getObs}(s') \\ st & \text{otherwise} \end{cases} \end{aligned}$$

Thus, the predicate  $\text{eqSec}(s, s')$  holds when the states  $s$  and  $s'$  have equal secrets, if any; and similarly for  $\text{eqAct}$  concerning actions. The function  $\text{newStat}(st, s, s')$  tracks any change to the status  $st$ , from eq to diff, that may have occurred due to the observations in  $s$  and  $s'$ .

**Def. 4** A relation  $\Delta : \mathbb{N}_\infty \rightarrow \text{SState}_{\text{opt}} \rightarrow \text{SState}_{\text{van}} \rightarrow \text{Bool}$  is said to be a *finitary unwinding* if the condition from Def. 3 holds when ignoring the two  $\mathbb{N}_\infty$  arguments highlighted in gray (including in Figs. 6 and 7).

**On finitary unwinding.** We first discuss finitary unwinding (Def. 4), ignoring everything highlighted in gray—this will be enough for finitary relative security. A finitary unwinding relation  $\Delta$  encodes the current states of the otraces,  $s_1$  and  $s_2$ , and their observation divergence status  $st$  (call them “ostates” and “ostatus”), along with the current states  $\hat{s}_1$  and  $\hat{s}_2$  and status  $\hat{st}$  of the vtraces being constructed (“vstates” and “vstatus”), and a timer parameter  $v \in \mathbb{N}_\infty$  (for bounding proactive growth, explained later), expressed as  $\Delta v(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st})$ . Def. 4 ensures that  $\Delta$  supports the secure growth of the vtraces, reactively or proactively.

A note on notation: In Figs. 6 and 7, we use superscripts to indicate which subset of the two ostates take transitions, and subscripts to indicate which of the two vstates are reacting. For example,  $\text{match}_2^{1,2}(\Delta)$  means that  $s_1$  and  $s_2$  both take transitions, and  $\hat{s}_2$  takes a matching transition in reaction. Superscripts refer to *demonic nondeterminism*: we must consider all three cases; thus,  $\text{react}(\Delta)$  is the *conjunction* of  $\text{match}_1^1(\Delta)$ ,  $\text{match}_2^2(\Delta)$  and  $\text{match}_1^{1,2}(\Delta)$ . By contrast, subscripts refer to *angelic nondeterminism*; we are free to choose between different ways to react: either ignore, or match with one counterpart, or match with both counterparts; thus, e.g.,  $\text{match}_1^{1,2}(\Delta)$  is a (conditional) *disjunction* involving  $\text{match}_1^{1,2}(\Delta)$ ,  $\text{match}_2^{1,2}(\Delta)$  and  $\text{match}_{1,2}^{1,2}(\Delta)$ .

**Reactive growth.** Predicate  $\text{react}(\Delta)$  (Fig. 6) describes how vstates  $\hat{s}_1$  and  $\hat{s}_2$  can match transitions of ostates  $s_1$  and  $s_2$ :

- either separately, via predicates  $\text{match}_1^1(\Delta)$  or  $\text{match}_2^2(\Delta)$ , in case one of the ostates transits without any interaction (i.e.,  $\neg \text{isInt}(s_1)$  or  $\neg \text{isInt}(s_2)$ ),

$$\begin{aligned}
\text{react}(\Delta)(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \text{match}^1(\Delta)(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) \wedge \\
&\quad \text{match}^2(\Delta)(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) \wedge \\
&\quad \text{match}^{1,2}(\Delta)(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) \\
\text{match}^1(\Delta)(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \\
\neg \text{isInt}(s_1) \longrightarrow & \\
(\forall s'_1. s_1 \Rightarrow s'_1 \longrightarrow & \\
(\exists (w_1, w_2) < (v_1, v_2). \neg \text{isSec}(s_1) \wedge \Delta \infty (w_1, w_2)(s'_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st})) \vee & \\
(\exists w_2 < v_2. \text{eqSec}(s_1, \hat{s}_1) \wedge \neg \text{isInt}(\hat{s}_1) \wedge \text{match}_1^1(\Delta)(\infty, w_2)(s'_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st})) \vee & \\
(\text{eqSec}(s_1, \hat{s}_1) \wedge \text{eqAct}(\hat{s}_1, \hat{s}_2) \wedge \neg \text{isSec}(\hat{s}_2) \wedge \text{match}_{1,2}^1(\Delta)(\infty, \infty)(s'_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}))) & \\
\text{match}^2(\Delta)(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \\
\neg \text{isInt}(s_2) \longrightarrow & \\
(\forall s'_2. s_2 \Rightarrow s'_2 \longrightarrow & \\
(\exists (w_1, w_2) < (v_1, v_2). \neg \text{isSec}(s_2) \wedge \Delta \infty (w_1, w_2)(s_1, s'_2, st)(\hat{s}_1, \hat{s}_2, \hat{st})) \vee & \\
(\exists w_1 < v_1. \text{eqSec}(s_2, \hat{s}_2) \wedge \neg \text{isInt}(\hat{s}_2) \wedge \text{match}_2^2(\Delta)(w_1, \infty)(s_1, s'_2, st)(\hat{s}_1, \hat{s}_2, \hat{st})) \vee & \\
(\text{eqSec}(s_2, \hat{s}_2) \wedge \text{eqAct}(\hat{s}_1, \hat{s}_2) \wedge \neg \text{isSec}(\hat{s}_1) \wedge \text{match}_{1,2}^2(\Delta)(\infty, \infty)(s_1, s'_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}))) & \\
\text{match}^{1,2}(\Delta)(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \\
\text{let } st' = \text{newStat}(st, s_1, s_2) \text{ in} & \\
\text{isInt}(s_1) \wedge \text{isInt}(s_2) \wedge \text{eqAct}(s_1, s_2) \longrightarrow & \\
(\forall s'_1, s'_2. s_1 \Rightarrow s'_1 \wedge s_2 \Rightarrow s'_2 \longrightarrow & \\
(\exists (w_1, w_2) < (v_1, v_2). \neg \text{isSec}(s_1) \wedge \neg \text{isSec}(s_2) \wedge (st' = st \vee \hat{st} = \text{diff}) \wedge \Delta \infty (w_1, w_2)(s'_1, s'_2, st')(\hat{s}_1, \hat{s}_2, \hat{st})) \vee & \\
(\exists w_2 < v_2. \text{eqSec}(s_1, \hat{s}_1) \wedge \neg \text{isSec}(s_2) \wedge \neg \text{isInt}(\hat{s}_1) \wedge (st' = st \vee \hat{st} = \text{diff}) \wedge \text{match}_1^{1,2}(\Delta)(\infty, w_2)(s'_1, s'_2, st')(\hat{s}_1, \hat{s}_2, \hat{st})) \vee & \\
(\exists w_1 < v_1. \neg \text{isSec}(s_1) \wedge \text{eqSec}(s_2, \hat{s}_2) \wedge \neg \text{isInt}(\hat{s}_2) \wedge (st' = st \vee \hat{st} = \text{diff}) \wedge \text{match}_2^{1,2}(\Delta)(w_1, \infty)(s'_1, s'_2, st')(\hat{s}_1, \hat{s}_2, \hat{st})) \vee & \\
(\text{eqSec}(s_1, \hat{s}_1) \wedge \text{eqSec}(s_2, \hat{s}_2) \wedge \text{eqAct}(\hat{s}_1, \hat{s}_2) \wedge \text{match}_{1,2}^{1,2}(\Delta)(\infty, \infty)(s'_1, s'_2, st')(\hat{s}_1, \hat{s}_2, \hat{st}))) & \\
\text{match}_1^1(\Delta)(v_1, v_2)(s'_1, s'_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \exists \hat{s}'_1 \in \text{State}_{\text{van.}}. \hat{s}_1 \Rightarrow \hat{s}'_1 \wedge \Delta \infty (v_1, v_2)(s'_1, s_2, st)(\hat{s}'_1, \hat{s}_2, \hat{st}) \\
\text{match}_{1,2}^1(\Delta)(v_1, v_2)(s'_1, s'_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \exists \hat{s}'_1, \hat{s}'_2 \in \text{State}_{\text{van.}}. \hat{s}_1 \Rightarrow \hat{s}'_1 \wedge \hat{s}_2 \Rightarrow \hat{s}'_2 \wedge \Delta \infty (v_1, v_2)(s'_1, s'_2, st)(\hat{s}'_1, \hat{s}'_2, \hat{st}) \\
\text{match}_2^2(\Delta)(v_1, v_2)(s_1, s'_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \exists \hat{s}'_2 \in \text{State}_{\text{van.}}. \hat{s}_2 \Rightarrow \hat{s}'_2 \wedge \Delta \infty (v_1, v_2)(s_1, s'_2, st)(\hat{s}_1, \hat{s}'_2, \hat{st}) \\
\text{match}_{1,2}^2(\Delta)(v_1, v_2)(s_1, s'_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \exists \hat{s}'_1, \hat{s}'_2 \in \text{State}_{\text{van.}}. \hat{s}_1 \Rightarrow \hat{s}'_1 \wedge \hat{s}_2 \Rightarrow \hat{s}'_2 \wedge \Delta \infty (v_1, v_2)(s_1, s'_2, st)(\hat{s}'_1, \hat{s}'_2, \hat{st}) \\
\text{match}_1^{1,2}(\Delta)(v_1, v_2)(s'_1, s'_2, st')(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \exists \hat{s}'_1 \in \text{State}_{\text{van.}}. \hat{s}_1 \Rightarrow \hat{s}'_1 \wedge \Delta \infty (v_1, v_2)(s'_1, s'_2, st')(\hat{s}'_1, \hat{s}_2, \hat{st}) \\
\text{match}_2^{1,2}(\Delta)(v_1, v_2)(s'_1, s'_2, st')(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \exists \hat{s}'_2 \in \text{State}_{\text{van.}}. \hat{s}_2 \Rightarrow \hat{s}'_2 \wedge \Delta \infty (v_1, v_2)(s'_1, s'_2, st')(\hat{s}_1, \hat{s}'_2, \hat{st}) \\
\text{match}_{1,2}^{1,2}(\Delta)(v_1, v_2)(s'_1, s'_2, st')(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \\
\text{let } \hat{st}' = \text{newStat}(\hat{st}, \hat{s}_1, \hat{s}_2) \text{ in} & \\
\exists \hat{s}'_1, \hat{s}'_2 \in \text{State}_{\text{van.}}. \hat{s}_1 \Rightarrow \hat{s}'_1 \wedge \hat{s}_2 \Rightarrow \hat{s}'_2 \wedge (st' = \text{diff} \longrightarrow \hat{st}' = \text{diff}) \wedge \Delta \infty (v_1, v_2)(s'_1, s'_2, st')(\hat{s}'_1, \hat{s}'_2, \hat{st}') &
\end{aligned}$$

Fig. 6: Definition of  $\text{react}(\Delta)$  – the reactive component of unwinding

$$\begin{aligned}
\text{proact}(\Delta) v(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \\
\neg \text{isSec}(s_1) \wedge \neg \text{isInt}(s_1) \wedge \text{imove}_1(\Delta) v(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) \vee & \\
\neg \text{isSec}(s_1) \wedge \neg \text{isInt}(s_1) \wedge \text{imove}_2(\Delta) v(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) \vee & \\
\neg \text{isSec}(s_1) \wedge \neg \text{isSec}(s_2) \wedge \text{eqAct}(s_1, s_2) \wedge \text{imove}_{1,2}(\Delta) v(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) & \\
\text{imove}_1(\Delta) v(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \exists \hat{s}'_1 \in \text{State}_{\text{van.}}. \hat{s}_1 \Rightarrow \hat{s}'_1 \wedge \Delta v(v_1, v_2)(s_1, s_2, st)(\hat{s}'_1, \hat{s}_2, \hat{st}) \\
\text{imove}_2(\Delta) v(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \exists \hat{s}'_2 \in \text{State}_{\text{van.}}. \hat{s}_2 \Rightarrow \hat{s}'_2 \wedge \Delta v(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}'_2, \hat{st}) \\
\text{imove}_{1,2}(\Delta) v(v_1, v_2)(s_1, s_2, st)(\hat{s}_1, \hat{s}_2, \hat{st}) &\doteq \text{let } \hat{st}' = \text{newStat}(\hat{st}, \hat{s}_1, \hat{s}_2) \text{ in} \\
&\quad \exists \hat{s}'_1, \hat{s}'_2 \in \text{State}_{\text{van.}}. \hat{s}_1 \Rightarrow \hat{s}'_1 \wedge \hat{s}_2 \Rightarrow \hat{s}'_2 \wedge \Delta v(v_1, v_2)(s_1, s_2, st)(\hat{s}'_1, \hat{s}'_2, \hat{st}')
\end{aligned}$$

Fig. 7: Definition of  $\text{proact}(\Delta)$  – the proactive component of unwinding

- or synchronously, via  $\text{match}^{1,2}(\Delta)$ , in case both ostates transit interactively (i.e.,  $\text{isInt}(s_1)$  and  $\text{isInt}(s_2)$ ).

There is ample flexibility when choosing the matching transitions. For example, let us detail the case of  $\text{match}^1(\Delta)$ ; those of  $\text{match}^2(\Delta)$  and  $\text{match}^{1,2}(\Delta)$  are similar. In the definition of  $\text{match}^1(\Delta)$ , we are free to match a transition  $s_1 \Rightarrow s'_1$  in one of three ways, as reflected by the disjunct on the righthand side of the implication from  $s_1 \Rightarrow s'_1$ :

- 1) either ignoring the transition, but this only provided no secret is produced by the target ( $\neg \text{isSec}(s_1)$ )—otherwise the vtrace would have been forced to also produce a secret to avoid breaching the secrecy contract;
- 2) or giving a matching transition by the counterpart vtrace,  $\hat{s}_1 \Rightarrow \hat{s}'_1$ , via  $\text{match}_1^1(\Delta)$ , provided the same secret (if any) is produced (i.e.,  $\text{eqSec}(s_1, \hat{s}_1)$ ), and no vtrace interaction happens (i.e.,  $\neg \text{isInt}(\hat{s}_1)$ )—to ensure that the vtrace does not breach the interaction contract;
- 3) or giving matching transitions by both vtraces,  $\hat{s}_1 \Rightarrow \hat{s}'_1$  and  $\hat{s}_2 \Rightarrow \hat{s}'_2$ , via  $\text{match}_{1,2}^1(\Delta)$ , provided  $\hat{s}_1$  respects the secrecy contract towards its counterpart (i.e.,  $\text{eqSec}(s_1, \hat{s}_1)$ ),  $\hat{s}_1$  and  $\hat{s}_2$  respect the interaction contract towards each other (i.e.,  $\text{eqAct}(\hat{s}_1, \hat{s}_2)$ ), and  $\hat{s}_2$  doesn't produce a secret (not to breach the secrecy contract).

The predicate  $\text{match}^{1,2}$  is different from  $\text{match}^1$  and  $\text{match}^2$  in that it allows interaction ( $\text{isInt}(s_1)$  and  $\text{isInt}(s_2)$ ), which can lead to a change in observation divergence status ( $st' = \text{newStat}(st, s_1, s_2)$ ), meaning that the otraces might *right now* diverge observationally. In this case, relative security requires that the vtraces also produce different observations, which is reflected in our unwinding condition. Indeed, according to the definition of  $\text{match}^{1,2}(\Delta)$ , if the status has not changed, *or* divergence had already been recorded before in the vtraces (i.e.,  $st' = st \vee \hat{st} = \text{diff}$ ), then one is allowed to react by either ignoring the transition or matching it with only one of the vtraces. However, if the status has changed then one must react with both vtraces transitioning, i.e., take the fourth,  $\text{match}_{1,2}^{1,2}(\Delta)$  option, and make sure that that vtrace divergence status has become diff in case the otraces status has ( $st' = \text{diff} \rightarrow \hat{st}' = \text{diff}$ )—as seen in the definition of  $\text{match}_{1,2}^{1,2}(\Delta)$ .

**Proactive growth.** Predicate  $\text{proact}(\Delta)$  (see Fig. 7) allows the vtraces to take transitions independently. Thus, proactive moves represent “extra help” when proving relative security. They have similar conditions as the reactive moves, but simpler since they involve only the vstates, not the ostates. Proactive, “independent” moves can be taken either by one vstate ( $\text{imove}_1(\Delta)$  or  $\text{imove}_2(\Delta)$ ), or by both ( $\text{imove}_{1,2}(\Delta)$ ), subject to restrictions of not producing secrets (to avoid breaching the secrecy contract) and producing the same action if any (or not at all if moving separately). In the case of mutual independent moves,  $\text{imove}_{1,2}(\Delta)$ , a possible change of observation status can occur—which is being recorded because, in a presumptive proof of unwinding, it is helpful to know if the vtrace observations have already diverged, so that in the future the proof no longer has to be “on watch” for the otraces to diverge (so to have the vtraces diverge at the same time, as discussed above for  $\text{match}_{1,2}^{1,2}(\Delta)$ ).

Another aspect that distinguishes proactive moves from reactive moves is that the former do not advance the otraces, which means that they should not be allowed to proceed indefinitely thus “filibustering” the unwinding game. Indeed, in that case, unwinding would fail to ensure relative security, because the otraces may end up not being entirely processed, rendering the secrecy-contract conditions  $S(\hat{\pi}_i) = S(\pi_i)$  uncertain. For this reason, we carry an additional timer parameter  $v \in \mathbb{N}_\infty$  that **1**) is forced to decrease each time we take a proactive move (as seen in Def. 3 when passing to  $\text{proact}(\Delta)$  a value  $w < v$ ) and **2**) is reset to  $\infty$  when we take a reactive move, and stays  $\infty$  during reactive moves. This ensures that the proactive moves cannot be taken continuously and infinitely, but eventually yield to reactive moves.

**Upgrading to general unwinding.** In the presence of infinite traces, there are more opportunities for “filibustering” in addition to the above discussed one coming from unbounded proactive moves (by the vtraces), namely from the unbalanced reactive moves where one of the vtraces  $\hat{\pi}_i$  grows indefinitely and “starves” the other one. This would allow us to successfully play the unwinding game without proving relative security, more precisely without proving that the constructed  $\hat{\pi}_1$  and  $\hat{\pi}_2$  have the same actions. To counter this, we ensure that either vtrace always eventually yields to the other one when taking reactive moves (unless that other one has terminated). This is the purpose of the additional timers  $v_1$  and  $v_2$  highlighted in Def. 3 and Figs. 6 and 7. Namely, we think of  $v_i$  as counting the time until  $\hat{\pi}_i$  will make progress in a reactive move, via matching—indeed, in Fig. 6 we see how  $v_i$  is decreased each time a matching action is taken without  $\hat{s}_i$  taking a transition, and is reset to  $\infty$  as soon as  $\hat{s}_i$  takes a transition.

### 5.3. Soundness of the unwinding proof method

Unwinding seeks to provide sufficient local (state-based) conditions ensuring relative security. This is indeed the case:

**Thm. 5 (The Proof Unwinding Theorem)** Assume that:

- $\Delta$  is a (finitary) unwinding relation.
- $\Delta \infty (\infty, \infty)$  (resp.  $\Delta \infty$ ) covers the initial states, i.e.: for all  $s_1, s_2 \in \text{State}_{\text{opt}}$  such that  $\text{istate}(s_1)$  and  $\text{istate}(s_2)$ , there exist  $\hat{s}_1, \hat{s}_2 \in \text{State}$  such that  $\text{istate}(\hat{s}_1)$  and  $\Delta \infty (\infty, \infty) (s_1, s_2, \text{eq}) (\hat{s}_1, \hat{s}_2, \text{eq})$  (resp.  $\Delta \infty (s_1, s_2, \text{eq}) (\hat{s}_1, \hat{s}_2, \text{eq})$ ).

Then (finitary) relative security holds, i.e.,  $(\mathcal{SM}_{\text{opt}}, \mathcal{AM}_{\text{opt}}) \geq_{\checkmark} (\mathcal{SM}_{\text{van}}, \mathcal{AM}_{\text{van}})$  (resp.  $(\mathcal{SM}_{\text{opt}}, \mathcal{AM}_{\text{opt}}) \geq_{\text{fin}} (\mathcal{SM}_{\text{van}}, \mathcal{AM}_{\text{van}})$ ).

The sound proof method enabled by Thm. 5 is the following: To prove that relative security holds, it suffices to provide an unwinding relation that covers the initial states.

Whereas relative security and finitary relative security are incomparable (neither implies the other), general (non-finitary) unwinding implies finitary unwinding (by ignoring the  $v_i$  timers). So the general unwinding proof method ensures both relative security and finitary relative security.

## 5.4. Compositional, distributed unwinding

In practice, one might prefer to not define a single unwinding relation, i.e., a relation  $\Delta$  shown to “unwind into itself”, but a network of relations that unwind into each other:

**Def. 6** An *unwinding network* is a tuple  $(n, \text{next}, \text{Init}, \overline{\Delta})$  where  $n \in \mathbb{N}$ ,  $\text{Init} \subseteq \{0, \dots, n\}$ ,  $\text{next} : \{0, \dots, n\} \rightarrow \mathcal{P}(\{0, \dots, n\})$  and  $\overline{\Delta} = (\Delta_i)_{i \in \{0, \dots, n\}}$  are such that, for all  $i \in \{0, \dots, n\}$  and  $v, v_1, v_2, s_1, s_2, st, \hat{s}_1, \hat{s}_2, \hat{st}$ , if  $\Delta_i v (v_1, v_2) (s_1, s_2, st) (\hat{s}_1, \hat{s}_2, \hat{st})$  holds, then the conditions from Def. 3 hold but for  $\Delta$  replaced with  $\bigvee_{j \in \text{next}(i)} \Delta_j$ .

This definition generalizes Def. 3 by allowing each  $\Delta_i$  to take reactive and proactive steps unwinding into any  $\Delta_j$  to which  $\Delta_i$  is connected, as described by the next operator.  $\text{Init}$  stands for set of the initial nodes in this network.

**Thm. 7 (The Proof Distributed-Unwinding Theorem)** Assume that  $(n, \text{next}, \overline{\Delta}, M)$  is an unwinding network and  $\bigcup_{i \in \text{Init}} (\Delta_i \infty)$  covers the initial states (as in Thm. 5). Then  $(SM_{\text{opt}}, AM_{\text{opt}}) \geq \checkmark (SM_{\text{van}}, AM_{\text{van}})$ .

(And a similar generalization holds for finitary unwinding.)

Thm. 7 was useful in our verification case studies (to be described in §6), where different unwinding components turned out to naturally correspond to the different phases that the considered programs’ executions go through.

## 5.5. Disproof Method for Relative Security

A counterexample for relative security requires **1**) a concrete step, providing a pair  $(\pi_1, \pi_2)$  of otraces that have the same actions and different observations, and **2**) an abstract reasoning step, showing that there is no similarly related pair  $(\hat{\pi}_1, \hat{\pi}_2)$  of vtraces producing the same secrets as  $(\pi_1, \pi_2)$ . For handling step 2, we can again use an unwinding-like technique. After collecting the sequences of secrets  $(\sigma_1, \sigma_2)$  of  $(\pi_1, \pi_2)$ , we make sure that no suitable vtraces  $(\hat{\pi}_1, \hat{\pi}_2)$  can cover these secrets. This is done by maintaining a four-place unwinding relation containing pairs  $(s_1, \sigma_1)$  and  $(s_2, \sigma_2)$ , where each  $s_i$  is the state currently reached by the presumptive trace  $\hat{\pi}_i$  and  $\sigma_i$  is the sequence of secrets still left to be covered by  $\hat{\pi}_i$ . We will state unwinding conditions ensuring that no vtraces starting in  $s_1$  and  $s_2$  can achieve both same actions and different observations if they are to stay on track, i.e., cover the given secrets  $\sigma_1$  and  $\sigma_2$ —in this sense, our unwinding relations will be secret-directed.

**Def. 8** A relation  $\Gamma : (\text{State}_{\text{van}} \times \text{Seq}(\text{Sec})) \rightarrow (\text{State}_{\text{van}} \times \text{Seq}(\text{Sec})) \rightarrow \text{Bool}$  is a *secret-directed (SD) unwinding* when for all  $s_1, \sigma_1, s_2, \sigma_2$ , if  $\Gamma(s_1, \sigma_1)(s_2, \sigma_2)$  then:

- $\text{isInt}(s_1) \leftrightarrow \text{isInt}(s_2)$
- $\neg \text{isInt}(s_1) \rightarrow \text{move}_1(\Gamma)(s_1, \sigma_1)(s_2, \sigma_2) \wedge \text{move}_2(\Gamma)(s_1, \sigma_1)(s_2, \sigma_2)$
- $\text{isInt}(s_1) \wedge \text{getAct}(s_1) = \text{getAct}(s_2) \rightarrow \text{getObs}(s_1) = \text{getObs}(s_2) \wedge \text{move}_{1,2}(\Gamma)(s_1, \sigma_1)(s_2, \sigma_2)$

The predicates  $\text{move}_1(\Gamma)$ ,  $\text{move}_2(\Gamma)$  and  $\text{move}_{1,2}(\Gamma)$  are defined in Fig. 8. They are based on the secret-directed transition relation  $\Rightarrow_{\text{isSec}}^{\text{getSec}}$ , also defined in Fig. 8.

As seen in the above bullet points, an SD unwinding guarantees that the two states  $s_1$  and  $s_2$  (reached by the presumptive vtraces) always have the same interaction status, and equal observations if equal actions; also, the secret-directed transition relation  $\Rightarrow_{\text{isSec}}^{\text{getSec}}$  shown in Fig. 8 makes sure that the states  $s_1$  and  $s_2$  can evolve (i.e., the vtraces can grow) only by respecting the remaining sequences of secrets—i.e., only producing the next secret in the corresponding sequence, if at all producing a secret. And indeed, maintaining an SD unwinding while starting with the sequences of secrets  $(S(\pi_1), S(\pi_2))$  given by two concrete otraces  $(\pi_1, \pi_2)$ , constitutes a sound disproof method:

**Thm. 9 (The Disproof Unwinding Theorem)** Assuming:

- $\pi_1, \pi_2 \in \text{Trace}_{\text{opt}}$ ,  $A(\pi_1) = A(\pi_2)$  and  $O(\pi_1) \neq O(\pi_2)$ .
- $\Gamma$  is an SD unwinding.
- $\Gamma$  covers the initial states w.r.t.  $(S(\pi_1), S(\pi_2))$ , in that  $\Gamma(s_1, S(\pi_1))(s_2, S(\pi_2))$  holds for all  $s_1, s_2 \in \text{State}_{\text{van}}$  such that  $\text{istate}(s_1)$  and  $\text{istate}(s_2)$ .

Then relative security fails, i.e.,  $(SM_{\text{opt}}, AM_{\text{opt}}) \not\geq \checkmark (AM_{\text{van}}, LM_{\text{van}})$ . And if  $\pi_1, \pi_2 \in \text{Trace}_{\text{fin}}$ , then finitary relative security also fails, i.e.,  $(SM_{\text{opt}}, AM_{\text{opt}}) \not\geq \checkmark^{\text{fin}} (SM_{\text{van}}, AM_{\text{van}})$ .

## 6. Unwinding (Dis)Proofs for our Examples

**Unwinding proofs.** We deployed the distributed unwinding method (Thm. 7) to prove the relative security of (§3.2’s IMP representations of) `fun2–fun6`. The proofs turn the informal intuition for why these programs are (relatively) secure into an unwinding argument involving ostates  $s_1$  and  $s_2$  (those reached by two otraces  $\pi_1$  and  $\pi_2$ ) and vstates  $\hat{s}_1$  and  $\hat{s}_2$  (those reached by two vtraces  $\hat{\pi}_1$  and  $\hat{\pi}_2$ ).

In all the proofs, the two ostates  $s_1$  and  $s_2$  always have the same PC (i.e., execute the statements lock-step synchronously) at all levels; and the two vstates  $\hat{s}_1$  and  $\hat{s}_2$  always have the same PC too. Moreover, often (but not always) the vstates have the same PC with the level-0 PC of the ostates. Since interaction is pervasive, the  $\text{match}^1$  and  $\text{match}^2$  predicates are vacuously true, and when checking  $\text{match}^{1,2}$  we only choose “ignore” (first disjunct) or the  $\text{match}_{1,2}^{1,2}$  (fourth disjunct) options—this is because of the aforementioned lock-step synchronization. Another common aspect will be that  $s_1$  and  $s_2$  always have the same value for the input variable  $x$ , and we also make sure that  $\hat{s}_1$  and  $\hat{s}_2$  have the same value for  $x$ —this is to respect the interaction contract, namely the “action” part of this contract, which assumes that  $A(\pi_1) = A(\pi_2)$  and guarantees  $A(\hat{\pi}_1) = A(\hat{\pi}_2)$ . Often (but not always) this value of  $x$  is the same across the board, i.e., the same for  $\hat{s}_i$  as for  $s_i$ .

Note that  $s_1$  and  $s_2$  are allowed to differ in the initial secrets, i.e., in the values stored in the arrays  $a$  and  $b$ . According to the secrecy contract, we always keep  $\hat{s}_1$  have the same values in  $a$  and  $b$  as  $s_1$ , and similarly for  $\hat{s}_2$  versus  $s_2$ . Our unwinding relations (unwindings for short) need to ensure that, whenever an output statement gets executed, if the outputted value or the set of read locations differs for



$$\begin{aligned}
(s, \sigma) &\Rightarrow_{\text{isSec}}^{\text{getSec}} (s', \sigma') \doteq s \Rightarrow s' \wedge ((\neg \text{isSec}(s) \wedge \sigma = \sigma') \vee (\text{isSec}(s) \wedge \sigma = \text{getSec}(s) \cdot \sigma')) \\
\text{move}_1(\Gamma) (s_1, \sigma_1) (s_2, \sigma_2) &\doteq \forall s'_1, \sigma'_1. (s_1, \sigma_1) \Rightarrow_{\text{isSec}}^{\text{getSec}} (s'_1, \sigma'_1) \longrightarrow \Gamma (s'_1, \sigma'_1) (s_2, \sigma_2) \\
\text{move}_2(\Gamma) (s_1, \sigma_1) (s_2, \sigma_2) &\doteq \forall s'_2, \sigma'_2. (s_2, \sigma_2) \Rightarrow_{\text{isSec}}^{\text{getSec}} (s'_2, \sigma'_2) \longrightarrow \Gamma (s_1, \sigma_1) (s'_2, \sigma'_2) \\
\text{move}_{1,2}(\Gamma) (s_1, \sigma_1) (s_2, \sigma_2) &\doteq \forall s'_1, \sigma'_1, s'_2, \sigma'_2. (s_1, \sigma_1) \Rightarrow_{\text{isSec}}^{\text{getSec}} (s'_1, \sigma'_1) \wedge (s_2, \sigma_2) \Rightarrow_{\text{isSec}}^{\text{getSec}} (s'_2, \sigma'_2) \longrightarrow \Gamma (s'_1, \sigma'_1) (s'_2, \sigma'_2)
\end{aligned}$$

Fig. 8: The defining predicates for SD unwinding

$s_1$  versus  $s_2$  then it must also differ for  $\hat{s}_1$  versus  $\hat{s}_2$ —this is here the observation part of the interaction contract.

`fun2` is secure essentially because the early fence on line 4 immediately inhibits any (mis)speculation that goes on the “then” branch of the “if” conditional. This is reflected by our unwindings keeping all PCs synchronized and all values of  $x$  the same, which means that the interaction contract is maintained straightforwardly.

`fun3`’s proof is slightly more subtle: In the only interesting scenario, speculation incorrectly takes the “then” branch, and now the otraces access locations that the vtraces cannot access, namely the  $x$ ’th location of  $a$ . But since  $s_1$  and  $s_2$  have the same value for  $x$ , the otraces access the same locations (among each other), so they are not observationally “more different” than the vtraces are—and this is the invariant we keep in our unwindings, which again validates the interaction contract.

For `fun4`, assuming  $N > 0$ , we must use two unwinding strategies depending on whether the value of the input is  $< N$  or not. The first case is not problematic, because there otraces can only diverge from vtraces in an immediately harmless way, i.e., by (mis)speculating on the “else” branch. For the second case, of  $x \geq N$ , (mis)speculation can go on the “then” branch and cause observational “damage” to the otraces, which access different locations from  $b$  depending on the value of  $a[0]$ , on which they may differ. To counter this, we choose a different input for the vtraces, namely 0. This means that the vtraces also take the “then” branch, keeping their PCs synchronized with the level-1 PCs of the otraces. At some point, speculation will be abandoned and the otraces will go back to the correct branch, reaching the output statement (i.e.,  $s_1$  and  $s_2$  will be back at speculation level 0 and PC the one of the output statement). This is the moment when, in our unwinding,  $\hat{s}_1$  and  $\hat{s}_2$  will take one or two proactive actions (depending on where the PCs were on the “then” branch when the speculation was abandoned by  $s_1$  and  $s_2$ ), using the `proact` option of the unwinding with the `imove1,2` option to also finish the “then” branch and reach the output statement. After this, the ostates are again PC-synchronized with the vstates, and the output statement can be taken in lockstep—while, importantly,  $\hat{s}_1$  and  $\hat{s}_2$  differ observationally at least as much as  $s_1$  and  $s_2$  do.

Implicit above is that we use networks of unwindings (as in Thm. 7) instead of single unwindings. We have one unwinding component for each phase in the program execution (e.g., “before entering the if branch”, or “at the fence”) and we transit between components upon relevant events (e.g., “the if branch is taken”, or “speculation starts”).

Since `fun5` is essentially the while-iteration of `fun3` (where input statements are replaced by `scanf` and `return`

statements by `printf`), its proof is obtained from that of `fun5` by adding corresponding cycles in the unwinding network that proves `fun4`, and additionally handling the speculation on the outer loop. This shows the potential of unwinding proofs following compositionally the syntactic structure.

Finally, the difference between `fun6` and `fun5` is that, in addition to the untrusted input  $x$ , `fun6` also takes a trusted input  $y$  which is considered to be secret. Our unwinding proof for `fun5` can be adapted to `fun6`, noticing that the required secrecy contract (enforced via the unwinding’s `eqSec` predicate) still holds because, in the unwinding strategy, the `InputT y` statement is always executed synchronously by the otraces and vtraces.

A note on unwinding timer parameters: The “ $v$ ” timer (which bounds proactive moves) stays in the range  $\{\infty, 2, 1, 0\}$  since we never need to take more than 3 proactive moves in a row (and this is only when vtraces must “catch-up” with mispredicting otraces by finishing their branch). As for the “ $v_1, v_2$ ” timers (which bound idle reactive moves), the same range suffices, since the only time this is needed is when mispredicting moves are ignored.

**SD unwinding disproof.** The intuitive reason why `fun1` is insecure is that, for  $x \geq N$ , speculation can access the  $(a[x] * 512)$ ’th location of  $b$  whereas normal execution cannot. With our disproof method, this argument is formalized by choosing a value  $i \geq N$  (as input for  $x$ ) and two initial arrays  $a$  that differ on  $a[i]$  but nowhere else. These give two (singleton) sequences of secrets (the  $S(\pi_1)$  and  $S(\pi_2)$  in Thm. 9), and two otraces  $\pi_1$  and  $\pi_2$  that produce these secrets (i.e., start in memories with the respective arrays). Then we show that any two vtraces that produce these secrets (i.e., start in the two given memories) must be observationally equal (i.e., output the same value) provided they are action-equal (i.e., input the same value  $x$ ). Indeed, here a trivial secret-directed unwinding keeps the two potential vtraces completely synchronized; they can only take the “then” branch if  $x$  is smaller than  $N$ , which means that  $x$  is different from  $i$ , which further means that what the two vtraces print cannot depend on  $a[i]$  hence must be equal.

## 7. Isabelle Mechanization

We used Isabelle/HOL [26] to mechanize the relative security framework [9, 28] and the examples [8].

**Theory structure and sizes.** The theory structure of the entire mechanization is shown in Fig. 9. The theories formalizing the abstract framework are shown in the middle-right part of the figure, e.g., `Relative_Security`, `Unwinding` etc. The suffix “fin” refers to the finitary (finite-trace) version of the concepts. These theories consist of 11K LOC, with

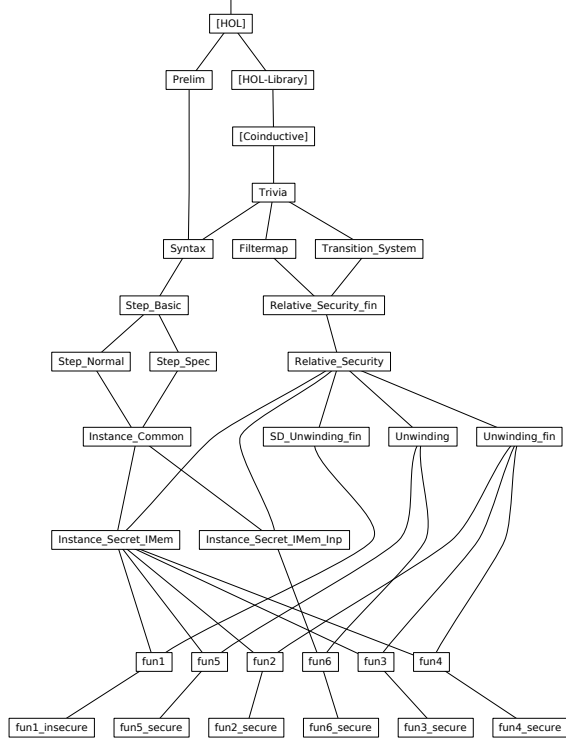


Fig. 9: Theory structure of our Isabelle mechanization

the largest part (7.5K LOC) taken by infinitary unwinding—which also required some substantial preparations (part of Trivia, 2.6K LOC) dedicated to extending the coinductive (lazy) list library with custom support for the inductive-coinductive mixture of reasoning called for by the timers.

The theories Syntax, Step\_Basic, Step\_Normal and Step\_Spec, taking 1.5K LOC, formalize IMP’s syntax and semantics straightforwardly, via datatypes and inductive predicates. The instantiation of relative security to IMP, taking 0.6K LOC, is formalized in theory Instance\_Common, covering the interaction infrastructure (which is common to all examples), and theories Instance\_Secret\_IMem and Instance\_Secret\_IMem\_Inp, covering the two secrecy infrastructures: one with secrets as the initial memories (for fun1–fun5), and one with additional secrets as trusted inputs and outputs (for fun6). We used finitary unwinding and SD-unwinding for verifying the programs that are known to be terminating, and used the heavier infinitary unwinding only for the possibly nonterminating fun5 and fun6. Verifying the (in)security of fun1–fun6 took an average of 2K LOC per example, which splits into 0.7K LOC for simplification-rules boilerplate (theories named fun(i)) and 1.3 LOC for the distributed unwinding proof (theories named fun(i)(in)secure). In the future, these large proof sizes (clearly unfeasible in practice) could be alleviated via custom tactics for automating the boilerplate.

**Locale structure.** Our mechanization relies heavily on Isabelle’s locales [17]. A locale fixes types, constants and assumptions, allowing one to perform definitions and proofs relative to these. A locale  $L$  can be *interpreted* at the top

level of an Isabelle theory by instantiating its types and constants, and verifying its assumptions. It can also be interpreted relatively to another (more concrete) locale  $L'$  by establishing a *sublocale* relationship  $L' \leq L$ .

We used (sub)locales to formalize our gradual move from abstract to concrete. We have a locale System\_Mod for system models, extended by a locale Leakage\_Mod that adds a type Leak and a predicate leakVia. Moreover, we have the locales Statewise\_Attacker\_Mod and Attacker\_Mod forming an abstract-towards-concrete sublocale hierarchy  $\text{Leakage\_Mod} < \text{Attacker\_Mod} < \text{Statewise\_Attacker\_Mod}$ .

Relative security is first defined abstractly by importing two (disjoint) copies of Leakage\_Mod (labeled “Van” and “Opt”) and then made more concrete along the above hierarchy, leading to the locale Rel\_Sec which imports two copies of Statewise\_Attacker\_Mod. It is in Rel\_Sec where §5’s unwinding (dis)proof methods are proved sound, following the ideas described in §5.1, §5.2, §5.5 when motivating (the soundness of) our unwinding conditions.

The instantiation of the locale Rel\_Sec, to IMP programs was performed by a sublocale relationship  $\text{Prog\_Mispred\_Init} < \text{Rel\_Sec}$ , where Prog\_Mispred\_Init is a locale that fixes a misprediction oracle (like in §3.1), a program and some initial memory, inputs and predictor state (like in §3.3). Finally, Prog\_Mispred\_Init is instantiated (via locale interpretation) to the fun1–fun6 examples; this makes available the unwinding (dis)proof methods from Rel\_Sec (via Prog\_Mispred\_Init), and the relative (in)security of these examples is proved along the lines sketched in §6.

## 8. Further Related Work

Our framework’s innovation compared to previous work is in regarding leaks as first-class citizens and allowing fully interactive attackers and dynamic secrets (see §4). Proof-theoretically, the framework generalizes the unwinding method [12], which was widely deployed in the security literature [21, 22, 25, 27, 29]. Our work is the first to develop an unwinding method not for proving the (absolute) security of a system, but for comparing the security of two systems (necessarily a four-trace property), and to propose unwinding for disproofs as well. This unwinding-based foundation could underpin automatic and compositional analysis/verification methods, such as type systems [2, 31, 33, 34].

Due to the need to consider alternative execution traces, proving information-flow security properties has similarities with proving program equivalence, where techniques involving symbolic execution and bisimulation were proposed [1, 10, 16]. Relative security seems related to Morgan’s refinement order for noninterference [23, 24]; working out the formal connection between the two notions could lead to more insights into compositional proofs for our notion.

**Acknowledgments.** We thank the reviewers for their excellent suggestions for improvement, and for catching several technical typos. We gratefully acknowledge support from the EPSRC grants EP/X015114/1 and EP/X015149/1, titled “Safe and secure COncurrent programming

for adVancEd aRchiTectures (COVERT)”. Dongol is additionally supported by EPSRC grants EP/Y036425/1, EP/X037142/1, EP/V038915/1, EP/R025134/2 and VeTSS. The authors are listed in alphabetical order.

## References

- [1] G. Barthe, J. M. Crespo, and C. Kunz, “Relational verification using product programs,” in *FM*, ser. LNCS, vol. 6664, 2011, pp. 200–214.
- [2] G. Boudol and I. Castellani, “Noninterference for concurrent programs and thread systems,” *Theor. Comp. Sci.*, vol. 281, no. 1-2, pp. 109–130, 2002.
- [3] C. Canella, J. V. Bulck, M. Schwarz, M. Lipp, B. von Berg, P. Ortner, F. Piessens, D. Evtushkin, and D. Gruss, “A systematic evaluation of transient execution attacks and defenses,” in *USENIX Sec.*, 2019, pp. 249–266.
- [4] S. Cauligi, C. Disselkoen, D. Moghimi, G. Barthe, and D. Stefan, “Sok: Practical foundations for software spectre defenses,” in *SP*. IEEE, 2022, pp. 666–680.
- [5] K. Cheang, C. Rasmussen, S. A. Seshia, and P. Subramanyan, “A formal approach to secure speculation,” in *CSF*. IEEE, 2019, pp. 288–303.
- [6] M. R. Clarkson and F. B. Schneider, “Hyperproperties,” *J. Com. Sec.*, vol. 18, no. 6, pp. 1157–1210, 2010.
- [7] B. Dongol, M. Griffin, A. Popescu, and J. Wright, Extended technical report associated to this paper. [https://www.andreipopescu.uk/pdf/TR\\_relative\\_security.pdf](https://www.andreipopescu.uk/pdf/TR_relative_security.pdf).
- [8] B. Dongol, M. Griffin, A. Popescu, and J. Wright, “Mechanization of relatively (in)secure example programs,” <https://www.andreipopescu.uk/form/RSEx.zip>.
- [9] B. Dongol, M. Griffin, A. Popescu, and J. Wright, “Secret-directed unwinding,” *Archive of Formal Proofs*, May 2024, [https://isa-afp.org/entries/Secret\\_Directed\\_Unwinding.html](https://isa-afp.org/entries/Secret_Directed_Unwinding.html), Formal proof development.
- [10] B. Godlin and O. Strichman, “Regression verification: proving the equivalence of similar programs,” *Softw. Test. Verif. Reliab.*, vol. 23, no. 3, pp. 241–258, 2013.
- [11] J. A. Goguen and J. Meseguer, “Security policies and security models,” in *SP*, 1982, pp. 11–20.
- [12] J. A. Goguen and J. Meseguer, “Unwinding and inference control,” in *SP*, 1984, pp. 75–87.
- [13] R. Guanciale, M. Balliu, and M. Dam, “Inspectre: Breaking and fixing microarchitectural vulnerabilities by formal analysis,” in *CCS*, 2020, pp. 1853–1869.
- [14] M. Guarnieri, B. Köpf, J. F. Morales, J. Reineke, and A. Sánchez, “Spectector: Principled detection of speculative information flows,” in *SP*. IEEE, 2020, pp. 1–19.
- [15] M. Guarnieri, B. Köpf, J. Reineke, and P. Vila, “Hardware-software contracts for secure speculation,” in *SP*. IEEE, 2021, pp. 1868–1883.
- [16] C. Hur, G. Neis, D. Dreyer, and V. Vafeiadis, “The power of parameterization in coinductive proof,” in *POPL*. ACM, 2013, pp. 193–206.
- [17] F. Kammüller, M. Wenzel, and L. C. Paulson, “Locales—a sectioning concept for Isabelle,” in *TPHOLs*, 1999, pp. 149–166.
- [18] P. Kocher, “Spectre mitigations in Microsoft’s C/C++ compiler,” 2018, <https://www.paulkocher.com/doc/MicrosoftCompilerSpectreMitigation.html>.
- [19] P. Kocher, J. Horn, A. Fogh, D. Genkin, D. Gruss, W. Haas, M. Hamburg, M. Lipp, S. Mangard, T. Prescher, M. Schwarz, and Y. Yarom, “Spectre attacks: Exploiting speculative execution,” in *SP*. IEEE, 2019, pp. 1–19.
- [20] M. Lipp, M. Schwarz, D. Gruss, T. Prescher, W. Haas, A. Fogh, J. Horn, S. Mangard, P. Kocher, D. Genkin, Y. Yarom, and M. Hamburg, “Meltdown: Reading kernel memory from user space,” in *USENIX Sec.*, 2018.
- [21] H. Mantel, “A uniform framework for the formal specification and verification of information flow security,” Ph.D. dissertation, University of Saarbrücken, 2003.
- [22] J. McLean, “A general theory of composition for trace sets closed under selective interleaving functions,” in *SP*, 1994, pp. 79–93.
- [23] C. Morgan, “The shadow knows: Refinement and security in sequential programs,” *Sci. Comput. Program.*, vol. 74, no. 8, pp. 629–653, 2009.
- [24] C. Morgan, “Compositional noninterference from first principles,” *Form. Asp. Comput.*, vol. 24, no. 1, pp. 3–26, 2012.
- [25] T. C. Murray, D. Matichuk, M. Brassil, P. Gammie, and G. Klein, “Noninterference for operating system kernels,” in *CPP*, 2012, pp. 126–142.
- [26] T. Nipkow, L. C. Paulson, and M. Wenzel, *Isabelle/HOL: A Proof Assistant for Higher-Order Logic*. Springer, 2002.
- [27] A. Popescu, P. Lammich, and T. Bauereiss, “CoCon: A confidentiality-verified conference management system,” *Arch. Formal Proofs*, 2021.
- [28] A. Popescu and J. Wright, “Relative security,” *Archive of Formal Proofs*, May 2024, [https://isa-afp.org/entries/Relative\\_Security.html](https://isa-afp.org/entries/Relative_Security.html), Formal proof development.
- [29] J. Rushby, “Noninterference, transitivity, and channel-control security policies,” SRI, Tech. Rep., Dec 1992, <http://www.csl.sri.com/papers/csl-92-2/>.
- [30] A. Sabelfeld and A. C. Myers, “Language-based information-flow security,” *IEEE Journal on Sel. Areas in Comm.*, vol. 21, no. 1, pp. 5–19, 2003.
- [31] A. Sabelfeld and D. Sands, “Probabilistic noninterference for multi-threaded programs,” in *IEEE Computer Security Foundations Workshop*, 2000, pp. 200–214.
- [32] D. Sangiorgi, “On the bisimulation proof method,” *Math. Struc. Com. Sci.*, vol. 8, no. 5, pp. 447–479, 1998.
- [33] B. A. Shivakumar, G. Barthe, B. Grégoire, V. Laporte, T. Oliveira, S. Priya, P. Schwabe, and L. Tabary-Maujean, “Typing high-speed cryptography against spectre v1,” in *SP*. IEEE, 2023, pp. 1094–1111.
- [34] D. Volpano and G. Smith, “Probabilistic noninterference in a concurrent language,” *Journal of Computer Security*, vol. 7, no. 2,3, pp. 231–253, 1999.
- [35] S. Zdancewic and A. C. Myers, “Observational determinism for concurrent program security,” in *IEEE CSF Workshop*, 2003, pp. 29–43.