# An Extended Account of Trace-Relating Compiler Correctness and Secure Compilation

CARMINE ABATE*, MPI-SP, Germany

ROBERTO BLANCO*, MPI-SP, Germany

ȘTEFAN CIOBÂCĂ, Alexandru Ioan Cuza University Iași, Romania

ADRIEN DURIER*, MPI-SP, Germany

DEEPAK GARG, Max Planck Institute for Software Systems, Germany

CĂTĂLIN HRIȚCU*, MPI-SP, Germany

MARCO PATRIGNANI, Stanford, USA and CISPA Helmholz Center for Information Security, Germany

ÉRIC TANTER*, University of Chile, Chile

JÉRÉMY THIBAULT*, MPI-SP, Germany

Compiler correctness, in its simplest form, is defined as the inclusion of the set of traces of the compiled program in the set of traces of the original program. This is equivalent to the preservation of all trace properties. Here, traces collect, for instance, the externally observable events of each execution. However, this definition requires the set of traces of the source and target languages to be the same, which is not the case when the languages are far apart or when observations are fine-grained. To overcome this issue, we study a generalized compiler correctness definition, which uses source and target traces drawn from potentially different sets and connected by an arbitrary relation. We set out to understand what guarantees this generalized compiler correctness definition gives us when instantiated with a non-trivial relation on traces. When this trace relation is not equality, it is no longer possible to preserve the trace properties of the source program unchanged. Instead, we provide a generic characterization of the target trace property ensured by correctly compiling a program that satisfies a given source property, and dually, of the source trace property one is required to show in order to obtain a certain target property for the compiled code. We show that this view on compiler correctness can naturally account for undefined behavior, resource exhaustion, different source and target values, side channels, and various abstraction mismatches. Finally, we show that the same generalization also

*Part of this work was conducted while these authors were employed at or visiting Inria Paris

applies to many definitions of *secure* compilation, which characterize the protection of a compiled program linked against adversarial code.

## 1   INTRODUCTION

Compiler correctness is an old idea [? ? ? ] that has seen a significant revival in recent times. This new wave was started by the creation of the CompCert verified C compiler [? ] and continued by the proposal of many significant extensions and variants of CompCert [? ? ? ? ? ? ? ? ? ? ? ] and the success of many other milestone compiler verification projects, including Vellvm [? ], Pilsner [? ], CakeML [? ], and CertiCoq [? ]. Verification through proof assistants allows the user of a compiler to trust the proofs without diving into all of the details. Still, in order to clearly understand the benefits and limitations of using a verified compiler, she has to deeply understand the statement of correctness. This is true not just for correct compilation, but also for secure compilation, which is the more recent idea that a compilation chain should not just provide correctness but also security against co-linked adversarial components [? ? ].

**Basic Compiler Correctness.** The gold standard for compiler correctness is *semantic preservation*, which intuitively says that the semantics of a compiled program (in the target language) is compatible with the semantics of the original program (in the source language). For practical verified compilers, such as CompCert [? ] and CakeML [? ], semantic preservation is stated extrinsically, by referring to *traces*. In these two settings, a trace is an ordered sequence of events—such as inputs from and outputs to an external environment—that are produced by the execution of a program.

A basic definition of compiler correctness can be given by the inclusion of the set of traces of the compiled program in the set of traces of the original program. Formally [? ]:

**Definition 1.1** (Basic Compiler Correctness (CC)).  A compiler ↓ is *correct* (CC) iff

$$\forall \mathsf{W} \; t. \; \mathsf{W}{\downarrow} \leadsto t \Rightarrow \mathsf{W} \leadsto t$$

This definition says that for any whole[1] source program $\mathsf{W}$, if we compile it (denoted $\mathsf{W}{\downarrow}$), execute it in the semantics of the target language, and observe a trace $t$, then the original $\mathsf{W}$ can produce *the same* trace $t$ in the semantics of the source language.[2] This definition is simple and easy to understand, since it only references a few familiar concepts: a compiler between a source and a target language, each equipped with a trace-producing semantics (usually nondeterministic).

**Beyond Basic Compiler Correctness.** Definition 1.1 implicitly assumes that the source and target traces are drawn from the very same set, and requires that any target trace produced by a compiled program can be faithfully reproduced by the source program. In practice, existing verified compiler adopt a less restrictive formulation of compiler correctness:

---

[1]For simplicity, for now we ignore separate compilation and linking, returning to it in §6.
[2]Typesetting convention [? ]: we use a blue, sans-serif font for source elements, an orange, **bold** font for **target** ones and a *black*, *italic* font for elements common to both languages.

**CompCert [? ]** The original compiler correctness theorem of CompCert [? ] can be seen as an instance of basic compiler correctness, but it does not provide any guarantees for programs that can exhibit undefined behavior [? ]. As allowed by the C standard, such unsafe programs are not even considered to be in the source language, so are not quantified over. This has important practical implications, since undefined behavior often leads to exploitable security vulnerabilities [? ? ? ] and serious confusion even among experienced C and C++ developers [? ? ? ? ]. As such, since 2010, CompCert provides an additional top-level correctness theorem[3] that better accounts for the presence of unsafe programs by providing guarantees for them up to the point when they encounter undefined behavior [? ]. This new theorem goes beyond the basic correctness definition above, as a target trace need only correspond to a source trace *up to the occurrence* of undefined behavior in the source trace.

**CakeML [? ]** Compiler correctness for CakeML accounts for memory exhaustion in target executions. Crucially, memory exhaustion events cannot occur in source traces, only in target traces. Hence, dually to CompCert, compiler correctness only requires source and target traces to coincide up to the occurrence of a memory exhaustion event in the target trace.

**Trace-Relating Compiler Correctness.** Generalized formalizations of compiler correctness like the ones above can be naturally expressed as instances of a uniform definition, which we call *trace-relating compiler correctness*. This generalizes basic compiler correctness by (a) considering that source and target traces belong to *possibly distinct* sets $\text{Trace}_S$ and $\text{Trace}_T$, and (b) being parameterized by an arbitrary *trace relation* $\sim$.

**Definition 1.2** (Trace-Relating Compiler Correctness (CC$^\sim$)). A compiler $\downarrow$ is *correct* with respect to a trace relation $\sim \subseteq \text{Trace}_S \times \text{Trace}_T$ iff

$$\forall W.\forall t.\ W{\downarrow}\rightsquigarrow t \Rightarrow \exists s \sim t.\ W\rightsquigarrow s$$

This definition requires that, for any target trace $t$ produced by the compiled program $W{\downarrow}$, there exist a source trace $s$ that can be produced by the original program $W$ and is *related* to $t$ according to $\sim$ (i.e., $s \sim t$). By choosing the trace relation appropriately, one can recover the different notions of compiler correctness presented above:

**Basic CC** Take $s \sim t$ to be $s = t$. Trivially, the basic CC of Definition 1.1 is CC$^=$.

**CompCert** Undefined behavior is modeled in CompCert as a trace-terminating event *Wrong* that can occur in any of its languages (source, target, and all intermediate languages), so for a given phase (or composition thereof), we have $\text{Trace}_S = \text{Trace}_T$. Nevertheless, the relation between source and target traces with which to instantiate CC$^\sim$ to obtain CompCert's current theorem is the following (note that we denote *finite* traces–or prefixes– as $m$)

$$s \sim t \equiv s = t \vee (\exists m \le t.\ s = m{\cdot}Wrong)$$

A compiler satisfying CC$^\sim$ for this trace relation can turn a source prefix ending in undefined behavior $m{\cdot}Wrong$ (where "$\cdot$" is concatenation) either into the same prefix in the target (first disjunct), or into a target trace that starts with the prefix $m$ but then continues *arbitrarily* (second disjunct, "$\le$" is the prefix relation).

**CakeML** Here, target traces are sequences of symbols from an alphabet $\Sigma_T$ that has a specific trace-terminating event, Resource_limit_hit, which is not available in the source alphabet $\Sigma_S$ (i.e., $\Sigma_T = \Sigma_S \cup \{\text{Resource\_limit\_hit}\}$). Then, the compiler correctness theorem of CakeML can be obtained by instantiating CC$^\sim$ with the following $\sim$ relation:

$$s \sim t \equiv s = t \vee (\exists m.\ m \le s.\ t = m{\cdot}\text{Resource\_limit\_hit})$$

---

[3]Stated at the top of the CompCert file `driver/Complements.v` and discussed by ? ].

The resulting $CC^\sim$ instance relates a target trace ending in Resource_limit_hit after executing prefix $m$ to a source trace that first produces $m$ and then continues in a way given by the semantics of the source program.

Beyond undefined behavior and resource exhaustion, there are many other practical uses for $CC^\sim$: in this paper we show that it also accounts for differences between source and target values, for a single source output being turned into a series of target outputs, and for side-channels.

On the flip side, the compiler correctness statement and its implications can be more difficult to understand for $CC^\sim$ than for $CC^=$. The full implications of choosing a particular $\sim$ relation can be subtle. In fact, using a bad relation can make the compiler correctness statement trivial or unexpected. For instance, it should be easy to see that if one uses the total relation, which relates all source traces to all target ones, the $CC^\sim$ property holds for every compiler, yet it might take one a bit more effort to understand that the same is true even for the following relation:

$$s \sim t \equiv \exists W. W \rightsquigarrow s \wedge W \downarrow \rightsquigarrow t$$

**Reasoning About Trace Properties.** To understand more about a particular $CC^\sim$ instance, we propose to also look at how it preserves *trace properties*—defined as sets of allowed traces [? ]—from the source to the target. For instance, it is well known that $CC^=$ is equivalent to the preservation of all trace properties (where $W \models \pi$ reads "$W$ satisfies property $\pi$" and stands for $\forall t. W \rightsquigarrow t \Rightarrow t \in \pi$):

$$CC^= \equiv \forall \pi \in 2^{\text{Trace}}. \forall W. W \models \pi \Rightarrow W \downarrow \models \pi$$

However, to the best of our knowledge, similar results have not been formulated for trace relations beyond equality, when it is no longer possible to preserve the trace properties of the source program unchanged. For trace-relating compiler correctness, where source and target traces can be drawn from different sets and related by an arbitrary trace relation, there are two crucial questions to ask:

(1) For a source trace property $\pi_S$ of a program—established for instance by formal verification—what is the strongest target property that any $CC^\sim$ compiler is guaranteed to ensure for the produced target program?

(2) For a target trace property $\pi_T$, what is the weakest source property we need to show of the original source program to obtain $\pi_T$ for the result of any $CC^\sim$ compiler?

Far from being mere hypothetical questions, they can help the developer of a verified compiler better understand the compiler correctness theorem they are proving, and we expect that any user of such a compiler will need to ask either one or the other if they are to make use of that theorem. In this work we provide a simple and natural answer to these questions, for any instance of $CC^\sim$. Building upon a bijection between relations and Galois connections [? ? ? ], we observe that any trace relation $\sim$ corresponds to two *property mappings* $\tilde{\tau}$ and $\tilde{\sigma}$, which are functions mapping source properties to target ones ($\tilde{\tau}$ standing for "to target") and target properties to source ones ($\tilde{\sigma}$ standing for "to source"):

$$\tilde{\tau}(\pi_S) = \{t \mid \exists s. s \sim t \wedge s \in \pi_S\}$$
$$\tilde{\sigma}(\pi_T) = \{s \mid \forall t. s \sim t \Rightarrow t \in \pi_T\}$$

The *existential image* of $\sim$, $\tilde{\tau}$, answers the first question above by mapping a given source property $\pi_S$ to the target property that contains all target traces for which *there exists a related source trace* that satisfies $\pi_S$. Dually, the *universal image* of $\sim$, $\tilde{\sigma}$, answers the second question by mapping a given target property $\pi_T$ to the source property that contains all source traces for which *all related target traces* satisfy $\pi_T$. We introduce two new correct compilation definitions in terms of *trace property preservation* (TP):

- $TP^{\tilde{\tau}}$ quantifies over all source trace properties and uses $\tilde{\tau}$ to obtain the corresponding target properties;

- $TP^{\tilde{\sigma}}$ quantifies over all target trace properties and uses $\tilde{\sigma}$ to obtain the corresponding source properties.

We prove that these two definitions are equivalent to $CC^{\sim}$, yielding a novel trinitarian view of compiler correctness (Figure 1).

$$\forall W. \ \forall t. \ W\downarrow\rightsquigarrow t \Rightarrow \exists s \sim t. \ W\rightsquigarrow s$$
$$\|\|$$
$$CC^{\sim}$$

$$\begin{array}{ccccc}
\forall \pi_T. \ \forall W. \ W \models \tilde{\sigma}(\pi_T) & & & & \forall \pi_S. \ \forall W. \ W \models \pi_S \\
\Rightarrow W\downarrow \models \pi_T & \equiv & TP^{\tilde{\sigma}} \Longleftrightarrow TP^{\tilde{\tau}} & \equiv & \Rightarrow W\downarrow \models \tilde{\tau}(\pi_S)
\end{array}$$

Fig. 1. The equivalent compiler correctness definitions forming our trinitarian view.

**Contributions.**

- We propose a new trinitarian view of compiler correctness that accounts for non-trivial relations between source and target traces. While, as discussed above, specific instances of the $CC^{\sim}$ definition have already been used in practice, we seem to be the first to propose assessing the meaningfulness of $CC^{\sim}$ instances in terms of how properties are preserved between the source and the target, and in particular by looking at the property mappings $\tilde{\sigma}$ and $\tilde{\tau}$ induced by the trace relation $\sim$. We prove that $CC^{\sim}$, $TP^{\tilde{\sigma}}$, and $TP^{\tilde{\tau}}$ are equivalent for any trace relation (§2.2), as illustrated in Figure 1. In the opposite direction, we show that for every trace relation corresponding to a given Galois connection [? ], an analogous equivalence holds.
- We extend these results from the preservation of trace properties to the larger class of subset-closed hyperproperties, e.g., noninterference (§3.1)[4], and to the classes of safety properties (§3.2) and all hyperproperties (§3.3).
- We use $CC^{\sim}$ compilers of various complexities to illustrate that our view on compiler correctness naturally accounts for undefined behavior (§4.1), resource exhaustion (§4.2), different source and target values (§4.3), and differences in the granularity of data and observable events (§4.4). We expect these ideas to extend to other discrepancies between source and target traces. For each compiler we show how to choose the relation between source and target traces and how the induced property mappings preserve interesting trace properties and subset-closed hyperproperties. We look at the way particular $\tilde{\sigma}$ and $\tilde{\tau}$ work on different kinds of properties and how the produced properties can be expressed for different kinds of traces.
- We analyze the impact of correct compilation on noninterference [? ], showing what can still be preserved (and thus also what is lost) when target observations are finer than source ones, e.g., side-channel observations (§5). We formalize the guarantee obtained by correct compilation of a noninterfering program as *abstract noninterference* [? ], a weakening of target noninterference. Dually, we identify a family of declassifications of target noninterference for which source reasoning is possible.

---

[4]Given the deterministic nature of our programs, we consider notions of noninterference that are often used in deterministic languages. We leave notions of noninterference in nondeterministic languages for future work.

- We show that the trinitarian view also extends to a large class of *secure compilation* definitions [? ], formally characterizing the protection of the compiled program against linked adversarial code (§6). For each secure compilation definition we again propose both a property-free characterization in the style of CC~, and two characterizations in terms of preserving a class of source or target properties satisfied against arbitrary adversarial contexts. The additional quantification over contexts allows for finer distinctions when considering different property classes, so we study mapping classes not only of trace properties and hyperproperties, but also of relational hyperproperties [? ].
- We provide instances of secure compilers that preserve three different classes of hyperproperties (trace, safety and hypersafety properties) when targeting a language with additional trace events that are not possible in the source (§7).

The results and insights that we provide often follow one's expected intuition and may be considered unsurprising. However our framework is the first to capture such expectations formally and precisely, and as such it provides a uniform way to discuss these and to formalise future (possibly surprising) ones. The paper closes with discussions of related (§8) and future work (§9). Some technical proofs can be found in the appendix (§A).

The traces considered in our examples are structured, usually as sequences of events. We notice however that unless explicitly mentioned, all our definitions and results are more general and make no assumption whatsoever about the structure of traces. Most of the theorems formally or informally mentioned in the paper were mechanized in the Coq proof assistant and are marked with ✎. This development has around 10k lines of code and is available at the following address: https://github.com/secure-compilation/different_traces.

## 2 TRACE-RELATING COMPILER CORRECTNESS

In this section, we start by generalizing the trace property preservation definitions at the end of the introduction to $\text{TP}^\sigma$ and $\text{TP}^\tau$, which depend on two *arbitrary* mappings $\sigma$ and $\tau$ (§2.1). We prove that, whenever $\sigma$ and $\tau$ form a Galois connection, $\text{TP}^\sigma$ and $\text{TP}^\tau$ are equivalent (Theorem 2.4). We then exploit a bijective correspondence between trace relations and Galois connections to close the trinitarian view (§2.2), with two main benefits: first, it helps us assess the meaningfulness of a given trace relation by looking at the property mappings it induces; second, it allows us to construct new compiler correctness definitions starting from a desired mapping of properties. Finally, we generalize the classic result that compiler correctness (i.e., CC=) is enough to preserve not just trace properties but also all subset-closed hyperproperties [? ]. For this, we show that CC~ is also equivalent to subset-closed hyperproperty preservation, for which we also define both a version in terms of $\tilde{\sigma}$ and a version in terms of $\tilde{\tau}$ (§3.1).

## 2.1 Property Mappings

As explained in §1, trace-relating compiler correctness CC~, by itself, lacks a crisp description of which trace properties are preserved by compilation. Since even the syntax of traces can differ between source and target, one can either focus on trace properties of the source (and then interpret them in the target), or on trace properties of the target (and then interpret them in the source). Formally we need two property mappings, $\tau : 2^{\text{Trace}_S} \to 2^{\text{Trace}_T}$ and $\sigma : 2^{\text{Trace}_T} \to 2^{\text{Trace}_S}$, which lead us to the following generalization of trace property preservation (TP).

**Definition 2.1** ($\text{TP}^\sigma$ and $\text{TP}^\tau$). Given two property mappings, $\tau : 2^{\text{Trace}_S} \to 2^{\text{Trace}_T}$ and $\sigma : 2^{\text{Trace}_T} \to 2^{\text{Trace}_S}$, for a compilation chain $\cdot\!\downarrow$ we define $\text{TP}^\tau$ and $\text{TP}^\sigma$ as follows:

$$\text{TP}^\tau \equiv \forall \pi_S. \ \forall W. \ W \models \pi_S \Rightarrow W\!\downarrow \models \tau(\pi_S)$$

$$\mathsf{TP}^{\sigma} \; \equiv \; \forall \boldsymbol{\pi}_{\mathrm{T}}.\; \forall \mathsf{W}.\; \mathsf{W} \models \sigma(\boldsymbol{\pi}_{\mathrm{T}}) \Rightarrow \mathsf{W}{\downarrow} \models \boldsymbol{\pi}_{\mathrm{T}}$$

For an arbitrary source program $\mathsf{W}$, $\tau$ interprets a source property $\pi_{\mathrm{S}}$ as the *target guarantee* for $\mathsf{W}{\downarrow}$. Dually, $\sigma$ defines a *source obligation* sufficient for the satisfaction of a target property $\boldsymbol{\pi}_{\mathrm{T}}$ after compilation. Ideally:

  i) Given $\boldsymbol{\pi}_{\mathrm{T}}$, the target interpretation of the source obligation $\sigma(\boldsymbol{\pi}_{\mathrm{T}})$ should actually guarantee that $\boldsymbol{\pi}_{\mathrm{T}}$ holds, i.e., $\tau(\sigma(\boldsymbol{\pi}_{\mathrm{T}})) \subseteq \boldsymbol{\pi}_{\mathrm{T}}$;

  ii) Dually for $\pi_{\mathrm{S}}$, we would not want the source obligation for $\tau(\pi_{\mathrm{S}})$ to be harder than $\pi_{\mathrm{S}}$ itself, i.e., $\sigma(\tau(\pi_{\mathrm{S}})) \supseteq \pi_{\mathrm{S}}$.

These requirements are satisfied when the two maps form a *Galois connection* between the posets of source and target properties ordered by inclusion. We briefly recall the definition and the characteristic property of Galois connections [? ? ].

**Definition 2.2** (Galois connection). Let $(X, \leq)$ and $(Y, \sqsubseteq)$ be two posets. A pair of maps, $\alpha : X \to Y$, $\gamma : Y \to X$ is a Galois connection *iff* it satisfies the *adjunction law*: $\forall x \in X.\; \forall y \in Y.\; \alpha(x) \sqsubseteq y \iff x \leq \gamma(y)$. $\alpha$ (resp. $\gamma$) is the lower (upper) adjoint or abstraction (concretization) function and $Y$ ($X$) the abstract (concrete) domain.

We will often write $\alpha : (X, \leq) \leftrightarrows (Y, \sqsubseteq) : \gamma$ to denote a Galois connection, or simply $\alpha : X \leftrightarrows Y : \gamma$, or even $\alpha \leftrightarrows \gamma$ when the involved posets are clear from context.

**Lemma 2.3** (Characteristic property of Galois connections). If $\alpha : (X, \leq) \leftrightarrows (Y, \sqsubseteq) : \gamma$ is a Galois connection, then $\alpha, \gamma$ are monotone and $id \leq \gamma \circ \alpha$ and $\alpha \circ \gamma \sqsubseteq id$, i.e.,

$$\forall x \in X.\; x \leq \gamma(\alpha(x))$$

$$\forall y \in Y.\; \alpha(\gamma(y)) \sqsubseteq y$$

If $X, Y$ are complete lattices, then $\alpha$ is continuous, i.e., $\forall F \subseteq X.\; \alpha(\bigsqcup F) = \bigsqcup \alpha(F)$.

If two property mappings, $\tau$ and $\sigma$, form a Galois connection on trace properties ordered by set inclusion, Lemma 2.3 (with $\alpha = \tau$ and $\gamma = \sigma$) tells us that they satisfy conditions *i*), *ii*) above, i.e., $\tau(\sigma(\boldsymbol{\pi}_{\mathrm{T}})) \subseteq \boldsymbol{\pi}_{\mathrm{T}}$ and $\sigma(\tau(\pi_{\mathrm{S}})) \supseteq \pi_{\mathrm{S}}$.[5] These conditions on $\tau$ and $\sigma$ are sufficient to show the equivalence of the criteria they define, respectively $\mathsf{TP}^{\tau}$ and $\mathsf{TP}^{\sigma}$.

**Theorem 2.4** ($\mathsf{TP}^{\tau}$ and $\mathsf{TP}^{\sigma}$ coincide 🐦). Let $\tau : 2^{\mathrm{Trace_S}} \rightleftarrows 2^{\mathrm{Trace_T}} : \sigma$ be a Galois connection, with $\tau$ and $\sigma$ the lower and upper adjoints (resp.). Then $\mathsf{TP}^{\tau} \iff \mathsf{TP}^{\sigma}$.

PROOF. Notice that if a program satisfies a property $\pi$, then it satisfies every less restrictive i.e., bigger property $\pi' \supseteq \pi$. Building on this:

  ($\Rightarrow$) Assume $\mathsf{TP}^{\tau}$ and that $\mathsf{W}$ satisfies $\sigma(\boldsymbol{\pi}_{\mathrm{T}})$. Apply $\mathsf{TP}^{\tau}$ to $\mathsf{W}$ and $\sigma(\boldsymbol{\pi}_{\mathrm{T}})$ and deduce that $\mathsf{W}{\downarrow}$ satisfies $\tau(\sigma(\boldsymbol{\pi}_{\mathrm{T}})) \subseteq \boldsymbol{\pi}_{\mathrm{T}}$.

  ($\Leftarrow$) Assume $\mathsf{TP}^{\sigma}$ and that $\mathsf{W}$ satisfies $\pi_{\mathrm{S}} \subseteq \sigma(\tau(\pi_{\mathrm{S}}))$. Apply $\mathsf{TP}^{\sigma}$ to $\mathsf{W}$ and $\sigma(\tau(\pi_{\mathrm{S}}))$ deducing $\mathsf{W}{\downarrow}$ satisfies $\tau(\pi_{\mathrm{S}})$.

□

## 2.2 Trace Relations and Property Mappings

We now investigate the relation between $\mathsf{CC}^{\sim}$, $\mathsf{TP}^{\tau}$ and $\mathsf{TP}^{\sigma}$. We show that for a trace relation and its corresponding Galois connection (Lemma 2.7), the three criteria are equivalent (Theorem 2.8). This equivalence offers interesting insights for both verification and the design of a correct compiler. For a $\mathsf{CC}^{\sim}$ compiler, the equivalence makes explicit both the guarantees one has after compilation

---

[5]While target traces are often *"more concrete"* than source ones, trace properties $2^{\mathrm{Trace}}$ (which in Coq we represent as the function type Trace→Prop) are contravariant in Trace and thus target properties correspond to the *abstract domain*.

($\tilde{\tau}$) and source proof obligations to ensure the satisfaction of a given target property ($\tilde{\sigma}$). On the other hand, a compiler designer might first determine the target guarantees the compiler itself must provide, i.e., $\tau$, and then prove an equivalent statement, $\mathrm{CC}^\sim$, for which more convenient proof techniques exist in the literature [**? ?** ].
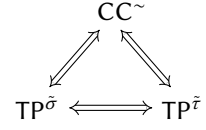
**Definition 2.5** (Existential and Universal Image [**?** ]). Given any two sets $X$ and $Y$ and a relation $\sim \subseteq A \times B$, define the relation's existential or direct image, $\tilde{\tau} : 2^X \to 2^Y$ and its universal image, $\tilde{\sigma} : 2^Y \to 2^X$ as follows:

$$\tilde{\tau} = \lambda\,\pi \in 2^X.\ \{y \mid \exists x.\, x \sim y \wedge x \in \pi\}$$

$$\tilde{\sigma} = \lambda\,\pi \in 2^Y.\ \{x \mid \forall y.\, x \sim y \Rightarrow y \in \pi\}$$

When trace relations are considered, the corresponding existential and universal images can be used to instantiate Definition 2.1 leading to the trinitarian view already mentioned in §1.

**Theorem 2.6** (Trinitarian View 🐿). For any trace relation $\sim$ and its existential and universal images $\tilde{\tau}$ and $\tilde{\sigma}$, we have:

$$
\begin{array}{ccc}
 & \mathrm{CC}^\sim & \\
 & \nearrow \quad \nwarrow & \\
\mathrm{TP}^{\tilde{\sigma}} & \Longleftrightarrow & \mathrm{TP}^{\tilde{\tau}}
\end{array}
$$

This result relies both on Theorem 2.4 and on the fact that the existential and universal images of a trace relation form a Galois connection (🐿 ). The theorem can be stated in a slightly more general form (Theorem 2.8), exploiting an isomorphism between the category of sets and relations and a sub category of monotonic predicate transformers [**?** ]. We specialize this isomorphism to what is of interest for our purposes and deduce a bijective correspondence between trace relations and Galois connections on properties.

**Lemma 2.7** (Trace relations $\cong$ Galois connections on trace properties). The function $\sim \mapsto \tilde{\tau} \leftrightarrows \tilde{\sigma}$ that maps a trace relation to its existential and universal images is a bijection between trace relations $2^{\mathrm{Trace}_S \times \mathrm{Trace}_T}$ and Galois connections on trace properties $2^{\mathrm{Trace}_S} \leftrightarrows 2^{\mathrm{Trace}_T}$. Its inverse is $\tau \leftrightarrows \sigma \mapsto \hat{\sim}$, where $\mathsf{s} \mathbin{\hat{\sim}} \mathsf{t} \equiv \mathsf{t} \in \tau(\{\mathsf{s}\})$.

The bijection just introduced allows us to generalize Theorem 2.6 and switch anytime between the three views of compiler correctness described earlier.

**Theorem 2.8** (Correspondence of Criteria). For any trace relation $\sim$ and corresponding Galois connection $\tau \leftrightarrows \sigma$, we have: $\mathrm{TP}^\tau \iff \mathrm{CC}^\sim \iff \mathrm{TP}^\sigma$.

Note that sometimes the lifted properties may be trivial: the target guarantee can be the true property (the set of all traces), or the source obligation the false property (the empty set of traces). This might be the case when source observations abstract away too much information (§4.2 presents an example).

# 3 PRESERVING OTHER (HYPER)PROPERTY CLASSES

In this section we investigate how to preserve other classes of (hyper)properties beyond trace properties: subset-closed hyperproperties (§3.1) safety properties (§3.2) and arbitrary hyperproperties that are not just subset-closed (§3.3). For each of these classes, we start by giving an intuition of what it means to preserve such a class in the equal-trace setting, then we study preservation of that class in the trace-relating setting. For subset-closed hyperproperties we have to refine the Galois connection to ensure the information "$H_S$ is subset-closed" is not lost with the application of $\tilde{\tau}$. Similarly, when looking at safety properties, we have to preserve the information that a propery is a safety property. For arbitrary hyperproperties one might instead require that no information at all is lost during the (pre or post) composition of $\tilde{\tau}$ and $\tilde{\sigma}$. The section concludes with a comparison of the criteria in terms of relative strengths (§3.4).

## 3.1  Preservation of Subset-Closed Hyperproperties

Hyperproperty preservation is a strong requirement in general. Fortunately, many interesting hyperproperties are *subset-closed* (*SCH* for short) (e.g., noninterference), and these are known to be preserved by refinement [? ]. When the trace semantics is common to source and target languages, a subset-closed hyperproperty is preserved if the behaviors of the compiled program refine the behaviors of the source program, which coincides with the statement of $CC^=$. We generalize this result to the trace-relating setting, introducing two other equivalent characterizations of $CC^\sim$ in terms of preservation of subset-closed hyperproperties (Theorem 3.3). In order to do so we close under subsets the images of both $\tilde{\tau}$ and $\tilde{\sigma}$ so that source subset-closed hyperproperties are mapped to target subset-closed ones and viceversa.

First, a hyperproperty $H$ is defined as a set of sets of traces, $H \in 2^{2^{\text{Trace}}}$ (recall that *Traces* is the set of all traces) [? ]. A program satisfies a hyperproperty when its complete set of traces, which from now on we will call its *behavior*, is a member of the hyperproperty.

**Definition 3.1** (Hyperproperty Satisfaction [? ]). A program $W$ satisfies a hyperproperty $H$, written $W \models H$,[6] iff $beh(W) \in H$, where $beh(W) = \{t \mid W \leadsto t\}$.

To talk about hyperproperty preservation in the trace-relating setting, we need an interpretation of source hyperproperties into the target and vice versa. The one we consider builds on top of the two trace property mappings $\tau$ and $\sigma$, which are naturally lifted to hyperproperty mappings. This way we are able to extract two hyperproperty mappings from a trace relation similarly to §2.2:

**Definition 3.2** (Lifting property mappings to hyperproperty mappings). Let $\tau : 2^{\text{Trace}_S} \to 2^{\text{Trace}_T}$ and $\sigma : 2^{\text{Trace}_T} \to 2^{\text{Trace}_S}$ be arbitrary property mappings. The images of $H_S \in 2^{2^{\text{Trace}_S}}$, $H_T \in 2^{2^{\text{Trace}_T}}$ under $\tau$ and $\sigma$ are, respectively:

$$\tilde{\tau}(H_S) = \{\tau(\pi_S) \mid \pi_S \in H_S\} \qquad \tilde{\sigma}(H_T) = \{\sigma(\pi_T) \mid \pi_T \in H_T\}$$

Formally, we are defining two new mappings, this time on hyperproperties, but with a small abuse of notation we still denote them $\tilde{\tau}$ and $\tilde{\sigma}$.

Interestingly, it is not possible to apply the argument used for $CC^=$ to show that a $CC^\sim$ compiler guarantees $W{\downarrow} \models \tilde{\tau}(H_S)$ whenever $W \models H_S$. This is because direct images do not necessarily preserve subset-closure [? ? ].[JT: not a sentence?] We therefore close the image of $\tilde{\tau}$ and $\tilde{\sigma}$ under subsets (denoted as $Cl_\subseteq$) and obtain the following result:

**Theorem 3.3** (Preservation of Subset-Closed Hyperproperties 🐾). For any trace relation $\sim$ and its existential and universal images lifted to hyperproperties, $\tilde{\tau}$ and $\tilde{\sigma}$, and for $Cl_\subseteq(H) = \{\pi \mid \exists \pi' \in H. \ \pi \subseteq \pi'\}$, we have the following:

$$\text{SCHP}^{Cl_\subseteq \circ \tilde{\tau}} \equiv \forall W \forall H_S \in \text{SCH}_S. W \models H_S \Rightarrow W{\downarrow} \models Cl_\subseteq(\tilde{\tau}(H_S));$$
$$\text{SCHP}^{Cl_\subseteq \circ \tilde{\sigma}} \equiv \forall W \forall H_T \in \text{SCH}_T. W \models Cl_\subseteq(\tilde{\sigma}(H_T)) \Rightarrow W{\downarrow} \models H_T.$$

$$CC^\sim$$
$$\text{SCHP}^{Cl_\subseteq \circ \tilde{\tau}} \Longleftrightarrow \text{SCHP}^{Cl_\subseteq \circ \tilde{\sigma}}$$

The use of $Cl_\subseteq$ in Theorem 3.3 implies a loss of precision in preserving subset-closed hyperproperties through compilation. In §5, we focus on a specific security-relevant subset-closed hyperproperty, noninterference, and show that such a loss of precision can be seen as a declassification. Instead, now we define the trinity and the related formal machinery for safety properties preservation.

---

[6]In case of ambiguity with property satisfaction the class of $H$ will be made explicit.

## 3.2 Preserving Safety Properties

The class of *Safety* properties collects all trace properties prescribing that *"something bad never happens"* or equivalently, all trace properties whose violation can be monitored and, once observed, no longer restored [? ]. More abstractly safety properties can be defined as the closed sets of a topology [? ? ], with no need to consider any particular structure on the traces. To ease the presentation, we consider the trace model adopted by ? ] where traces resemble lists and streams of events. This model naturally comes with a notion of *prefixes* and a relation between a prefix $m$ and a trace $t$, written $m \leq t$. Intuitively, $\pi$ is a safety property if any trace $t$ *violating* the property extends a "bad prefix" $m$ that witnesses such a violation. Every safety property is therefore uniquely defined by the set of its "bad prefixes". We recall below the definition and the characterization of safety properties in terms of sets of finite prefixes $m$.

**Definition 3.4** (Safety Properties [? ]). Let $\pi$ be a trace property. Then,

$\pi \in Safety$ iff $\forall t \notin \pi.\ \exists m \leq t.\ \forall t'.\ m \leq t' \Rightarrow t' \notin \pi$.

Equivalently, $\pi \in Safety$ iff there exists a set of finite prefixes $M$, such that

$$\forall t.\ t \notin \pi \iff (\exists m \in M.\ m \leq t)$$

Due to this characterization of safety properties through finite prefixes (Definition 3.4), the preservation of all and only the safety properties is equivalent to $CC^=$ restricted to finite prefixes. [CA: I don't have a reference for this, but it is folklore and immediate to proof. – If I have time I will add a paper proof in the appendix]

**Theorem 3.5.** The following are equivalent:

$$SC^= \equiv \forall W, m.\ W{\downarrow}\rightsquigarrow^* m \Rightarrow W \rightsquigarrow^* m$$

$$SP \equiv \forall \pi \in Safety.\ W \models \pi \Rightarrow W{\downarrow} \models \pi$$

where $W \rightsquigarrow^* m$ stands for $\exists t.\ m \leq t \wedge W \rightsquigarrow t$.

Unfolding $\rightsquigarrow^*$ we can interpret $SC^=$ as follows. Whenever $W{\downarrow}$ produces a trace $t \geq m$ that violates a specific safety property, namely, the one defined by the singleton prefix set $\{m\}$, then $W$ violates the *same* safety property, producing a trace $t' \geq m$ but possibly distinct from $t$.

The generalization we propose of $SC^=$ to the trace-relating setting, states that whenever $W{\downarrow}$ produces a trace $t$ that violates a target safety property, then $W$ violates the source *interpretation* of the property, i.e., its image through $\tilde{\sigma}$.[7] The following theorem defines $SC^\sim$ and its two equivalent formulations.

**Theorem 3.6** (Trinitarian view for Safety). For a trace relation $\sim\ \subseteq \text{Traces}_S \times \text{Trace}_T$ and its corresponding property mappings $\tilde{\sigma}$ and $\tilde{\tau}$, the following are equivalent:

$$SC^\sim \equiv \forall W \forall t \forall m \leq t.\ W{\downarrow}\rightsquigarrow t \Rightarrow \exists t' \geq m \exists s \sim t'.\ W \rightsquigarrow s$$

$$SP^{\tilde{\sigma}} \equiv \forall W\ \forall \pi_T \in Safety_T.\ W \models \tilde{\sigma}(\pi_T) \Rightarrow W{\downarrow} \models \pi_T$$

$$TP^{Safe \circ \tilde{\tau}} \equiv \forall W\ \forall \pi_S \in 2^{Traces_S}.\ W \models \pi_S \Rightarrow W{\downarrow} \models (Safe \circ \tilde{\tau})(\pi_S)$$

$$\begin{array}{c} SC^\sim \\ \nearrow \quad \nwarrow \\ SP^{\tilde{\sigma}} \iff TP^{Safe \circ \tilde{\tau}} \end{array}$$

Coherent with the informal meaning we aimed to give to $SC^\sim$, $SP^{\tilde{\sigma}}$ quantifies over target safety properties, while $TP^{Safe \circ \tilde{\tau}}$ quantifies over *arbitrary* source properties, but imposes the composition of $\tilde{\tau}$ with *Safe*, which maps an arbitrary target property $\pi_T$ to the target safety property that

---

[7]At least one other symmetric generalization is possible: For $\pi_S \in Safety_S$ defined by $M = \{m\}$, if $W{\downarrow}$ produces a trace $t$ that violates the target interpretation of $\pi_S$, i.e., $\tilde{\tau}(\pi_S)$, then $W$ produces $s \geq m$ thus violating $\pi_S$.

best over-approximates $\pi_T$.[8] More precisely, $Safe$ is a closure operator on target properties, with $\mathbf{Safety_T} = \left\{ Safe(\pi_T) \mid \pi_T \in 2^{\mathrm{Trace_T}} \right\}$ being the class of target safety properties.

In Figure 2 the blue and red ellipses represent source and target properties properties respectively and are connected by $\tilde{\tau} \leftrightarrows \tilde{\sigma}$. The red ellipse is the class of all target safety properties. $Safe \leftrightarrows id$ is a Galois connection between target properties and the target safety properties, as $Safe$ is a closure operator [? ]. Finally, the composition of Galois connections is still a Galois connection [? ]. Hence,

$$Safe \circ \tilde{\tau} : 2^{\mathrm{Trace_S}} \leftrightarrows \mathbf{Safety_T} : \tilde{\sigma}$$

is a Galois connection between source properties and target safety properties, that we used to prove the equivalence $\mathrm{TP}^{Safe \circ \tilde{\tau}} \iff \mathrm{SP}^{\tilde{\sigma}}$ (🐿). We notice that this argument generalizes to arbitrary



Fig. 2. Composition of $\tilde{\tau} \leftrightarrows \tilde{\sigma}$ and $Safe \leftrightarrows id$.

closure operators on target properties (🐿). We come back to this in §6, where more such results will be needed when considering other classes of properties being preserved by secure compilers. Now, we define the trinity for arbitrary hyperproperties, not just the subset-closed ones.

### 3.3 Preserving Non-Subset Closed Hyperproperties

Subset-closed hyperproperties are not expressive enough to all capture interesting properties, e.g., possibilistic notions of information-flow [? ], so we aim to briefly discuss the preservation of *arbitrary* hyperproperties. In general, one cannot lift a Galois connection over trace properties to a Galois connection over arbitrary hyperproperties.

While two out of three of the criteria we introduce in this section are equivalent under no assumptions ($\mathrm{HC}^{\sim} \iff \mathrm{HP}^{\tilde{\tau}}$), for a comparison with the third one we require that no information is lost in the pre or post composition of $\tau$ and $\sigma$. For this, we label the trinity in Theorem 3.8 as *weak*.

To start, we note that the following strengthening of $\mathrm{CC}^{=}$, denoted $\mathrm{HC}^{=}$, is equivalent to the preservation of arbitrary hyperproperties. Here, $\mathrm{beh}(\mathsf{W})$ is the set of all traces of $\mathsf{W}$. [CA: I don't have a reference for this, but it is folklore and immediate to proof. – If I have time I will add a paper proof in the appendix]

**Theorem 3.7** ($\mathrm{HC}^{=}$, HP). The followings are equivalent

$$\mathrm{HC}^{=} \equiv \forall \mathsf{W}.\ \mathbf{beh}(\mathsf{W}{\downarrow}) = \mathrm{beh}(\mathsf{W})$$

$$\mathrm{HP} \equiv \forall \mathsf{W}\ \forall H \in 2^{2^{\mathrm{Trace}}}.\ \mathsf{W} \models H \iff \mathsf{W}{\downarrow} \models H$$

---

[8] $Safe(\pi_T) = \cap \{S_T \mid \pi_T \subseteq S_T \wedge S_T \in \mathbf{Safety_T}\}$ is the topological closure in the topology where safety properties coincide with the closed sets (see, e.g., ? ] and ? ]).

$HC^=$ requires that the behavior of $W\!\downarrow$ is exactly the same as the behavior of $W$. We generalize this to the trace-relating setting, by requiring that the behavior of $W\!\downarrow$ coincide with the target interpretation of the source properties describing the behavior of $W$.[9]

**Theorem 3.8** (Weak Trinity for Hyperproperties ✍). *For a trace relation $\sim\ \subseteq\ \text{Traces}_S \times \mathbf{Trace_T}$ and induced property mappings $\tilde{\sigma}$ and $\tilde{\tau}$, we have:*
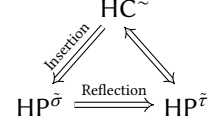
$HC^\sim \iff HP^{\tilde{\tau}}$;

*if $\tilde{\tau} \leftrightarrows \tilde{\sigma}$ is a Galois insertion (i.e., $\tilde{\tau} \circ \tilde{\sigma} = id$), then $HC^\sim \Rightarrow HP^{\tilde{\sigma}}$,*

*if $\tilde{\sigma} \leftrightarrows \tilde{\tau}$ is a Galois reflection (i.e., $\tilde{\sigma} \circ \tilde{\tau} = id$), then $HP^{\tilde{\sigma}} \Rightarrow HP^{\tilde{\tau}}$,*

$$HC^\sim \equiv \forall W.\ beh(W\!\downarrow) = \tilde{\tau}(beh(W))$$

$$HP^{\tilde{\tau}} \equiv \forall W\ \forall H_S.\ W \models H_S \Rightarrow W\!\downarrow\ \models \tilde{\tau}(H_S)$$

$$HP^{\tilde{\sigma}} \equiv \forall W\ \forall \mathbf{H_T}.\ W \models \tilde{\sigma}(\mathbf{H_T}) \Rightarrow W\!\downarrow\ \models \mathbf{H_T}$$



In other words, it is still possible (and sound) to deduce a source obligation for a given target hyperproperty $\mathbf{H_T}$ ($HC^\sim \Rightarrow HP^{\tilde{\sigma}}$) when no information is lost in the composition $\tilde{\tau} \circ \tilde{\sigma}$. Dually, $HP^{\tilde{\tau}}$ (and hence $HC^\sim$) is a consequence of $HP^{\tilde{\sigma}}$ when no information is lost in composing in the other direction, $\tilde{\sigma} \circ \tilde{\tau}$.

### 3.4 Comparing the Presented Criteria

At this point we have presented four trinities of criteria that preserve trace properties, subset-closed hyperproperties, safety properties and arbitrary hyperproperties. Figure 3 sums up our trinities and orders them according their relative strength.



Fig. 3. Generalization of Compiler Correctness and its trace-relating variations.

[CA: I think this is necessary. Still no reference nor a proof but immediate. ] In §6 we will also consider, in the setting of *secure* compilation, the class of safety hyperproperties or hyper-safety, and relational hyperproperties. In the setting of *correct* compilation – that focuses only on whole programs – it is straightforward to show that the trinity for hypersafety coincides with the one for safety properties in the same way the trinity of trace properties and subset-closed hyperproperties coincide. Similarly the trinity for relational hyperproperties coincides with the one for hyperproperties.

## 4 INSTANCES OF TRACE-RELATING COMPILER CORRECTNESS

The trace-relating view of compiler correctness above can serve as a unifying framework for studying a range of interesting compilers. This section provides several representative instantiations of the framework: source languages with undefined behavior that compilation can turn into arbitrary

---

[9] At least one generalization is possible: $\tilde{\sigma}(beh(W\!\downarrow)) = beh(W)$. In this case, $HC^\sim \iff HP^{\tilde{\sigma}}$ holds unconditionally while the other two implications hold under the same, but swapped, hypotheses from Theorem 3.8.

target behavior (§4.1), target languages with resource exhaustion that cannot happen in the source (§4.2), changes in the representation of values (§4.3), and differences in the granularity of data and observable events (§4.4).

## 4.1 Undefined Behavior

We start by expanding upon the discussion of undefined behavior in §1. We first study the model of CompCert, where source and target alphabets are the same, including the event for undefined behavior. The trace relation weakens equality by allowing undefined behavior to be replaced with an arbitrary sequence of events.

**Example 4.1** (CompCert-like Undefined Behavior Relation). Source and target traces are sequences of events drawn from $\Sigma$, where $Wrong \in \Sigma$ is a terminal event that represents an undefined behavior. We then use the trace relation defined in the introduction:

$$\mathsf{s} \sim \mathsf{t} \equiv \mathsf{s} = \mathsf{t} \vee \exists m \leq \mathsf{t}.\ \mathsf{s} = m \cdot Wrong$$

Each trace of a target program produced by a $CC^\sim$ compiler either also is a trace of the original source program or has a finite prefix that the source program also produces, immediately before encountering undefined behavior. As explained in §1, one of the correctness theorems in CompCert can be rephrased as this variant of $CC^\sim$.

We proved that the property mappings induced by the relation can be written as (🐿):

$$\tilde{\sigma}(\pi_{\mathsf{T}}) = \{\mathsf{s} \mid \mathsf{s} \in \pi_{\mathsf{T}} \wedge \mathsf{s} \neq m \cdot Wrong\} \cup \{m \cdot Wrong \mid \forall \mathsf{t}.\ m \leq \mathsf{t} \Rightarrow \mathsf{t} \in \pi_{\mathsf{T}}\}$$

$$\tilde{\tau}(\pi_{\mathsf{S}}) = \{\mathsf{t} \mid \mathsf{t} \in \pi_{\mathsf{S}}\} \cup \{\mathsf{t} \mid \exists m \leq \mathsf{t}.\ m \cdot Wrong \in \pi_{\mathsf{S}}\}$$

These two mappings explain what a $CC^\sim$ compiler ensures for the $\sim$ relation above. The target-to-source mapping $\tilde{\sigma}$ states that to prove that a compiled program has a property $\pi_{\mathsf{T}}$ using source-level reasoning, one has to prove that any trace produced by the source program must either be a target trace satisfying $\pi_{\mathsf{T}}$ or have undefined behavior, but only provided that *any continuation* of the trace substituted for the undefined behavior satisfies $\pi_{\mathsf{T}}$. The source-to-target mapping $\tilde{\tau}$ states that by compiling a program satisfying a property $\pi_{\mathsf{S}}$ we obtain a program that produces traces that satisfy the same property or that extend a source trace that ends in undefined behavior.

These definitions can help us reason about programs. For instance, $\tilde{\sigma}$ specifies that, to prove that an event does not happen in the target, it is not enough to prove that it does not happen in the source: it is also necessary to prove that the source program does not have any undefined behavior (second disjunct). Indeed, if it had an undefined behavior, its continuations could exhibit the unwanted event.                                                                                    ⊡

This relation can be easily generalized to other settings. For instance, consider the setting in which we compile down to a low-level language like machine code. Target traces can now contain new events that cannot occur in the source: indeed, in modern architectures like x86 a compiler typically uses only a fraction of the available instruction set. Some instructions might even perform dangerous operations, such as writing to the hard drive, or controlling a device that is hidden from the source language. Formally, the source and target do not have the same events any more. Thus, we consider a source alphabet $\Sigma_{\mathsf{S}} = \Sigma \cup \{Wrong\}$, and a target alphabet $\Sigma_{\mathsf{T}} = \Sigma \cup \Sigma'$. The trace relation is defined in the same way and we obtain the same property mappings as above, except that target traces now have more events (some of which may be dangerous), the arbitrary continuations of target traces get more interesting. For instance, consider a new event that represents writing data on the hard drive, and suppose we want to prove that this event cannot happen for a compiled program. Then, proving this property requires exactly proving that the source program exhibits no undefined behavior [? ]. More generally, what one can prove about target-only events can only be

either that they cannot appear (because there is no undefined behavior) or that any of them can appear (in the case of undefined behavior).

In §7.1 we study a similar example, showing that even in a safe language linked adversarial contexts can cause dangerous target events that have no source correspondent.

## 4.2 Resource Exhaustion

Let us return to the discussion about resource exhaustion in §1.

**Example 4.2** (Resource Exhaustion). We consider traces made of events drawn from $\Sigma_S$ in the source, and $\Sigma_T = \Sigma_S \cup \{\text{Resource\_Limit\_Hit}\}$ in the target. Recall the trace relation for resource exhaustion:

$$s \sim t \equiv s = t \vee \exists m \le s.\ t = m \cdot \text{Resource\_Limit\_Hit}$$

Formally, this relation is similar to the one for undefined behavior, except this time it is the target trace that is allowed to end early instead of the source trace.

The induced trace property mappings $\tilde{\sigma}$ and $\tilde{\tau}$ are the following (🐓):

$$\tilde{\sigma}(\pi_T) = \{s \mid s \in \pi_T\} \cap \{s \mid \forall m \le s.\ m \cdot \text{Resource\_Limit\_Hit} \in \pi_T\}$$

$$\tilde{\tau}(\pi_S) = \pi_S \cup \{m \cdot \text{Resource\_Limit\_Hit} \mid \exists s \in \pi_S.\ m \le s\}$$

These capture the following intuitions. The target-to-source mapping $\tilde{\sigma}$ states that to prove a property of the compiled program one has to show that the traces of the source program satisfy two conditions: (1) they must also satisfy the target property; and (2) the termination of every one of their prefixes by a resource exhaustion error must be allowed by the target property. This is rather restrictive: any property that prevents resource exhaustion cannot be proved using source-level reasoning. Indeed, if $\pi_T$ does not allow resource exhaustion, then $\tilde{\sigma}(\pi_T) = \varnothing$. This is to be expected since resource exhaustion is simply not accounted for at the source level. The source-to-target mapping $\tilde{\tau}$ states that a compiled program produces traces that either belong to the same properties as the traces of the source program or end early due to resource exhaustion.

In this example, safety properties [? ] are mapped (in both directions) to other safety properties (🐓). This can be desirable for a relation: since safety properties are usually easier to reason about, one interested only in safety properties at the target can reason about them using source-level reasoning tools for safety properties. To reason about safety, one would use the criteria presented in §3.2

Since it focuses on traces and not just safety, the compiler correctness theorem in CakeML is an instance of CC$^\sim$ for the $\sim$ relation above. We have also implemented two small compilers that are correct for this relation. The full details can be found in the Coq development.The first compiler (🐓) goes from a simple expression language (similar to the one in §4.3 but without inputs) to the same language except that execution is bounded by some amount of fuel: each execution step consumes some amount of fuel and execution immediately halts when it runs out of fuel. The compiler is the identity.

The second compiler (🐓) is more interesting: we proved this CC$^\sim$ instance for a variant of a compiler from a WHILE language to a simple stack machine by Xavier Leroy [? ]. We enriched the two languages with outputs and modified the semantics of the stack machine so that it falls into an error state if the stack reaches a certain size. The proof uses a standard forward simulation modified to account for failure: if the source execution takes a step from a configuration to another configuration emitting some event (which can be a silent event), then there are two possibilities for a related target configuration: either (i) it can take some steps to another configuration related to the second source configuration and emit the same event (as in a standard simulation); or (ii) it can take some steps to an error state without emitting any events. The latter corresponds to the case of

a resource exhaustion error: the target execution can terminate early, producing only a prefix of the source execution trace, as allowed by the relation.                                                                          ⊡

We conclude this subsection by noting that the resource exhaustion relation and the undefined behavior relation from the previous subsection can easily be combined. Indeed, given a relation $\sim_{UB}$ and a relation $\sim_{RE}$ defined as above on the same sets of traces, we can build a new relation $\sim$ that allows both refinement of undefined behavior and resource exhaustion by taking their union: $s \sim t \equiv s \sim_{UB} t \lor s \sim_{RE} t$. A compiler that is $CC^{\sim_{UB}}$ or $CC^{\sim_{RE}}$ is trivially $CC^{\sim}$, though the converse is not true.

## 4.3 Different Source and Target Values

This section first presents the common language formalisation (§4.3.1) that the following (§4.3.2) and later instances (§4.4 and §7.1) build upon. This shared language formalisation does not contain a key language feature, namely the expressions that generate actions and thus labels. This is because each instance deals with specific ways to generate actions, so each instance will define its own extension to each of the languages defined below. Additionally, each instance will define its own compiler and the trace relation used to attain $CC^{\sim}$.

*4.3.1  Shared Source and Target Language Formalisation.* The source language is a pure, statically typed expression language whose expressions e include naturals, booleans, a boolean conditional and a conditional for expressions that reduce to $0$, arithmetic and relational operations and sequencing.

$$e ::= n \mid b \mid \text{if } e \text{ then } e \text{ else } e \mid \text{ifz } e \text{ then } e \text{ else } e \mid e \text{ op } e \mid e; e'$$

$$\text{op} ::= + \mid \times \mid \leq \mid == \qquad \text{ty} ::= B \mid N$$

Types ty are either $N$ (naturals) or $B$ (booleans) and typing is standard.

$$
\frac{}{\vdash n : N} \text{(Type-nat)}
\qquad
\frac{}{\vdash b : B} \text{(Type-bool)}
\qquad
\frac{\vdash e_1 : N \quad \vdash e_2 : N \quad \cdot = + \text{ or } \times}{\vdash e_1 \cdot e_2 : N} \text{(Type-plus-times)}
\qquad
\frac{\vdash e_1 : N \quad \vdash e_2 : N}{\vdash e_1 \leq e_2 : B} \text{(Type-le)}
$$

$$
\frac{\vdash e_1 : B \quad \vdash e_2 : ty \quad \vdash e_3 : ty}{\vdash \text{if } e_1 \text{ then } e_2 \text{ else } e_3 : ty} \text{(Type-ite)}
\qquad
\frac{\vdash e_1 : N \quad \vdash e_2 : ty \quad \vdash e_3 : ty}{\vdash \text{ifz } e_1 \text{ then } e_2 \text{ else } e_3 : ty} \text{(Type-izte)}
$$

The language semantics deal with actions i, lists of actions is and expression results r. A list of actions is is a list of individual actions i, which are instance-dependant and thus presented later; the same holds for source traces s.

$$r ::= n \mid b \qquad\qquad i, s ::= \text{instance-specific} \qquad\qquad \text{is} ::= i \cdot \text{is} \mid \varnothing$$

The source language has a standard big-step operational semantics ($e \rightsquigarrow \langle \text{is}, r \rangle$) which tells how an expression e generates a list of actions and a result $\langle \text{is}, r \rangle$.

$$
\frac{}{n \rightsquigarrow \langle \varnothing, n \rangle} \text{(Sem-nat)}
\qquad
\frac{}{b \rightsquigarrow \langle \varnothing, b \rangle} \text{(Sem-bool)}
\qquad
\frac{e_1 \rightsquigarrow \langle \text{is}_1, n_1 \rangle \quad e_2 \rightsquigarrow \langle \text{is}_2, n_2 \rangle \quad \text{op} \in \{+, \times\}}{e_1 \text{ op } e_2 \rightsquigarrow \langle \text{is}_1 \cdot \text{is}_2, (n_1 \text{ op } n_2) \rangle} \text{(Sem-op-nat)}
$$

$$
\frac{e_1 \rightsquigarrow \langle \text{is}_1, n_1 \rangle \quad e_2 \rightsquigarrow \langle \text{is}_2, n_2 \rangle}{e_1 \leq e_2 \rightsquigarrow \langle \text{is}_1 \cdot \text{is}_2, (n_1 \leq n_2) \rangle} \text{(Sem-le)}
\qquad
\frac{e \rightsquigarrow \langle \text{is}, b \rangle \quad b ? i = 1 : i = 2 \quad e_i \rightsquigarrow \langle \text{is}_i, r_i \rangle}{\text{if } e \text{ then } e_1 \text{ else } e_2 \rightsquigarrow \langle \text{is} \cdot \text{is}_i, r_i \rangle} \text{(Sem-ite)}
$$

$$
\frac{e \rightsquigarrow \langle \text{is}, n \rangle \quad n == 0 ? i = 1 : i = 2 \quad e_i \rightsquigarrow \langle \text{is}_i, r_i \rangle}{\text{ifz } e \text{ then } e_1 \text{ else } e_2 \rightsquigarrow \langle \text{is} \cdot \text{is}_i, r_i \rangle} \text{(Sem-izte)}
\qquad
\frac{e \rightsquigarrow \langle \text{is}, r \rangle \quad e' \rightsquigarrow \langle \text{is}', r' \rangle}{e; e' \rightsquigarrow \langle \text{is} \cdot \text{is}', r' \rangle} \text{(Sem-seq)}
$$

The target language is analogous to the source one, except that it is untyped, it only has naturals **n** and its only conditional is **ifz e then e else e**.

$$e ::= \mathbf{n} \mid \mathbf{e\ op\ e} \mid \mathbf{ifz\ e\ then\ e\ else\ e} \mid \mathbf{e; e'} \qquad \mathbf{op} ::= \mathbf{+} \mid \mathbf{\times} \qquad \mathbf{r} ::= \mathbf{n}$$

$$\mathbf{i, t} ::= \text{instance-specific} \qquad\qquad \mathbf{is} ::= \mathbf{i \cdot is} \mid \varnothing$$

The semantics of the target language is also given in big-step style; since its rules are a subset of the source rules, they are omitted. Since we only have naturals and all expressions operate on them, no error result is possible in the target.

*4.3.2  Different Source and Target Values.* In this instance, we extend the source language with expressions to perform booleans and natural inputs, while the target only has expressions to input naturals. To compile the $\leq$, the target is also extended with a conditional that checks if an expression is less than another.

$$e ::= \cdots \mid \text{in-b} \mid \text{in-n} \qquad\qquad i ::= n \mid b \qquad\qquad s ::= \langle is, r \rangle$$

$$e ::= \cdots \mid \text{in-n} \mid \text{if e} \leq \text{e then e else e} \qquad i ::= n \qquad\qquad t ::= \langle is, r \rangle$$

Source actions are boolean **b** and natural inputs **n** and source traces **s** are lists of actions **is** together with a final result **r**. Target actions are just natural inputs **n**.

The source extensions respect typing and thus well-typed programs never produce error (🦤). The semantics of the extensions adds elements to the traces.

$$\frac{}{\vdash \text{in-b} : B} \text{(Type-in-b)} \qquad \frac{}{\vdash \text{in-n} : N} \text{(Type-in-n)} \qquad \Bigg| \qquad \frac{}{\text{in-n} \rightsquigarrow \langle n \cdot \varnothing, n \rangle} \text{(Sem-in-nat)} \qquad \frac{}{\text{in-b} \rightsquigarrow \langle b \cdot \varnothing, b \rangle} \text{(Sem-in-bool)}$$

$$\frac{e_1 \rightsquigarrow \langle is_1, n_1 \rangle \qquad e_2 \rightsquigarrow \langle is_2, n_2 \rangle \qquad n_1 \leq n_2 ? i = 3 : i = 4 \qquad e_i \rightsquigarrow \langle is_i, n_i \rangle}{\text{if } e_1 \leq e_2 \text{ then } e_3 \text{ else } e_4 \rightsquigarrow \langle is_1 \cdot is_2 \cdot is_i, n_i \rangle} \text{(Sem-itele)}$$

The compiler is homomorphic, translating a source expression to the same target expression; the only differences are natural numbers (and conditionals).

$$n{\downarrow} = \mathbf{n} \qquad\qquad \text{true}{\downarrow} = \mathbf{1} \qquad\qquad\qquad e_1 + e_2{\downarrow} = e_1{\downarrow}\mathbf{+}e_2{\downarrow}$$

$$\text{in-n}{\downarrow} = \text{in-n} \qquad \text{false}{\downarrow} = \mathbf{0} \qquad\qquad\qquad e_1 \leq e_2{\downarrow} = \mathbf{if}\ e_1{\downarrow} \leq e_2{\downarrow}\ \mathbf{then\ 1\ else\ 0}$$

$$\text{in-b}{\downarrow} = \text{in-n} \qquad e_1 \times e_2{\downarrow} = e_1{\downarrow}\mathbf{\times}e_2{\downarrow} \qquad \text{if } e_1 \text{ then } e_2 \text{ else } e_3{\downarrow} = \mathbf{ifz}\ e_1{\downarrow}\ \mathbf{then}\ e_3{\downarrow}\ \mathbf{else}\ e_2{\downarrow}$$

$$e; e'{\downarrow} = e{\downarrow}; e'{\downarrow} \qquad\qquad\qquad \text{ifz } e_1 \text{ then } e_2 \text{ else } e_3{\downarrow} = \mathbf{ifz}\ e_1{\downarrow}\ \mathbf{then}\ e_2{\downarrow}\ \mathbf{else}\ e_3{\downarrow}$$

When compiling an *if-then-else* the *then* and *else* branches of the source are swapped in the target because of the compilation of booleans.

**Relating Traces.** We relate basic values (naturals and booleans) in a non-injective fashion as noted below. Then, we extend the relation to lists of inputs pointwise (Rules Empty and Cons) and lift that relation to traces (Rules Nat and Bool).

$$n \sim \mathbf{n} \qquad\qquad \text{true} \sim \mathbf{n} \quad \text{if } \mathbf{n} > 0 \qquad\qquad \text{false} \sim \mathbf{0}$$

$$\frac{}{\varnothing \sim \varnothing} \text{(Empty)} \qquad \frac{i \sim \mathbf{i} \qquad is \sim \mathbf{is}}{i \cdot is \sim \mathbf{i \cdot is}} \text{(Cons)} \qquad \Bigg| \qquad \frac{is \sim \mathbf{is} \qquad n \sim \mathbf{n}}{\langle is, n \rangle \sim \langle \mathbf{is, n} \rangle} \text{(Nat)} \qquad \frac{is \sim \mathbf{is} \qquad b \sim \mathbf{n}}{\langle is, b \rangle \sim \langle \mathbf{is, n} \rangle} \text{(Bool)}$$

**Property mappings.** The property mappings $\tilde{\sigma}$ and $\tilde{\tau}$ induced by the trace relation $\sim$ defined above capture the intuition behind encoding booleans as naturals:
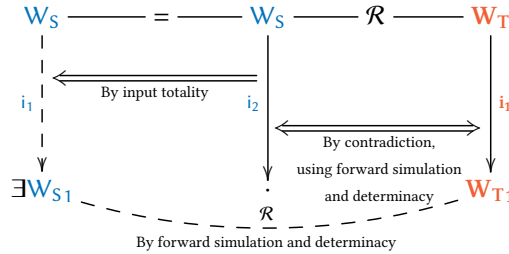- the source-to-target mapping allows true to be encoded by any non-zero number;
- the target-to-source mapping requires that **0** be replaceable by *both* **0** and false.

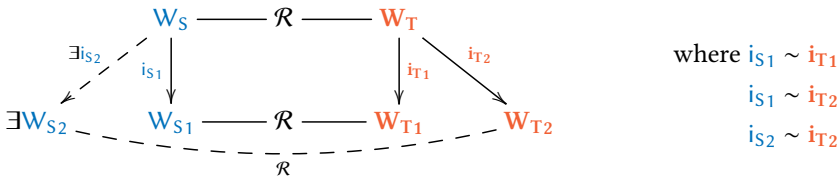**Compiler correctness.** With the relation above, the compiler is proven to satisfy CC~.

**Theorem 4.3** ( $\cdot\downarrow$ is correct 🐦 ).   $\cdot\downarrow$ is CC~.

**Simulations with different traces.** In the settings where $\mathsf{Traces_S} = \boldsymbol{Trace_T}$, it is customary to prove compiler correctness showing a forward simulation (i.e., a simulation between source and target transition system); then, using determinacy [? ? ] of the target language and input totality [? ? ] (receptiveness) of the source, this forward simulation is flipped into a backward simulation (a simulation between target and source transition system), as described by [? ? ]. This *"flipping"* is useful because forward simulations are often much easier to prove (by induction on the transitions of the source) than backward ones. For the proof of Theorem 4.3 we had to show a *backward* simulation as it was not possible to define a forward one and then flip it. Hereafter we show the reason lies in the shape of trace relation itself and disccus when is possible to generalize the flipping to the trace-relating setting.

We first give the main idea of the flipping proof, when the inputs are the same in the source and the target [? ? ]. We only consider inputs, as it is the most interesting case, since with determinacy, nondeterminism only occurs on inputs. Given a forward simulation $\mathcal{R}$, and a target program $\mathbf{W_T}$ that simulates a source program $\mathsf{W_S}$, $\mathbf{W_T}$ is able to perform an input iff so is $\mathsf{W_S}$: otherwise, say for instance that $\mathsf{W_S}$ performs an output, by forward simulation $\mathbf{W_T}$ would also perform an output, which is impossible because of determinacy. By input totality of the source, $\mathsf{W_S}$ must be able to perform the exact same input as $\mathbf{W_T}$; using forward simulation and determinacy, the resulting programs must be related.



The trace relation from §4.3.2 is not injective (both $0$ and false are mapped to $0$), therefore these arguments do not apply: not all possible inputs of target programs are accounted for in the forward simulation. In order to flip a forward simulation into a backward one it's necessary that, for any source program $\mathsf{W_S}$ and target program $\mathbf{W_T}$ related by the forward simulation $\mathcal{R}$, the following diagram is satisfied



We say that a forward simulation for which this property holds is *flippable*. For our example compiler, a flippable forward simulation works as follows: whenever a boolean input occurs in the source, the target program must perform every strictly positive input $\mathbf{n}$ (and not just $\mathbf{1}$, as suggested by the compiler). Using this property, determinacy of the target, input totality of the source, as well as the fact that any target input has an inverse image through the relation, we can indeed show that the forward simulation can be turned into a backward one: starting from $\mathsf{W_S}\ \mathcal{R}\ \mathbf{W_T}$ and an input $\mathbf{i_{T2}}$, we show that there is $\mathsf{i_{S1}}$ and $\mathbf{i_{T2}}$ as in the diagram above, using the same arguments as

when the inputs are the same; because the simulation is flippable, we can close the diagram, and obtain the existence of an adequate $i_{S_2}$. From this we obtain CC~.

In fact we showed that the flippable hypothesis is also sufficient to flip a forward simulation into a backward one, even in the trace-relating setting, and proved it in a general (i.e., language independent) 'flipping theorem' (✎). We have also shown that if the relation ~ defines a bijection between the inputs of the source and the target, then any forward simulation is flippable, hence reobtaining the usual proof technique [? ?] as a special case.

### 4.4 Abstraction Mismatches

We now consider how to relate traces where a single source action is compiled to multiple target ones. To illustrate this, we extend our source language to output (nested) pairs of arbitrary size, and our target language to send values that have a fixed size. Concretely, the source is analogous to the language of §4.3, except that it does not have inputs (nor booleans for simplicity) but it has pairs. Additionally, it has an expression send e which can emit a (nested) pair e of values in a single action. Given that e reduces to a pair, e.g., ⟨v1, ⟨v2, v3⟩⟩, expression send ⟨v1, ⟨v2, v3⟩⟩ emits action ⟨v1, ⟨v2, v3⟩⟩. That expression is eventually compiled into a sequence of individual sends in the target language **send v1 ; send v2 ; send v3**, since in the target, **send e** sends the value that **e** reduces to, but the language cannot send pairs (although it has pair constructs).

The source and target languages are formally extended (resp. in the first and second lines below) with pairs and sending constructs as follows. For reasons that we explain when the compiler is presented, we extend the target language with a let-in construct and variables. Finally, source traces are sequences of sent values i (which include nested pairs) and target traces are only sequences of natural numbers.

$$e ::= \cdots \mid \langle e, e \rangle \mid e.1 \mid e.2 \mid \text{send } e \qquad ty ::= N \mid ty \times ty \qquad i ::= n \mid \langle i, i \rangle \qquad s ::= is$$

$$e ::= \cdots \mid \langle e, e \rangle \mid e.1 \mid e.2 \mid \text{let } x = e \text{ in } e \mid x \mid \text{send } e \qquad i ::= n \qquad t ::= is$$

The source additions are well-typed and their semantics is unsurprising; the semantics relies on the usual capture-avoiding substitution $[r/x]$ of a result $r$ for a variable $x$

| (Type-send) | (Type-pair) | (Type-p1) | (Type-p2) | (Type-send) |
|---|---|---|---|---|
| $\dfrac{\vdash e : \tau \times \tau'}{\vdash \text{send } e}$ | $\dfrac{\vdash e : \tau \quad \vdash e' : \tau'}{\vdash \langle e, e' \rangle : \tau \times \tau'}$ | $\dfrac{\vdash e : \tau \times \tau'}{\vdash e.1 : \tau}$ | $\dfrac{\vdash e : \tau \times \tau'}{\vdash e.2 : \tau'}$ | $\dfrac{\vdash e : N}{\vdash \text{send } e}$ |

| (Eval-P1) | (Eval-P2) | (Eval-Pair) |
|---|---|---|
| $\dfrac{e \rightsquigarrow \langle is, \langle r_1, r_2 \rangle \rangle}{e.1 \rightsquigarrow \langle is, r_1 \rangle}$ | $\dfrac{e \rightsquigarrow \langle is, \langle r_1, r_2 \rangle \rangle}{e.1 \rightsquigarrow \langle is, r_2 \rangle}$ | $\dfrac{e \rightsquigarrow \langle is, r \rangle \quad e' \rightsquigarrow \langle is', r' \rangle}{\langle e, e' \rangle \rightsquigarrow \langle is \cdot is', \langle r, r' \rangle \rangle}$ |

| (Eval-Send) | (Sem) | (Eval-letin) |
|---|---|---|
| $\dfrac{e \rightsquigarrow \langle is, r \rangle}{\text{send } e \rightsquigarrow \langle is \cdot r, r \rangle}$ | $\dfrac{e \rightsquigarrow \langle is, r \rangle}{e \rightsquigarrow is}$ | $\dfrac{e \rightsquigarrow \langle is, r \rangle \quad e'[r/x] \rightsquigarrow \langle is', r' \rangle}{\text{let } x = e \text{ in } e' \rightsquigarrow \langle is \cdot is', r' \rangle}$ |

The compiler is defined inductively on the type derivation of a source expression $(\cdot \downarrow \; : \vdash e : \tau \rightarrow e)$. The only interesting case is when compiling a send e, where we use the source type information concerning the message (i.e., a pair) being sent to deconstruct that pair into a sequence of natural numbers, which is what is sent in the target. This is the reason we need the let-in construct in the target, since we run the pair once (as the argument of the let-in) and then we send all of its projection, to avoid duplicating side effects. Technically, since it is defined on the type derivations of terms, the compiler is defined inductively on type derivations (and not simply on terms). Thus,

compiling $e; e'$ would look like the following (using $D$ as a metavariable to range over derivations).

$$\left( \frac{\dfrac{D}{\vdash e} \quad \dfrac{D'}{\vdash e'}}{\vdash e; e'} \right)\Big\downarrow = \left( \frac{D}{\vdash e} \right)\Big\downarrow; \left( \frac{D'}{\vdash e'} \right)\Big\downarrow$$

However, note that each judgment uniquely identifies which typing rule is being applied and the underlying derivation. Thus, for compactness, we only write the judgment in the compilation and implicitly apply the related typing rule to obtain the underlying judgments for recursive calls. To differentiate this from the compiler of Section 4.3.2, this compiler has parentheses over its input.

$$(\vdash n : N)\downarrow = n \qquad\qquad (\vdash e.1 : \tau)\downarrow = (\vdash e : \tau \times \tau')\downarrow.1$$

$$(\vdash e \oplus e' : N)\downarrow = (\vdash e : N)\downarrow \oplus (e' : N)\downarrow \qquad (\vdash e.2 : \tau')\downarrow = (\vdash e : \tau \times \tau')\downarrow.2$$

$$(\vdash \langle e, e' \rangle : \tau \times \tau')\downarrow = \langle (\vdash e : \tau)\downarrow, (\vdash e' : \tau')\downarrow \rangle$$

$$\left( \vdash \begin{array}{c} \text{if } e \\ \text{then } e \text{ else } e' \end{array} \right)\Big\downarrow = \begin{array}{c} \text{if } (\vdash e : N)\downarrow \\ \text{then } (\vdash e)\downarrow \text{ else } (\vdash e')\downarrow \end{array} \qquad (\vdash \text{send } e)\downarrow = \begin{array}{c} \text{let } x = (\vdash e : \tau \times \tau')\downarrow \\ \text{in gensend } (x, \tau \times \tau') \end{array}$$

$$\left( \vdash \begin{array}{c} \text{ifz } e \\ \text{then } e \text{ else } e' \end{array} \right)\Big\downarrow = \begin{array}{c} \text{ifz } (\vdash e : N)\downarrow \\ \text{then } (\vdash e)\downarrow \text{ else } (\vdash e')\downarrow \end{array} \qquad (\vdash e; e')\downarrow = (\vdash e)\downarrow; (\vdash e')\downarrow$$

$$\text{gensend } (x, \tau) = \begin{cases} \text{send } x & \text{if } \tau = N \\ \text{gensend } (x, \tau').1; \text{gensend } (x, \tau'').2 & \text{if } \tau = \tau' \times \tau'' \end{cases}$$

**Relating Traces.** We start with the trivial relation between numbers: $n \sim^0 n$, i.e., numbers are related when they are the same. We cannot build a relation between single actions since a single source action is related to multiple target ones. Therefore, we define a relation between a source action $i$ and a target trace $t$ (a list of numbers), inductively on the structure of $i$.

$$\begin{array}{cccc} \text{(Trace-Rel-N-N)} & \text{(Trace-Rel-N-M)} & \text{(Trace-Rel-M-N)} & \text{(Trace-Rel-M-M)} \\[4pt] \dfrac{n \sim^0 n \quad n' \sim^0 n'}{\langle n, n' \rangle \sim n \cdot n'} & \dfrac{n \sim^0 n \quad i \sim t}{\langle n, i \rangle \sim n \cdot t} & \dfrac{i \sim t \quad n \sim^0 n}{\langle i, n \rangle \sim t \cdot n} & \dfrac{i \sim t \quad i' \sim t'}{\langle i, i' \rangle \sim t \cdot t'} \end{array}$$

A pair of naturals is related to the two actions that send each element of the pair (Rule Trace-Rel-N-N). If a pair is made of sub-pairs, we require all such sub-pairs to be related (Rules Trace-Rel-N-M to Trace-Rel-M-M).

We build on these rules to define the $s \sim t$ relation between source and target traces for which the compiler is correct (Theorem 4.5). Trivially, traces are related when they are both empty. Alternatively, given related traces, we can concatenate a source action and a second target trace provided that they are related (Rule Trace-Rel-Single). Before proving that the compiler is correct we need Lemma 4.4. Intuitively, that lemma tells us that the way we break down a source sent value $r$ into multiple target sends is correct.

$$\begin{array}{c} \text{(Trace-Rel-Single)} \\[4pt] \dfrac{s \sim t \quad i \sim t'}{s \cdot i \sim t \cdot t'} \end{array}$$

**Lemma 4.4** (gensend $(\cdot, \cdot)$ works). *if* gensend $(x, \tau \times \tau')[(\vdash r : \tau \times \tau')\downarrow/x] \rightsquigarrow t$ *then* $r \sim t$ *(since $r$ is necessarily a sent value $i$, that can be related to $t$).*

**Theorem 4.5** $((\cdot)\downarrow$ is correct). $(\cdot)\downarrow$ is $CC^\sim$.

With our trace relation, the trace property mappings capture the following intuitions:
- The target-to-source mapping states that a source property can reconstruct target action as it sees fit. For example, trace $4 \cdot 6 \cdot 5 \cdot 7$ is related to $\langle 4, 6 \rangle \cdot \langle 5, 7 \rangle$ and $\langle \langle 4, \langle 6, \langle 5, 7 \rangle \rangle \rangle \rangle$ (and

many more variations). This gives freedom to the source implementation of a target behavior, which follows from the non-injectivity of $\sim$.[10]

- The source-to-target mapping "forgets" about the way pairs are nested, but is faithful w.r.t. the values $v_i$ contained in a message. Notice that source safety properties are always mapped to target safety properties. For instance, if $\pi_S \in \mathsf{Safety}_S$ prescribes that some bad number is never sent, then $\tilde{\tau}(\pi_S)$ prescribes the same number is never sent in the target and $\tilde{\tau}(\pi_S) \in \mathsf{Safety_T}$. Of course if $\pi_S \in \mathsf{Safety}_S$ prescribes that a particular nested pairing like $\langle 4, \langle 6, \langle 5, 7 \rangle \rangle \rangle$ never happens, then $\tilde{\tau}(\pi_S)$ is still a target safety property, but the trivial one, since $\tilde{\tau}(\pi_S) = \top \in \mathsf{Safety_T}$.

## 5 TRACE-RELATING COMPILATION AND NONINTERFERENCE PRESERVATION

We now study the relation between trace-relating compilation and noninterference preservation. As mentioned earlier (§3.1), in the particular case where source and target observations are drawn from the same set, a correct compiler ($CC^=$) is enough to ensure the preservation of all subset-closed hyperproperties, in particular of *noninterference* (NI) [? ]. But in the scenario where target observations are strictly more informative than source observations, this is not the case. In fact, as we will show, the best guarantee one may expect from a correct trace-relating compiler ($CC^\sim$) in such a setting is a *weakening* (or *declassification*) of target noninterference that matches the noninterference property satisfied in the source. In certain scenarios, it turns out that the noninterference property of interest in the target comes "for free", while in others, it does not, and therefore establishing noninterference requires an additional proof effort beyond $CC^\sim$. To formalize this reasoning, this section applies the trinitarian view of trace-relating compilation to the general framework of abstract noninterference (ANI) [? ], clarifying the kind of noninterference preservation that follows from a given trace relation and correct compilation.

We first define NI and explain the issue of preserving source NI via a $CC^\sim$ compiler (§5.1). We then introduce ANI, which allows characterizing various forms of noninterference (§5.2), and formulate a theory of ANI preservation via $CC^\sim$, both with respect to a *timing insensitive* declassification (§5.3) and in general (§5.4). We also study how to deal with cases such as undefined behavior in the target (§5.5). We then answer the dual question, i.e., which source NI should be satisfied to guarantee that compiled programs are noninterfering with respect to target observers (§5.6). Finally, we use this formal development to analyze recent work on correct compilers with interesting noninterference guarantees [? ? ], clarifying whether these guarantees follow from correctness alone or not (§5.7).

### 5.1 Noninterference and Trace-Relating Compilation

Intuitively, noninterference (NI) requires that publicly observable outputs do not reveal information about private inputs. To define this formally, we need a few additions to our setup. We indicate the (disjoint) *input* and *output* projections of a trace $t$ as $t^\circ$ and $t^\bullet$ respectively.[11] Denote with $[t]_{low}$ the equivalence class of a trace $t$, obtained using a standard low-equivalence relation that relates low (public) events only if they are equal, and ignores any difference between private events. Then, NI for source traces can be defined as:

$$\mathsf{NI_S} = \{\pi_S \mid \forall \mathsf{s_1 s_2} \in \pi_S. \ [\mathsf{s_1^\circ}]_{low} = [\mathsf{s_2^\circ}]_{low} \Rightarrow [\mathsf{s_1^\bullet}]_{low} = [\mathsf{s_2^\bullet}]_{low} \}$$

---

[10]Making $\sim$ injective is a matter of adding open and close parenthesis actions in target traces.

[11]The exact shape of inputs and outputs depends on the scenario. For instance, inputs can be initial memories and outputs trace semantics of programs as in [? , Section 7], while for interactive programs one may want to consider streams like ? ]. We only require the sets of input and output projections to be disjoint. Further information, such as the ordering of events, is part of the attacker/observer model or the declassification of noninterference itself.

That is, source NI comprises the sets of traces that have equivalent low output projections as long as their low input projections are equivalent.

When additional observations are possible in the target, it is unclear whether a noninterfering source program is compiled to a noninterfering target program or not, and if so, whether the notion of NI in the target is the expected (or desired) one. We illustrate this issue by considering a scenario where target traces extend source traces by exposing the execution time. While source noninterference $NI_S$ requires that private inputs do not affect public outputs, $NI_T$ additionally requires that the execution time is not affected by varying private inputs.

To model the scenario described, we represent target traces as pairs of a source trace and a natural number that denotes the time spent to produce the trace (using $\omega$ for infinite time units). Formally, if $Trace_S$ denotes the set of source traces, then $Trace_T = Trace_S \times \mathbb{N}^\omega$ is the set of target traces, where $\mathbb{N}^\omega \triangleq \mathbb{N} \cup \{\omega\}$.

Notice that if two source traces $s_1, s_2$ are low-equivalent then $\{s_1, s_2\} \in NI_S$ and $\{(s_1, 42), (s_1, 42)\} \in NI_T$, but $\{(s_1, 42), (s_2, 43)\} \notin NI_T$ and $\{(s_1, 42), (s_2, 42), (s_1, 43), (s_2, 43)\} \notin NI_T$.

Consider the following straightforward trace relation, which relates a source trace to any target trace whose first component is equal to it, irrespective of execution time:

$$s \sim t \equiv \exists n.\ t = (s, n)$$

A compiler is $CC^\sim$ for this trace relation if any trace that can be exhibited in the target can be simulated in the source in some amount of time. For such a compiler Theorem 3.3 says that if $W$ satisfies $NI_S$, then $W\!\!\downarrow$ satisfies $Cl_\subseteq \circ \tilde{\tau}(NI_S)$. This hyperproperty is however strictly weaker than $NI_T$, as it contains for example $\{(s_1, 42), (s_2, 42), (s_1, 43), (s_2, 43)\}$, and one cannot conclude that $W\!\!\downarrow$ is noninterfering in the target. It is easy to check that

$$Cl_\subseteq \circ \tilde{\tau}(NI_S) = Cl_\subseteq (\{\ \pi_S \times \mathbb{N}^\omega\ \mid\ \pi_S \in NI_S\}) = \{\ \pi_S \times \mathcal{I}\ \mid\ \pi_S \in NI_S \wedge \mathcal{I} \subseteq \mathbb{N}^\omega\},$$

the first equality coming from $\tilde{\tau}(\pi_S) = \pi_S \times \mathbb{N}^\omega$, and the second from $NI_S$ being subset-closed. As we will see, this hyperproperty *can* be characterized as a form of NI, which one might call *timing-insensitive noninterference*, i.e., ensured only against attackers that cannot measure execution time. For this characterization, and to describe different forms of noninterference as well as formally analyze their preservation by a $CC^\sim$ compiler, we rely on the general framework of *abstract noninterference* [? ].

## 5.2 Abstract Noninterference

Abstract noninterference (ANI) [? ] is a generalization of NI whose formulation relies on *abstractions* (in the sense of Abstract Interpretation [? ]) in order to encompass arbitrary variants of NI. ANI is parameterized by an *observer abstraction* $\rho$, which denotes the distinguishing power of the attacker, and a *selection abstraction* $\phi$, which specifies when to check NI, and therefore captures a form of declassification [? ].[12] Formally:

$$ANI_\phi^\rho = \{\pi \mid \forall t_1 t_2 \in \pi.\ \phi(t_1^\circ) = \phi(t_2^\circ) \Rightarrow \rho(t_1^\bullet) = \rho(t_2^\bullet)\}$$

By picking $\phi = \rho = [\cdot]_{low}$, we recover the standard noninterference defined above, where NI must hold for all low inputs (i.e., no declassification of private inputs), and the observational power of the attacker is limited to distinguishing low outputs. The observational power of the attacker can be weakened by choosing a more liberal relation for $\rho$. For instance, one may limit the attacker to observe the *parity* of output integer values. Another way to weaken ANI is to use $\phi$ to specify that noninterference is only required to hold for a subset of low inputs.

---

[12]To be precise, the original formulation of ANI by [? ] includes a third parameter $\eta$, which describes the maximal input variation that the attacker may control. Here we omit $\eta$ (i.e., take it to be the identity) in order to simplify the presentation.

The operators $\phi$ and $\rho$ are defined over sets of (input and output projections of) traces, explicitly $\phi : 2^{Trace^{\circ}} \to 2^{Trace^{\circ}}$ and $\rho : 2^{Trace^{\bullet}} \to 2^{Trace^{\bullet}}$. When we write $\phi(t)$ like above, this should be understood as a convenience notation for $\phi(\{t\})$. Likewise, $\phi = [\cdot]_{low}$ should be understood as $\phi = \lambda\pi. \bigcup_{t \in \pi}[t]_{low}$, i.e., the powerset lifting of $[\cdot]_{low}$. Additionally, $\phi$ and $\rho$ are required to be upper-closed operators (*uco*)—i.e., monotonic, idempotent and extensive (i.e., $\forall\pi^{\bullet}. \pi^{\bullet} \subseteq \rho(\pi^{\bullet})$) —on the poset that is the powerset of (input and output projections of) traces ordered by inclusion [? ].

## 5.3 Trace-Relating Compilation and ANI for Timing

We can now reformulate our example with observable execution times in target traces in terms of ANI. We have $\mathsf{NI_S} = ANI_{\rho}^{\phi}$ with $\phi = \rho = [\cdot]_{low}$. In this case, the hyperproperty that a compiled program $\mathsf{W}{\downarrow}$ satisfies whenever $\mathsf{W}$ satisfies $\mathsf{NI_S}$ can be described as an instance of ANI:

$$Cl_{\subseteq} \circ \tilde{\tau}(\mathsf{NI_S}) = ANI_{\phi}^{\rho}$$
$$\text{for } \phi = \phi \text{ and } \rho(\pi) = \{(\mathsf{s}, \mathsf{n}) \mid \exists(\mathsf{s_1}, \mathsf{n_1}) \in \pi.\ [\mathsf{s}^{\bullet}]_{low} = [\mathsf{s_1^{\bullet}}]_{low}\}$$

The definition of $\phi$ tells us that the trace relation does not affect the selection abstraction, i.e., declassification is unaffected. The definition of $\rho$ characterizes an observer that cannot distinguish execution times for noninterfering traces (notice that $\mathsf{n_1}$ in the definition of $\rho$ is discarded). For instance, $\rho(\{(\mathsf{s}, \mathsf{n_1})\}) = \rho(\{(\mathsf{s}, \mathsf{n_2})\})$, for any $\mathsf{s}, \mathsf{n_1}, \mathsf{n_2}$. Therefore, in this setting, we know explicitly through $\rho$ that a CC$^{\sim}$ compiler degrades source noninterference to target *timing-insensitive* noninterference.

## 5.4 Trace-Relating Compilation and ANI in General

While the particular $\phi$ and $\rho$ above can be discovered by intuition, we want to know whether there is a systematic way of obtaining them in general. In other words, for *any* trace relation $\sim$ and *any* notion of source NI, what property is guaranteed on noninterfering source programs by any CC$^{\sim}$ compiler?

We can now answer this question generally (Theorem 5.1): any source notion of noninterference expressible as an instance of ANI is mapped to a corresponding instance of ANI in the target, whenever source traces are an abstraction of target ones (i.e., when $\sim$ is a total and surjective map). For this result we consider trace relations that can be split into input and output trace relations (denoted as $\sim \triangleq \langle\overset{\circ}{\sim}, \overset{\bullet}{\sim}\rangle$) such that $\mathsf{s} \sim \mathsf{t} \iff \mathsf{s}^{\circ} \overset{\circ}{\sim} \mathsf{t}^{\circ} \wedge \mathsf{s}^{\bullet} \overset{\bullet}{\sim} \mathsf{t}^{\bullet}$. The trace relation $\sim$ corresponds to a Galois connection between the sets of trace properties $\tilde{\tau} \leftrightarrows \tilde{\sigma}$ as described in §2.2. Similarly, the pair $\overset{\circ}{\sim}$ and $\overset{\bullet}{\sim}$ corresponds to a pair of Galois connections, $\tilde{\tau}^{\circ} \leftrightarrows \tilde{\sigma}^{\circ}$ and $\tilde{\tau}^{\bullet} \leftrightarrows \tilde{\sigma}^{\bullet}$, between the sets of input and output properties. In the timing example, time is an output so we have $\sim \triangleq \langle=, \overset{\bullet}{\sim}\rangle$ and $\overset{\bullet}{\sim}$ is defined as $\mathsf{s}^{\bullet} \overset{\bullet}{\sim} \mathsf{t}^{\bullet} \equiv \exists\mathsf{n}.\ \mathsf{t}^{\bullet} = (\mathsf{s}^{\bullet}, \mathsf{n})$.

**Theorem 5.1** (Compiling ANI). Assume traces of source and target languages are related via $\sim \subseteq \mathsf{Traces_S} \times \mathbf{Trace_T}$, $\sim \triangleq \langle\overset{\circ}{\sim}, \overset{\bullet}{\sim}\rangle$ such that $\overset{\circ}{\sim}$ and $\overset{\bullet}{\sim}$ are both total maps from target to source traces, and $\overset{\circ}{\sim}$ is surjective. Assume $\downarrow$ is a CC$^{\sim}$ compiler, and $\phi \in uco(2^{\mathsf{Trace_S^{\circ}}})$, $\rho \in uco(2^{\mathsf{Trace_S^{\bullet}}})$. If $\mathsf{W}$ satisfies $ANI_{\phi}^{\rho}$, then $\mathsf{W}{\downarrow}$ satisfies $ANI_{\phi^{\#}}^{\rho^{\#}}$, where $\phi^{\#}$ and $\rho^{\#}$ are defined as:

$$\phi^{\#} = g^{\circ} \circ \phi \circ f^{\circ} \qquad\qquad \rho^{\#} = g^{\bullet} \circ \rho \circ f^{\bullet}$$
$$f^{\circ}(\pi^{\circ}) = \{\mathsf{s}^{\circ} \mid \exists\mathsf{t}^{\circ} \in \pi^{\circ}.\ \mathsf{s}^{\circ} \overset{\circ}{\sim} \mathsf{t}^{\circ}\} \qquad g^{\circ}(\pi_{\mathsf{S}}^{\circ}) = \left\{\mathsf{t}^{\circ} \mid \forall\mathsf{s}^{\circ}.\ \mathsf{s}^{\circ} \overset{\circ}{\sim} \mathsf{t}^{\circ} \Rightarrow \mathsf{s}^{\circ} \in \pi_{\mathsf{S}}^{\circ}\right\}$$

(and both $f^{\bullet}$ and $g^{\bullet}$ are defined analogously).

Moreover, we can prove that if $\stackrel{\sim}{\cdot}$ is surjective, then $ANI^{\rho^{\#}}_{\phi^{\#}} \subseteq Cl_{\subseteq} \circ \tilde{\tau}(ANI^{\rho}_{\phi})$. Therefore, the derived guarantee $ANI^{\rho^{\#}}_{\phi^{\#}}$ is at least as strong as the hyperproperty (a priori different from some noninterference) that follows by just knowing that the compiler $\downarrow$ is CC$^{\sim}$.

The target abstract noninterference has to be intended as the *best correct approximation* of the source one. The mappings $f^{\circ} \leftrightarrows g^{\circ}$ are the existential and universal images of the relation $\stackrel{\sim}{\cdot}_{swap} \subseteq \text{Trace}_T \times \text{Trace}_S$, defined by $t^{\circ} \stackrel{\sim}{\cdot}_{swap} s^{\circ}$ if and only if $s^{\circ} \stackrel{\sim}{\cdot} t^{\circ}$. Therefore $f^{\circ}$ and $g^{\circ}$ are lower and upper adjoints, respectively (§2). The operator $\phi^{\#}$ is the best correct approximation of $\phi$ w.r.t to $f^{\circ} \leftrightarrows g^{\circ}$ [?] (hence the choice of the (_)$^{\#}$ notation). A similar result holds for $\rho^{\#}$.

Coming back to our example above, we can formally recover the intuitively-justified definitions, i.e., $\phi^{\#} = g^{\circ} \circ \phi \circ f^{\circ} = \phi$ and $\rho^{\#} = g^{\bullet} \circ \rho \circ f^{\bullet} = \rho$.

## 5.5 Noninterference and Undefined Behavior

As stated above, Theorem 5.1 does not apply to several scenarios from §4 such as undefined behavior (§4.1). Indeed, in these cases, the relation $\stackrel{\sim}{\cdot}$ is not a total map. Nevertheless, we can still exploit our framework to reason about the impact of compilation on noninterference.

Let us consider $\sim \triangleq \langle \stackrel{\sim}{\cdot}, \stackrel{\sim}{\cdot} \rangle$ where $\stackrel{\sim}{\cdot}$ is any total and surjective map from target to source inputs (e.g., equality) and $\stackrel{\sim}{\cdot}$ is defined as $s^{\bullet} \stackrel{\sim}{\cdot} t^{\bullet} \equiv s^{\bullet} = t^{\bullet} \vee \exists m^{\bullet} \leq t^{\bullet}. s^{\bullet} = m^{\bullet} \cdot Wrong$. Intuitively, a CC$^{\sim}$ compiler guarantees noninterference for the compiled program, provided that the target attacker cannot exploit undefined behavior to learn private information. This intuition can be made formal by the following theorem.

**Theorem 5.2** (Relaxed Compiling ANI). Relax the assumptions of Theorem 5.1 by allowing $\stackrel{\sim}{\cdot}$ to be *any* output trace relation. If W satisfies $ANI^{\rho}_{\phi}$, then W$\downarrow$ satisfies $ANI^{\rho^{\#}}_{\phi^{\#}}$ where $\phi^{\#}$ is defined as in Theorem 5.1, and $\rho^{\#}$ is such that:

$$\forall s\, t.\, s^{\bullet} \stackrel{\sim}{\cdot} t^{\bullet} \Rightarrow \rho^{\#}(t^{\bullet}) = \rho^{\#}(\tilde{\tau}^{\bullet}(\rho(s^{\bullet}))) \tag{1}$$

Technically, instead of giving us a *definition* of $\rho^{\#}$, the theorem gives a *property* of it. The property states that, given a target output trace $t^{\bullet}$, the attacker cannot distinguish it from any other target output traces produced by other possible compilations ($\tilde{\tau}^{\bullet}$) of the source trace $s$ it relates to, up to the observational power of the source level attacker $\rho$. Therefore, given a source attacker $\rho$, the theorem characterizes a *family* of attackers that cannot observe any interference for a correctly compiled noninterfering program. Notice that the target attacker $\rho^{\top} \triangleq \lambda\_. \top$ satisfies the premise of the theorem, but defines a trivial hyperproperty, so that we cannot prove in general that $ANI^{\rho^{\#}}_{\phi^{\#}} \subseteq Cl_{\subseteq} \circ \tilde{\tau}(ANI^{\rho}_{\phi})$. Also, this degenerate attacker $\rho^{\top}$ shows that the family of attackers described in Theorem 5.2 is nonempty, which ensures the existence of a most powerful attacker among them [?].

## 5.6 From Target NI to Source NI

We now explore the dual question: under what hypotheses does trace-relating compiler correctness alone allow target noninterference to be reduced to source noninterference? This is of practical interest, as one would be able to protect from target attackers by ensuring noninterference in the source. This task can be made easier if the source language has some static enforcement mechanism [? ?].

Let us consider the languages from §4.4 extended with the ability to accept inputs as (pairs of) values. It is easy to show that the compiler described in §4.4 (extended to treat the new input expressions homomorphically) is still CC$^{\sim}$: given a target trace $t$ with the same inputs of the

source one (i.e., $s^\circ = t^\bullet$), the compiler of §4.4 ensures that $t$ simulates the same outputs of $s$ (i.e., $s^\bullet \mathrel{\tilde{\smile}} t^\bullet$). Assume that we want to satisfy a given notion of target noninterference after compilation, i.e., $\mathbb{W}\!\downarrow\models ANI_\phi^\rho$. Recall that the observational power of the target attacker, $\rho$, is expressed as a property of sequences of values. To express the same property (or attacker) in the source, we have to abstract the way pairs of values are nested. For instance, the source attacker should not distinguish $\langle v_1, \langle v_2, v_3\rangle\rangle$ and $\langle\langle v_1, v_2\rangle, v_3\rangle$. In general (i.e., when $\tilde{\smile}$ is not the identity), this argument is valid only when $\phi$ can be represented in the source. More precisely, $\phi$ must consider as equivalent all target inputs that are related to the same source input, because in the source it is not possible to have a finer distinction of inputs. This intuitive correspondence can be formalized as follows.

**Theorem 5.3** (Target ANI by source ANI). Let $\phi \in uco(2^{\mathrm{Trace}_\mathrm{T}^\circ})$, $\rho \in uco(2^{\mathrm{Trace}_\mathrm{T}^\bullet})$ and $\tilde{\smile}$ a total and surjective map from source outputs to target ones and assume that

$$\forall s\ t.\ s^\circ \mathrel{\tilde{\smile}} t^\circ \Rightarrow \phi(t^\circ) = \phi(\tilde{\tau}^\circ(s^\circ))$$

If $\cdot\!\downarrow$ is a $CC^\sim$ compiler and $\mathbb{W}$ satisfies $ANI_{\phi^\#}^{\rho^\#}$, then $\mathbb{W}\!\downarrow$ satisfies $ANI_\phi^\rho$ for

$$\phi^\# = \tilde{\sigma}^\circ \circ \phi \circ \tilde{\tau}^\circ \qquad\qquad\qquad\qquad \rho^\# = \tilde{\sigma}^\bullet \circ \rho \circ \tilde{\tau}^\bullet$$

## 5.7 Analyzing Noninterference Preserving Compilers

The results presented in this section formalize and generalize some intuitive facts about compiler correctness and noninterference, clarifying which noninterference property follows "for free" from trace-relating compiler correctness. Of course, in the general case, compiler correctness alone is not a strong enough criterion for dealing with many security properties [? ? ]. This section exploits our ANI-based framework and results to analyze two compilers from the recent literature [? ? ] that are both proven to be correct and to preserve two interesting notions of noninterference: cryptographic constant time (§5.7.1) and value-dependent noninterference (§5.7.2). For each, we explain how to express compiler correctness as an instance of $CC^\sim$, describe the noninterference property that is implied by the trace relation and the correctness result, and compare it with the noninterference properties of interest as established by their authors.

*5.7.1 A Correct Compiler Preserving Cryptographic Constant Time.* ? ] provide a correct compiler (as an extension of CompCert) that also preserves cryptographic constant time (CT). CT is a security property stating that the runtime of a program does not depend on its secret, and thus an attacker cannot extrude secrets of a program by observing its execution time. A CT-preserving compiler takes code that is CT and generates code that also is CT. Thus, a CT-preserving compiler must translate runtime-equivalent source programs into runtime-equivalent target ones. Notice that it is not necessary for the leakage of target programs to be the same of their source counterparts, rather: source programs with the same leakage must be compiled to target programs with the same leakage.

 ? ] prove CT preservation for seventeen passes of CompCert. The authors partition the seventeen steps in four categories depending on the proof technique they use to show CT preservation. Every category proves an instance of $CC^\sim$ by improving on the existing CompCert simulation. In three out of the four cases this is sufficient to also prove CT preservation, while for the last category a further proof is necessary. In what follows, we first encode CT as an instance of abstract noninterference, i.e., show for which operators $CT = ANI_{\phi_{CT}}^{\rho_{CT}}$ and then use our framework to understand why modifying CompCert simulation is sufficient in the first three categories but not in the last one. For each category Theorem 5.2 applies, so that no $\rho$ that respects Equation 1 can notice any interference on compiled programs that were source constant-time. In the first three categories the attacker

that defines CT – $\rho_{CT}$ – respects the equation [13] i.e.,

$$\forall s^\bullet t^\bullet. \; s^\bullet \overset{\cdot}{\sim} t^\bullet \Rightarrow \rho_{CT}(t^\bullet) = \rho_{CT}(\tilde{\tau}^\bullet(\rho_{CT}(s^\bullet))) \tag{2}$$

and CT preservation is therefore a consequence of CC$^\sim$. In the last category $\rho_{CT}$ does not respect Equation 2 and the authors have to prove an additional theorem, the *CT-diagram*.

**Trace Model and CT as an instance of ANI.** The formal definition of CT is given by extending the semantics of the languages in CompCert and enriching the traces of input and output events with leakages. Leakages are results of execution steps that involve conditional branching or memory access. A program is CT w.r.t. a certain relation over program states $\varphi$ [? , Definition 3.2] iff for every two initial states $i, i'$ such that $\varphi(i, i')$, the leakages that can be observed are the same. Notice that in [? , Definition 3.2] the secret is stored in the program states and defined by $\varphi$, therefore in order to regard CT as an instance of abstract noninterference program states will be regarded as inputs and events together with their leakages as outputs. More precisely a trace $t$ is a sequence of of triples $(i, e, j)$ where $i$ and $j$ are program states and $e$ an event in the instrumented semantics, i.e., input/output event and associated leakage.

We consider:

- $\phi_{CT}$ to be (the uco corresponding to) the relation defined by $t_1^\circ \; \phi_{CT} \; t_2^\circ$ iff $t_1^\circ, t_2^\circ$ have the same length with $t_1^\circ = (i_0, i_1), (i_1, i_2), \ldots, t_2^\circ = (j_0, j_1), (j_1, j_2), \ldots$ and $\forall n. \; \varphi(i_n, j_n)$.

- $\rho_{CT}$ to be (the uco corresponding to) the relation defined by $t_1^\bullet \; \rho_{CT} \; t_2^\bullet$ iff $t_1^\bullet, t_2^\bullet$ have the same length with $t_1^\bullet = e_0, e_1, \ldots, t_2^\bullet = f_0, f_1, \ldots$ and $\forall n. \; leak(e_n) = leak(f_n)$, where $leak(e)$ denotes the leakage in the event $e$ ( projection of $e$ on the leak-only semantics [? ]).

It is easy to check that $CT = ANI_{\phi_{CT}}^{\rho_{CT}}$ for the $\phi_{CT}$ and $\rho_{CT}$ given above.

We now present more details for each of the four proof techniques adopted by ? ]. Since CT is defined only for *safe* programs [? , Definition 3.1] we can assume no undefined behavior is ever encountered and have a simpler presentation. We also omit $\phi^\#$ coming from the application of Theorem 5.2, as it always coincides with $\phi_{CT}$. [CA: double check this fact]

**Constant-time security preservation by leakage preservation (? , Section 5.2).** For compilation passes that belong to this category, the authors prove that the source leakage is preserved exactly in the target. Thus in this simple case, the theorem proved is CC$^\sim$ where $\overset{\cdot}{\sim}$ is point-wise equality of events together with leakages, $\tilde{\tau}^\bullet$ the identity and $\rho_{CT}$ satisfies Equation 2 by idempotency of $\rho_{CT}$,

$$
\begin{array}{ll}
\rho_{CT}(\tilde{\tau}^\bullet(\rho_{CT}(s^\bullet))) = & [s^\bullet \overset{\cdot}{\sim} t^\bullet \Rightarrow s^\bullet = t^\bullet] \\
\rho_{CT}(\tilde{\tau}^\bullet(\rho_{CT}(t^\bullet))) = & [\tilde{\tau}^\bullet = \lambda \, x.x] \\
\rho_{CT}(\rho_{CT}(t^\bullet)) = & [\rho_{CT} \text{ idempotent}] \\
\rho_{CT}(t^\bullet)
\end{array}
$$

**CT preservation from leakage-erasing simulation (? , Section 5.3).** In this case, CC$^\sim$ is proved for a relation that erases source leakage-only events, i.e., those events that do not contain inputs or outputs, but only the amount of leakage revealed. More precisely (see also [? , Fig. 8]) for $s^\bullet = e_0, e_1, \ldots$ and $t^\bullet = e_0, e_1, \ldots$ of the same length, $s^\bullet \overset{\cdot}{\sim} t^\bullet$ iff

$$\forall k, e_k = e_k \lor (e_k = \epsilon \land e_k \text{ is leak only})$$

---

[13] In each compilation step source and target traces are drawn from the same set so that $\rho_{CT}$ can be applied to both source and target traces.

The property mapping associated to the above relation, $\tilde{\tau}^\bullet$, erases all leak-only events from the traces of a source property. If an attacker cannot notice at any point any difference in the leakages of two traces and we erase the leak-only events from them, the attacker will still not notice any difference on leakages, therefore it is easy to check that Equation 2 holds also in this case.

**CT preservation via memory injection (? , Section 5.4]).** This case is analogous to the one above, save that it rests on a more complex relation $\tilde{\sim}$ involving a *memory injection* relation (see ? , Definition 5.8]). [CA: Roughly speaking the first disjunct of the above $\tilde{\sim}$ would use a map $\iota$ that is indeed complex as it needs also information about the states (i.e., inputs)] Intuitively $\tilde{\sim}$ relates source and target traces that differ at most in leakages due to memory accesses. While in the previous case, leakages where simply erased, here they are modified and crucially with some uniformity. Reasoning as in the previous case, if an attacker cannot notice a difference in the leakages of two traces and we modify equal leakages of the same factor, the attacker will still not notice any difference on leakages, thus Equation 2 holds.

**CT preservation from CT-diagram (? , Section 5.5]).** In this case $\rho_{CT}$ does not satisfy Equation 2 because the *counting simulation* ([? , Definition 5.10]) does not necessarily relate source and target leakages but only the inputs and outputs[14]. $CC^\sim$ alone does not ensure that an attacker cannot observe any interference in the target leakages, in order to show preservation of CT the authors need to prove an extra condition, the so-called CT diagram [? ].

*5.7.2  Value-dependent noninterference.* ? ] introduce a compiler that provably preserves value-dependent noninterference (VDNI) for a concurrent language with shared variables. *Value-dependent* means that the secrecy level of a variable – *low* or *high* – may depend on the value of some other variable, called the control variable of the first, and therefore could change throughout its lifetime.

Preservation of VDNI for concurrent programs enjoys *compositionality*, meaning that it follows from the preservation of VDNI for each single thread [? ] under certain conditions. As the compositionality result is orthogonal to our framework, we can study either (1) the preservation of VDNI for one local thread or for (2) the whole-program,

In the remainder of this section we focus on the preservation of VDNI for a single thread, that is proven by showing a *secure refinement* relation between source and compiled threads. Similarly to the previous section, the secure refinement is expressed via a cube diagram ([? ], Figure 1), and can be proven directly [? ] or split into more obligations [? ].

As ? ] use a state transition based semantics, we first show how to encode this semantics into a trace model by defining the $\sim$ relation based on the secure refinement relation. We then show how to encode VDNI as an instance of abstract noninterference( i.e., both $\text{VDNI}_S = ANI^\rho_\phi$ and $\text{VDNI}_T = ANI^\rho_\phi$). Finally we apply Theorem 5.2 and conclude that if $W$ satisfies $\text{VDNI}_S$ then $W{\downarrow}$ satisfies $\text{VDNI}_T$ given that the trace relation $\sim$ has properties defined in [? , Theorem 5.1].

Source (WHILE) and target (RISC-like assembly) languages are equipped with a determined evaluation step semantics (i.e., a semantics where the only source of nondeterminism are external inputs, [? ], Section 2) between thread-local configurations, which are triples of the form $\langle tps, mds, mem \rangle$. In such a configuration, *mds* is the access mode state for program variables and *mem* is a map relating global program variables to their values. Both of these components are common to the source and target language. The *tps* component denotes the thread-private state. In the source language, it is the program to be executed. In the target language, *tps* consists of the target program (labelled assembly-language instructions), of a program counter and of the set of thread-local registers. We

---

[14]The interested reader will notice the difference from the previous category by comparing condition (1) of Definition 5.10 and condition (1) of Definition 5.8 by ? ].

denote WHILE configurations by tuples of the form: $\langle \text{tps}, mds, mem \rangle$ and RISC configurations by tuples of the form: $\langle \mathbf{tps}, mds, mem \rangle$.

**Trace Model and Trace relation.** We consider traces that are (possibly infinite) sequences of configurations. The traces produced by a program are the sequences of local configurations that the program may encounter during execution, according to the evaluation semantics. Let $\mathsf{s} = \langle \text{tps}_1, mds_1, mem_1 \rangle, \langle \text{tps}_2, mds_2, mem_2 \rangle \ldots$ be a source trace. The input projection is defined by $\mathsf{s}^\circ = \langle mds_1, mem_1 \rangle$ (the tuple consisting of the access modes and the memory in the first state) and the output projection is defined by $\mathsf{s}^\bullet = \mathsf{s}$ (the trace itself). Input/output projections are defined similarly for target traces.

We take the trace relation $\sim \ \subseteq \mathsf{Trace}_\mathsf{S} \times \boldsymbol{Trace}_\mathsf{T}$ to be the point-wise lifting of a secure refinement relation $\mathcal{R}$ ([? ], Definition 6). Source and target configurations $\langle \text{tps}, mds, mem \rangle \ \mathcal{R} \ \langle \mathbf{tps}, mds', mem' \rangle$ that are related coincide on the access mode and memory part (i.e., $mds = mds'$ and $mem = mem'$, ([? ], Definition 4), so that $\overset{\circ}{\sim}$ is simply the identity and $\overset{\bullet}{\sim}$ coincides with $\sim$.

**VDNI as abstract noninterference.** A program satisfies VDNI ([? ], Definition 2) if any two of its executions starting in low equivalent memories are related via a *strong low bisimulation modulo modes* (strong low bisimulation mm). Intuitively, a strong low bisimulation mm is a bisimulation that preserves low-equivalence. Preservation of VDNI is proved by ? ] by showing that for every strong low-bisimulation mm $\mathcal{B}$ for source threads, there exists a target strong low bisimulation mm $\mathcal{B}$ such that if two source threads are related by $\mathcal{B}$, then the compiled threads are related by $\mathcal{B}$ ([? ], Theorem 5.1).

The intuition for the encoding of VDNI as an instance of abstract noninterference is to model low equivalence through the operator $\phi$, and bisimilarity through $\rho$. More rigorously, $\mathrm{VDNI}_\mathsf{S} = ANI_\phi^\rho$, where $\phi$ and $\rho$ are defined as following.

For $\mathsf{s}^\circ = \langle mds_1, mem_1 \rangle$,

$$\phi(\mathsf{s}^\circ) = \left\{ \langle mds_1, mem_1' \rangle \ \middle| \ mem_1 =_{mds_1}^{Low} mem_1' \right\},$$

where $=_{mds}^{Low}$ is the low-equivalence modulo $mds$ ([? ], Definition 1).

For $\mathsf{s}^\bullet = \langle \text{tps}_1, mds_1, mem_1 \rangle, \langle \text{tps}_2, mds_2, mem_2 \rangle, \ldots,$

$$\rho(\mathsf{s}^\bullet) = \{ \langle \text{tps}_1', mds_1', mem_1' \rangle, \langle \text{tps}_2', mds_2', mem_2' \rangle, \ldots \mid$$
$$\forall i. \exists \mathcal{B}_i. \ (\langle \text{tps}_i, mds_i, mem_i \rangle, \langle \text{tps}_i', mds_i', mem_i' \rangle) \in \mathcal{B}_i \}$$

where $\mathcal{B}_i$ denotes a strong low bisimulation modulo modes. Similarly $\mathrm{VDNI}_\mathsf{T} = ANI_\phi^\rho$ where

$$\phi(\mathsf{t}^\circ) = \left\{ \langle mds_1, mem_1' \rangle \ \middle| \ mem_1 =_{mds_1}^{Low} mem_1' \right\},$$

$$\rho(\mathsf{t}^\bullet) = \{ \langle \mathbf{tps}_1', mds_1', mem_1' \rangle, \langle \mathbf{tps}_2', mds_2', mem_2' \rangle, \ldots \mid$$
$$\forall i. \exists \mathcal{B}_i. \ (\langle \mathbf{tps}_i, mds_i, mem_i \rangle, \langle \mathbf{tps}_i', mds_i', mem_i' \rangle) \in \mathcal{B}_i \}$$

The relation $\mathcal{R}$ is a simulation, and therefore CC$^\sim$ holds. In order to apply Theorem 5.2 and conclude that whenever a source program $\mathsf{W}$ satisfies $\mathrm{VDNI}_\mathsf{S} = ANI_\phi^\rho$, then $\mathsf{W}{\downarrow}$ satisfies $\mathrm{VDNI}_\mathsf{T} = ANI_\phi^\rho$, it is sufficient for $\rho$ to satisfy Equation 1, that is

$$\rho(\mathsf{t}^\bullet) = \rho(\tilde{\tau}^\bullet(\rho(\mathsf{s}^\bullet)))$$

for $\mathsf{s}^\bullet \overset{\bullet}{\sim} \mathsf{t}^\bullet$. If one is willing to unfold all definitions, this amounts to show the set of traces "bismilar" to $\mathsf{t}^\bullet$ coincides with the set of traces that are bisimilar to some $\mathsf{t}'^\bullet$ and $\mathsf{s}'^\bullet \overset{\bullet}{\sim} \mathsf{t}'^\bullet$ for some $\mathsf{s}'^\bullet$ bisimilar to $\mathsf{s}^\bullet$. The "$\subseteq$" is immediate, while for the other one has to prove some properties of $\mathcal{R}$, the ones in the definition of secure $-$ refinement (? , inlined above Theorem 5.1]) which entails preservation of low-equivalence as shown in ? , Theorem 5.1].

1324    [CA: (not relevant for the moment) I am starting to conjecture that the equation from our theorem
1325    5.2 describes what is the cube diagram to prove if one wants to preserve some noninterference. It
1326    generalizes the cube diagram when ANI is not expressed via equiv relations but via uco]

1328    In summary, our framework makes it possible to precisely characterize the target noninterference
1329    properties that are implied by (trace-relating) correct compilation of source noninterfering programs.
1330    As we have shown, such properties are not necessarily as strong as desired. Crucially, the target
1331    noninterference property one gets *for free* for a given trace-relating correct compiler is a function of
1332    the trace relation under consideration. By considering more sophisticated trace relations, [CA: e.g.,
1333    the one in the VDNI example] one could be able to get more interesting noninterference properties
1334    in the target *for free* —but this would likely come at the expense of a more challenging trace-relating
1335    compiler correctness proof.

## 6    TRACE-RELATING SECURE COMPILATION

1338    So far we have studied compiler correctness criteria for whole, standalone programs. However,
1339    in practice, programs do not exist in isolation, but in a context where they interact with other
1340    programs, libraries, etc. In many cases, this context cannot be assumed to be benign and could
1341    instead behave maliciously to try to disrupt a compiled program.

1342        Hence, in this section we consider the following *secure compilation* scenario: a source program is
1343    compiled and linked with an arbitrary target-level context, i.e., one that may not be expressible as
1344    the compilation of a source context. Compiler correctness does not address this case, as it does not
1345    consider arbitrary target contexts, looking instead at whole programs (empty context [? ]) or well-
1346    behaved target contexts that behave like source ones (as in compositional compiler correctness [? ?
1347    ? ? ]).

1348    **Summary of the work of ? ].** To account for this scenario, ? ] describe several secure compilation
1349    criteria based on the preservation of classes of (hyper)properties (e.g., trace properties, safety,
1350    hypersafety, hyperproperties, etc.) against arbitrary target contexts. For each of these criteria, they
1351    give an equivalent "property-free" criterion, analogous to the equivalence between TP and CC$^=$. For
1352    instance, their *robust* trace property preservation criterion (RTP) states that, for any trace property
1353    $\pi$, if a source *partial* program $P$ plugged into any context $C_S$ satisfies $\pi$, then the compiled program
1354    $P{\downarrow}$ plugged into any target context $C_T$ satisfies $\pi$. Their equivalent criterion to RTP is RTC, which
1355    states that for any trace produced by the compiled program, when linked with any target context,
1356    there is a source context that produces the same trace. Formally (writing $C[P]$ to mean the whole
1357    program that results from linking partial program $P$ with context $C$) they define:

$$\text{RTP} \;\equiv\; \forall P.\,\forall \pi.\,(\forall C_S.\,\forall t.C_S\,[P]{\rightsquigarrow}t \Rightarrow t \in \pi) \Rightarrow (\forall C_T.\,\forall t.\;C_T\,[\,P{\downarrow}\,]{\rightsquigarrow}t \Rightarrow t \in \pi)$$

$$\text{RTC} \;\equiv\; \forall P.\,\forall C_T.\forall t.C_T\,[\,P{\downarrow}\,]{\rightsquigarrow}t \Rightarrow \exists C_S.\;C_S\,[P]{\rightsquigarrow}t$$

1362    In the following we adopt the notation $P \models_R \pi$ to mean "$P$ *robustly* satisfies $\pi$," i.e., $P$ satisfies $\pi$
1363    irrespective of the contexts $(C)$ it is linked with. Formally, $P \models_R \pi \overset{\text{def}}{=} \forall C, C\,[P] \models \pi$, where $\models$ is the
1364    same as before. Thus, we write more compactly:

$$\text{RTP} \;\equiv\; \forall \pi.\,\forall P.\;P \models_R \pi \Rightarrow P{\downarrow} \models_R \pi$$

1367    All the criteria of ? ] share this flavor of stating the existence of some source context that sim-
1368    ulates the behavior of any given target context, with some variations depending on the class of
1369    (hyper)properties under consideration. For trace properties, they also have criteria that preserve
1370    safety properties plus their version of liveness properties. For hyperproperties, they have criteria
1371    that preserve hypersafety properties, subset-closed hyperproperties, and arbitrary hyperproperties.

Finally, they define *relational* hyperproperties, which are relations between the behaviors of multiple programs for expressing, e.g., that a program *always* runs faster than another. For relational hyperproperties, they have criteria that preserve arbitrary relational properties, relational safety properties, relational hyperproperties and relational subset-closed hyperproperties.

Each category of criteria provides different kinds of security guarantees (confidentiality or integrity) for the code and data segments of programs. Roughly speaking, the security guarantees due to robust preservation of trace properties regard only protecting the integrity of the program from the context, the guarantees of hyperproperties also regard data confidentiality, and the guarantees of relational hyperproperties may even regard code confidentiality. Naturally, these stronger guarantees are increasingly harder to enforce and prove.

All the criteria of [? ] are stated in a setting where source and target traces are the same. In this section, we extend their results to the trace-relating setting, obtaining trintarian views for secure compilation. There are many similarities with §2 which show up in the secure compilation setting too, but also some crucial differences. As in §2, the application of $\tilde{\sigma}$ or $\tilde{\tau}$, may lose the information that a property belongs to the class *Safety*, or that a hyperproperty is subset-closed, which are both crucial for the equivalence with the property-free criterion of [? ]. As in §2, we solve this problem by interpreting classes of properties as an *abstraction* of another class of properties induced by a closure operator. Differently from §2, the presence of adversarial contexts makes the criteria for subset-closed hyperproperties and trace properties distinct. [? ] show that the criterion for robust preservation of hypersafety is distinct from robust safety preservation and all criteria about classes of trace properties are distinct from their relational counterparts e.g., robust preservation of relational safety and robust preservation of safety properties are different. We therefore further generalize the argument from §3.2 to safety hyperproperties as well as to relational hyperproperties.

Specifically, we provide a trinity for the preservation of trace properties and subset-closed hyperproperties (§6.1), of safety properties and hypersafety hyperproperties (§6.2), of hyperproperties (§6.3), and for 2-relational (hyper)properties (§6.4). We conclude the section by studying the relative expressiveness of these criteria (§6.5).

*Robustness and Compositional Compilation.* Before diving into the criteria for robust compilation, it is worth noting the relationship between these and compositional compiler correctness. Compositional compiler correctness (*CCC*) is a statement of compiler correctness for programs that are linked against *some* contexts. Unlike robustness, which imposes no constraints on the contexts, *CCC* imposes conditions on the target contexts that compiled programs can be linked against: they need to be related (in ways that vary from work to work [? ? ]) to the source contexts [? ]. As [? ] also point out, the notions of *CCC* and of robust compilation are incomparable: neither can be proven stronger than the other. This is not surprising since robust compilation criteria are used to prove compiler security while *CCC* is used to prove correctness.[15]

The criteria we adopt could be generalised further by adding an extra parameter that qualifies the relation between source and target contexts. Such a general statement would let us express both *CCC* and robust compilation by picking the correct extra parameter. However, we refrain from presenting such general statements, as the implications in terms of preservation of classes of (hyper)properties has not been studied for them.

---

[15] We remark *CCC* has been used to conclude security of compilation in the previously discussed work of [? ] (and in its predecessor [? ]). However, there is a key difference in the 'role' of contexts: in robust compilation criteria, contexts model attackers while in [? ] contexts are other bits of compiled code. This treatment lets [? ] reason compositionally about the concurrently-executing compiled code.

### 6.1 Trace-Relating Secure Compilation: Trace Properties and Subset-closed Hyperproperties

This section shows the simple generalization of RTC to the trace-relating setting (RTC$^\sim$) and its corresponding trinitarian view (Theorem 6.1). Then, it presents the trinitarian view for criteria that preserve subset-closed hyperproperties (Theorem 6.2).

**Theorem 6.1** (Trinity for Robust Trace Properties �explanatory). For any trace relation $\sim$ and induced property mappings $\tilde{\tau}$ and $\tilde{\sigma}$, we have: $\text{RTP}^{\tilde{\tau}} \iff \text{RTC}^\sim \iff \text{RTP}^{\tilde{\sigma}}$, where

$$\text{RTC}^\sim \equiv \forall \mathsf{P} \; \forall \mathsf{C_T} \; \forall \mathsf{t}. \; \mathsf{C_T}\,[\,\mathsf{P}{\downarrow}\,] \rightsquigarrow \mathsf{t} \Rightarrow \exists \mathsf{C_S} \; \exists \mathsf{s} \sim \mathsf{t}. \; \mathsf{C_S}\,[\,\mathsf{P}\,] \rightsquigarrow \mathsf{s}$$

$$\text{RTP}^{\tilde{\tau}} \equiv \forall \mathsf{P} \; \forall \pi_\mathsf{S} \in 2^{\text{Trace}_\mathsf{S}}. \; \mathsf{P} \models_\mathsf{R} \pi_\mathsf{S} \Rightarrow \mathsf{P}{\downarrow} \models_\mathsf{R} \tilde{\tau}(\pi_\mathsf{S})$$

$$\text{RTP}^{\tilde{\sigma}} \equiv \forall \mathsf{P} \; \forall \pi_\mathsf{T} \in 2^{\text{Trace}_\mathsf{T}}. \; \mathsf{P} \models_\mathsf{R} \tilde{\sigma}(\pi_\mathsf{T}) \Rightarrow \mathsf{P}{\downarrow} \models_\mathsf{R} \pi_\mathsf{T}$$

The trinity for robust trace property preservation is the straightforward adaptation of the concepts of §2 to the definitions of [? ]. Intuitively, these criteria simply deal with partial programs $\mathsf{P}$ instead of whole programs $\mathsf{W}$. Necessarily, these criteria then consider arbitrary program contexts linked with $\mathsf{P}$; the universal quantification over $\mathsf{C_S}$ and $\mathsf{C_T}$ are tacit in the expression $\models_R$.

We can also generalize §2 to *robust* subset-closed hyperproperties (Theorem 6.2). However, unlike the correct compilation case of §2, the equivalent property-free criterion (RSCHC$^\sim$) does not coincide with RSC$^\sim$, but states the existence of a single source context for all the target traces produced by a program in a given context.

**Theorem 6.2** (Trinity for Robust Subset-closed Hyperproperties ✑). Let $\mathsf{SCH_S}$ and $\mathsf{SCH_T}$ denote the sets of all subset-closed hyperproperties in the source and target languages, respectively. For any trace relation $\sim$ and its existential and universal images lifted to hyperproperties (that is, the lifting of the respective functions from Definition 2.5), $\tilde{\tau}$ and $\tilde{\sigma}$, and for $Cl_\subseteq(H) = \{\pi \mid \exists \pi' \in H. \; \pi \subseteq \pi'\}$, we have: $\text{RSCHP}^{Cl_\subseteq \circ \tilde{\tau}} \iff \text{RSCHC}^\sim \iff \text{RSCHP}^{Cl_\subseteq \circ \tilde{\sigma}}$, where

$$\text{RSCHC}^\sim \equiv \forall \mathsf{P} \; \forall \mathsf{C_T} \; \exists \mathsf{C_S} \; \forall \mathsf{t} \; \mathsf{C_T}\,[\,\mathsf{P}{\downarrow}\,] \rightsquigarrow \mathsf{t} \Rightarrow \exists \mathsf{s} \sim \mathsf{t}'. \; \mathsf{C_S}\,[\,\mathsf{P}\,] \rightsquigarrow \mathsf{s}$$

$$\text{RSCHP}^{Cl_\subseteq \circ \tilde{\tau}} \equiv \forall \mathsf{P} \; \forall \mathsf{H_S} \in \mathsf{SCH_S}. \; \mathsf{P} \models_\mathsf{R} \mathsf{H_S} \Rightarrow \mathsf{P}{\downarrow} \models_\mathsf{R} Cl_\subseteq(\tilde{\tau}(\mathsf{H_S}))$$

$$\text{RSCHP}^{Cl_\subseteq \circ \tilde{\sigma}} \equiv \forall \mathsf{P} \; \forall \mathsf{H_T} \in \mathsf{SCH_T}. \; \mathsf{P} \models_\mathsf{R} Cl_\subseteq(\tilde{\sigma}(\mathsf{H_T})) \Rightarrow \mathsf{P}{\downarrow} \models_\mathsf{R} \mathsf{H_T}$$

### 6.2 Trace-Relating Secure Compilation: Safety and Hypersafety

In this section we elaborate the robust preservation of safety (Theorem 6.3) and hypersafety properties (Theorem 6.4). Similar to §3.2, we consider the trace model adopted by [? ] to ease the presentation. Our starting point is the two equivalent criteria for preservation of robust satisfaction of *all* and *only* the safety properties [? ],

$$\text{RSP} \equiv \forall \mathsf{P}. \; \forall \pi \in Safety. \; \mathsf{P}{\models_\mathsf{R}}\pi \Rightarrow \mathsf{P}{\downarrow}{\models_\mathsf{R}}\pi$$

$$\text{RSC} \equiv \forall \mathsf{P}. \; \forall \mathsf{C_T}. \forall m. \mathsf{C_T}\,[\,\mathsf{P}{\downarrow}\,] \rightsquigarrow^* m \Rightarrow \exists \mathsf{C_S}. \; \mathsf{C_S}\,[\,\mathsf{P}\,] \rightsquigarrow^* m$$

where $\mathsf{C_T}\,[\,\mathsf{P}{\downarrow}\,] \rightsquigarrow^* m$ is a shorthand for $\exists t \geq m. \mathsf{C_T}\,[\,\mathsf{P}{\downarrow}\,] \rightsquigarrow t$.

RSP differs from RTP as it only quantifies over safety properties, and RSC differs from RTC as it quantifies over finite prefixes $m$, rather than complete traces $t$. This comes from the fact that safety properties can be characterized in terms of sets of *bad* prefixes (as in Definition 3.4). Unfolding $\rightsquigarrow^*$ we can interpret RSC as follows. If $\mathsf{C_T}\,[\,\mathsf{P}{\downarrow}\,]$ produces a trace $t \geq m$ that violates a specific safety property, namely, the one defined by $M = \{m\}$, then there exists $\mathsf{C_S}$ in which $\mathsf{P}$ violates the *same* safety property, producing a trace $t' \geq m$ but possibly distinct from $t$.

Our generalization of RSC to the trace-relating setting states that whenever $C_T[P\downarrow]$ produces a trace $t$ that violates a target safety property, there exists a source context $C_S$ in which $P$ violates the source *interpretation* of the property, i.e., its image through $\tilde{\sigma}$. The following theorem defines $RSC^\sim$ and its two equivalent formulations.

**Theorem 6.3** (Trinity for Robust Safety Properties 🦜). *For any trace relation $\sim$ and for the corresponding property mappings $\tilde{\tau}$ and $\tilde{\sigma}$, we have:* $RTP^{Safe\circ\tilde{\tau}} \iff RSC^\sim \iff RSP^{\tilde{\sigma}}$, *where*

$$RSC^\sim \equiv \forall P \, \forall C_T \, \forall t \, \forall m \le t. C_T[P\downarrow] \rightsquigarrow t \Rightarrow \exists C_S \, \exists t' \ge m \, \exists s \sim t'. \, C_S[P] \rightsquigarrow s$$

$$RTP^{Safe\circ\tilde{\tau}} \equiv \forall P \forall \pi_S \in 2^{Traces}. P \models_R \pi_S \Rightarrow P\downarrow \models_R (Safe\circ\tilde{\tau})(\pi_S)$$

$$RSP^{\tilde{\sigma}} \equiv \forall P \forall \pi_T \in Safety_T. P \models_R \tilde{\sigma}(\pi_T) \Rightarrow P\downarrow \models_R \pi_T$$

where the closure operator *Safe* is the one introduced in §3.2.

Exactly like §3.2, Theorem 6.3 exploits the fact that

$$Safe\circ\tilde{\tau} : 2^{Traces} \leftrightarrows Safety_T : \tilde{\sigma}$$

is a Galois connection between source properties and target safety properties and the argument generalizes to arbitrary closure operators on target properties (🦜). More interestingly, we can further generalize this idea to hypersafety. Hypersafety lifts the idea of safety with another level of sets (just like hyperproperties do w.r.t. trace properties) in order to talk about multiple runs of the same program. Just like for safety, hypersafety is concerned with a set of bad prefixes (called $M$) that no program upholding the hypersafety property should extend. Formally, a hyperproperty $H$ is hypersafety if: $\forall \pi. \pi \notin H \Rightarrow (\exists M. M \prec \pi \wedge (\forall \pi'M \prec \pi' \Rightarrow \pi' \notin H))$. In Theorem 6.4, we indeed exploit the following Galois connection between source subset-closed hyperproperties and target

$$HSafe\circ\tilde{\tau} : SCH_S \leftrightarrows HSafety_T : Cl_\subseteq \circ \tilde{\sigma}$$

where $HSafety_T = \left\{ HSafe(H_T) \mid H_T \in 2^{2^{Trace_T}} \right\}$ and *HSafe* is the closure operator that maps an arbitrary target hyperproperty $H_T$ to the target hypersafety that best over-approximates $H_T$.[16]

**Theorem 6.4** (Trinity for Robust Hypersafety 🦜). *For any trace relation $\sim$ and for the induced property mappings $\tilde{\tau}$ and $\tilde{\sigma}$, we have:* $RSCHP^{HSafe\circ\tilde{\tau}} \iff RHSC^\sim \iff RHSP^{Cl_\subseteq\circ\tilde{\sigma}}$, *where*

$$RHSC^\sim \equiv \forall P \, \forall C_T \, \forall M \in \mathcal{M}^{fin}. M \le beh(C_T[P\downarrow]) \Rightarrow$$
$$\exists C_S \, \forall m \in M \, \exists t \ge m. \, \exists s \sim t. \, C_S[P] \rightsquigarrow s$$

$$RSCHP^{HSafe\circ\tilde{\tau}} \equiv \forall P \, \forall H_S \in SCH_S. \, P \models_R H_S \Rightarrow P\downarrow \models_R HSafe(\tilde{\tau}(H_S))$$

$$RHSP^{Cl_\subseteq\circ\tilde{\sigma}} \equiv \forall P \, \forall H_T \in HSafety_T. P \models_R Cl_\subseteq(\tilde{\sigma}(H_T)) \Rightarrow P\downarrow \models_R H_T$$

and $\mathcal{M}^{fin}$ is the set of *finite* sets of prefixes.

We conclude this section with the following remark. The reader might wonder about extracting a "new" trace relation from the Galois connection $Safe\circ\tilde{\tau} : 2^{Traces} \leftrightarrows Safety_T : \tilde{\sigma}$ and get another formulation of $RSC^\sim$. We note that this is not possible in general, as the class of safety properties, i.e., closed sets, is not necessarily a powerset and hence Lemma 2.7 cannot be applied.

---

[16] $HSafe(H_T) = \cap \left\{ H'_T \mid H_T \subseteq H'_T \wedge H'_T \in HSafety_T \right\}$. See, e.g., ? ] and ? ].

## 6.3 Trace-Relating Secure Compilation: Arbitrary Hyperproperties

We already mentioned that some properties of interest for security e.g., possibilistic information-flow are not subset closed [? ]. In this section we lift the results from §3.3 to the *secure* compilation setting. Once again, the trinity is *weak* as the equivalence to $\text{RHP}^{\tilde{\sigma}}$ requires an extra assumption.

**Theorem 6.5** (Weak Trinity for Robust Hyperproperties ✎). *For a trace relation $\sim \subseteq \text{Traces}_S \times \text{Trace}_T$ and induced property mappings $\tilde{\sigma}$ and $\tilde{\tau}$, we have:*

$\text{RHC}^\sim \iff \text{RHP}^{\tilde{\tau}}$;

*if $\tilde{\tau} \leftrightarrows \tilde{\sigma}$ is a Galois insertion (i.e., $\tilde{\tau} \circ \tilde{\sigma} = id$), then $\text{RHC}^\sim \Rightarrow \text{RHP}^{\tilde{\sigma}}$,*

*if $\tilde{\sigma} \leftrightarrows \tilde{\tau}$ is a Galois reflection (i.e., $\tilde{\sigma} \circ \tilde{\tau} = id$), then $\text{RHP}^{\tilde{\sigma}} \Rightarrow \text{RHP}^{\tilde{\tau}}$,*

$$\text{where } \text{RHC}^\sim \equiv \forall P\, \forall C_T\, \exists C_S\, \forall t.\, C_T\,[\,P\!\downarrow\,]\rightsquigarrow t \iff (\exists s \sim t.\, C_S\,[P]\rightsquigarrow s)$$

$$\text{RHP}^{\tilde{\tau}} \equiv \forall P\, \forall H_S.\, P \models_R H_S \Rightarrow P\!\downarrow\, \models_R \tilde{\tau}(H_S)$$

$$\text{RHP}^{\tilde{\sigma}} \equiv \forall P\, \forall H_T.\, P \models_R \tilde{\sigma}(H_T) \Rightarrow P\!\downarrow\, \models_R H_T$$

It is therefore possible and correct to deduce a source obligation for a given target hyperproperty $H_T$ ($\text{RHC}^\sim \Rightarrow \text{RHP}^{\tilde{\sigma}}$) when no information is lost in the composition $\tilde{\tau} \circ \tilde{\sigma}$. On the other hand, $\text{RHP}^{\tilde{\tau}}$ is a consequence of $\text{RHP}^{\tilde{\sigma}}$ when no information is lost in composing in the other direction, $\tilde{\sigma} \circ \tilde{\tau}$.

## 6.4 Trace-Relating Secure Compilation: 2-Relational Hyperproperties

Finally, we turn to *relational* properties and hyperproperties. Relational hyperproperties, as defined by [? ], are predicates on a sequence of behaviors; a sequence of programs has the relational hyperproperty if their behaviors collectively satisfy the predicate. Depending on the arity of the sequence, there exist different subclasses of relational hyperproperties, though for simplicity here we only study relational hyperproperties of arity 2. A key example of a relational hyperproperty is trace equivalence, which holds if two programs have identical behaviors.

All the trinities in this section follow the pattern of their non-relational counterparts. We first explain how one can get a Galois connection between source and target relational properties from a trace relation.

Given a trace relation $\sim \subseteq \text{Traces}_S \times \text{Trace}_T$, we can relate pairs of source traces with pairs of target traces point-wise,

$$(s_1, s_2) \sim (t_1, t_2) \iff s_1 \sim t_1 \wedge s_2 \sim t_2$$

Formally this is $\sim^2 \subseteq \text{Traces}_S{}^2 \times \text{Trace}_T{}^2$, the product of the relation $\sim$ with itself. Therefore by Lemma 2.7 it corresponds to a Galois connection between source and target relational properties (✎), that with a little abuse of notation[17] we still denote by

$$\tilde{\tau} : 2^{\text{Traces}_S \times \text{Traces}_S} \leftrightarrows 2^{\text{Trace}_T \times \text{Trace}_T} : \tilde{\sigma}$$

Explicitly, for $r_S \in 2^{\text{Traces}_S \times \text{Traces}_S}$ and $r_T \in 2^{\text{Trace}_T \times \text{Trace}_T}$,

$$\tilde{\tau}(r_S) = \{(t_1, t_2) \mid \exists (s_1, s_2).\, s_1 \sim t_1 \wedge s_2 \sim t_2 \wedge (s_1, s_2) \in r_S\}$$

$$\tilde{\sigma}(r_T) = \{(s_1, s_2) \mid \forall (t_1, t_2).\, s_1 \sim t_1 \wedge s_2 \sim t_2 \Rightarrow (t_1, t_2) \in r_T\}$$

$\tilde{\tau}$ and $\tilde{\sigma}$ are then lifted to relational hyperproperties similarly to Definition 3.2. Explicitly, for $R_S \in 2^{2^{\text{Traces}_S \times \text{Traces}_S}}$ and $R_T \in 2^{2^{\text{Trace}_T \times \text{Trace}_T}}$,

$$\tilde{\tau}(R_S) = \{\tilde{\tau}(r_S) \mid r_S \in R_S\}$$

$$\tilde{\sigma}(R_T) = \{\tilde{\sigma}(r_T) \mid r_T \in R_T\}$$

---

[17]Technically, we should write: $\tilde{\tau}^2 \leftrightarrows \tilde{\sigma}^2$

Given a relational property $r \in 2^{Trace \times Trace}$ and two programs $P_1, P_2$, we write $P_1, P_2 \models_R r$ for

$$\forall C. \forall t_1 t_2. \ C\,[P_1] \rightsquigarrow t_1 \ \wedge \ C\,[P_2] \rightsquigarrow t_2 \Rightarrow (t_1, t_2) \in r$$

Given a relational hyperproperty $R \in 2^{2^{Trace \times Trace}}$, by $P_1, P_2 \models_R R$ we mean

$$\forall C.(beh(C\,[P_1]), beh(C\,[P_2])) \in R$$

**Theorem 6.6** (Trinity for Robust 2-Relational Trace Properties ✎). For any trace relation $\sim$ and for the corresponding property mappings $\tilde{\tau}$ and $\tilde{\sigma}$, we have: R2rTP$^{\tilde{\tau}}$ $\iff$ R2rTC$^{\sim}$ $\iff$ R2rTP$^{\tilde{\sigma}}$, where

$$\text{R2rTC}^{\sim} \ \equiv \ \forall C_T \ \forall P_1 \ \forall P_2 \ \forall t_1 \ \forall t_2. \ (C_T\,[P_1\downarrow]\rightsquigarrow t_1 \wedge C_T\,[P_2\downarrow]\rightsquigarrow t_2) \Rightarrow$$
$$\exists C_S \ \exists s_1 \sim t_1 \ \exists s_2 \sim t_2. \ C_S\,[P_1]\rightsquigarrow s_1 \wedge C_S\,[P_2]\rightsquigarrow s_2$$
$$\text{R2rTP}^{\tilde{\tau}} \ \equiv \ \forall P_1 P_2 \ \forall r_S \in 2^{Trace_S \times Trace_S}. \ P_1, P_2 \models_R r_S \Rightarrow \ P_1\downarrow, P_2\downarrow \models_R \tilde{\tau}(r_S)$$
$$\text{R2rTP}^{\tilde{\sigma}} \ \equiv \ \forall P_1 P_2. \ \forall r_T \in 2^{Trace_T \times Trace_T}. \ P_1, P_2 \models_R \tilde{\sigma}(r_T) \Rightarrow \ P_1\downarrow, P_2\downarrow \models_R r_T$$

Next, we propose the trinity for 2-relational subset-closed hyperproperties, i.e., elements of $2^{2^{Trace \times Trace}}$ that are closed under subsets. Exactly as in the case of subset-closed hyperproperties, the application of $\tilde{\tau}$ and $\tilde{\sigma}$ may lose the information of being subset-closed. In order to guarantee the equivalence of the three criteria, we compose the two mappings with a closure operator that we still denote by $Cl_{\subseteq}$.

**Theorem 6.7** (Trinity for 2-Relational Robust Subset-Closed Hyperproperties ✎). For any trace relation $\sim$ and for the corresponding property mappings $\tilde{\tau}$ and $\tilde{\sigma}$, we have R2rSCHP$^{Cl_{\subseteq} \circ \tilde{\tau}}$ $\iff$ R2rSCHC$^{\sim}$ $\iff$ R2rSCHP$^{Cl_{\subseteq} \circ \tilde{\sigma}}$, where

$$\text{R2rSCHC}^{\sim} \ \equiv \ \forall C_T \ \forall P_1 \ \forall P_2 \ \exists C_S \ \forall t_1 \ \forall t_2. \ (C_T\,[P_1\downarrow]\rightsquigarrow t_1 \wedge C_T\,[P_2\downarrow]\rightsquigarrow t_2) \Rightarrow$$
$$\exists s_1 \sim t_1 \ \exists s_2 \sim t_2. \ C_S\,[P_1]\rightsquigarrow s_1 \wedge C_S\,[P_2]\rightsquigarrow s_2$$
$$\text{R2rSCHP}^{Cl_{\subseteq} \circ \tilde{\tau}} \ \equiv \ \forall P_1 \ \forall P_2 \ \forall R_S \in \text{2RelSCH}_S. \ P_1, P_2 \models_R R_S \Rightarrow P_1\downarrow, P_2\downarrow \models_R \tilde{\tau}(R_S)$$
$$\text{R2rSCHP}^{Cl_{\subseteq} \circ \tilde{\sigma}} \ \equiv \ \forall P_1 \ \forall P_2 \ \forall R_T \in \text{2RelSCH}_T. P_1, P_2 \models_R \tilde{\sigma}(R_T) \Rightarrow P_1\downarrow, P_2\downarrow \models_R R_T$$

We move now to the class of relational safety properties, a notion that generalizes safety properties to relations on programs. Similarly to Theorem 6.3, R2rSP$^{\tilde{\sigma}}$ quantifies over target relational safety properties, while R2rTP$^{2rSafe \circ \tilde{\tau}}$ quantifies over all source relational property and compose $\tilde{\tau}$ with $2rSafe$ a closure operator that best approximates a relational property with a relational safety property.

**Theorem 6.8** (Trinity for Robust 2-Relational Safety Properties ✎). For any trace relation $\sim$ and for the corresponding property mappings $\tilde{\tau}$ and $\tilde{\sigma}$, we have: R2rTP$^{2rSafe \circ \tilde{\tau}}$ $\iff$ R2rSC$^{\sim}$ $\iff$ R2rSP$^{\tilde{\sigma}}$, where

$$\text{R2rSC}^{\sim} \ \equiv \ \forall C_T \ \forall P_1 P_2 \ \forall t_1 t_2 \ \forall m_1 \le t_1 \ \forall m_2 \le t_2. \ C_T\,[P_1\downarrow]\rightsquigarrow t_1 \Rightarrow C_T\,[P_2\downarrow]\rightsquigarrow t_2 \Rightarrow$$
$$\exists C_S \ \exists t_1' \ge m_1 \ \exists s_1 \sim t_1' \ \exists t_2' \ge m_2 \ \exists s_2 \sim t_2'. \ C_S\,[P_1]\rightsquigarrow s_1 \wedge C_S\,[P_2]\rightsquigarrow s_2$$
$$\text{R2rTP}^{2rSafe \circ \tilde{\tau}} \ \equiv \ \forall P_1 P_2 \ \forall r_S \in 2^{Trace_S \times Trace_S}. \ P_1, P_2 \models_R r_S \Rightarrow P_1\downarrow, P_2\downarrow \models_R (2rSafe \circ \tilde{\tau})(r_S)))$$
$$\text{R2rSP}^{\tilde{\sigma}} \ \equiv \ \forall P_1 P_2 \ \forall r_T \in \text{2rel-Safety}_T. \ P_1, P_2 \models_R \tilde{\sigma}(r_T) \Rightarrow P_1\downarrow, P_2\downarrow \models_R r_T$$

Finally, we present the most general criterion: preservation of *arbitrary* 2-relational hyperproperties. As for the preservation of arbitrary hyperproperties, this (weak) trinity requires additional assumptions to hold, namely that the Galois connection is an insertion or a reflection.

**Theorem 6.9** (Weak trinity for Robust 2-Relational Hyperproperties 🕊). For a trace relation $\sim \subseteq \text{Traces}_S \times \textbf{Trace}_T$ and the corresponding property mappings $\tilde{\sigma}$ and $\tilde{\tau}$, we have:

R2rHC$^\sim$ $\iff$ R2rHP$^{\tilde{\tau}}$;

if $\tilde{\tau} \leftrightarrows \tilde{\sigma}$ is a Galois insertion (i.e., $\tilde{\tau} \circ \tilde{\sigma} = id$), then R2rHC$^\sim \Rightarrow$ R2rHP$^{\tilde{\sigma}}$,

if $\tilde{\sigma} \leftrightarrows \tilde{\tau}$ is a Galois reflection (i.e., $\tilde{\sigma} \circ \tilde{\tau} = id$), then R2rHP$^{\tilde{\sigma}} \Rightarrow$ R2rHP$^{\tilde{\tau}}$,

$$\text{where R2rHC}^\sim \equiv \forall P_1 P_2 \, \forall C_T \, \exists C_S.$$
$$(\forall t. \, C_T \, [\, P_1\downarrow] \leadsto t \iff (\exists s \sim t. \, C_S \, [P_1] \leadsto s)) \wedge$$
$$(\forall t. \, C_T \, [\, P_2\downarrow] \leadsto t \iff (\exists s \sim t. \, C_S \, [P_2] \leadsto s))$$
$$\text{R2rHP}^{\tilde{\tau}} \equiv \forall P_1 \, \forall P_2 \, \forall R_S. \, P_1, P_2 \models_R R_S \Rightarrow P_1\downarrow, P_2\downarrow \models_R \tilde{\tau}(R_S)$$
$$\text{R2rHP}^{\tilde{\sigma}} \equiv \forall P_1 \, \forall P_2 \, \forall R_T. P_1, P_2 \models_R \tilde{\sigma}(R_T) \Rightarrow P_1\downarrow, P_2\downarrow \models_R R_T$$

## 6.5 Relating the Secure Compilation Trinities

Figure 4 orders criteria referring to the same trace relation $\sim$ according to their relative strength. If a trinity entails another (denoted by $\Rightarrow$), then the former provides stronger security for a compilation chain than the latter.

The hypotheses of insertion and reflection mentioned in Theorem 6.9 and Theorem 6.5 are highlighted with the labels 'Ins' and 'Refl'. Recall that when composing $\tilde{\tau}$ with *Safe* we quantify over the whole class of source trace properties rather than only safety properties. This is represented by the blue background in RTP$^{Safe \circ \tilde{\tau}}$. The trinity for the robust preservation of arbitrary trace properties is on the same blue background. Red and green backgrounds are reserved for subset-closed hyperproperties and arbitrary relational properties and serve the same purpose.

We now describe how to interpret the acronyms in Figure 4. All criteria start with R meaning they refer to robust preservation (secure compilation criteria). Criteria for relational hyperproperties—here only arity 2 is shown for simplicity—contain 2r. Next, criteria names spell the class of hyperproperties they preserve: H for hyperproperties, SCH for subset-closed hyperproperties, HS for hypersafety, T for trace properties, and S for safety properties. Finally, property-free criteria end with a C while property-full ones involving $\tilde{\sigma}$ and $\tilde{\tau}$ end with P. Thus, *robust (R) subset-closed hyperproperty-preserving (SCH) compilation (C)* is RSCHC$^\sim$, *robust (R) two-relational (2r) safety-preserving (S) compilation (C)* is R2rSC$^\sim$, etc.

## 7 INSTANCES OF TRACE-RELATING SECURE COMPILATION

This section presents instances of compilers that adopt our framework for secure compilation purposes. We provide three illustrative cases, for compilers that respectively robustly-preserve trace properties (§7.1), safety properties (§7.2) and hypersafety properties (§7.3). The last two examples are not novel instances we devise but rather existing work whose results we recount as instantiations of our framework.

### 7.1 An Instance of Trace-Relating Robust Preservation of Trace Properties

This subsection illustrates trace-relating secure compilation when the target events are strictly more events than the source ones

The source and target languages used here extend the syntax of the source language of §4.3.1. Both languages have *outputs of naturals*, and the expressions that generate them: out$_S$ n and **out$_S$ e**. Additionally, the target has a different output action and its related expression **out$_T$ n**; this is the only difference between the languages. The extra events in the target model the ability of target language to perform potentially-dangerous operations (e.g., writing to the hard drive), which cannot

R2rHC$^{\sim}$

R2rHP$^{\tilde{\sigma}}$  $\xrightarrow{\text{Refl.}}$  R2rHP$^{\tilde{\tau}}$

Ins.

R2rSCHC$^{\sim}$

R2rSCHP$^{Cl_\subseteq \circ \tilde{\sigma}}$  $\longleftrightarrow$  R2rSCHP$^{Cl_\subseteq \circ \tilde{\tau}}$

RHC$^{\sim}$

Ins.

RHP$^{\tilde{\sigma}}$  $\xrightarrow{\text{Refl.}}$  RHP$^{\tilde{\tau}}$

R2rTC$^{\sim}$

R2rTP$^{\tilde{\sigma}}$  $\longleftrightarrow$  R2rTP$^{\tilde{\tau}}$

RSCHC$^{\sim}$

RSCHP$^{Cl_\subseteq \circ \tilde{\sigma}}$  $\Leftrightarrow$  RSCHP$^{Cl_\subseteq \circ \tilde{\tau}}$

R2rSC$^{\sim}$

R2rSP$^{\tilde{\sigma}}$  $\longleftrightarrow$  R2rTP$^{2rSafe\circ\tilde{\tau}}$

RTC$^{\sim}$

RTP$^{\tilde{\sigma}}$  $\longleftrightarrow$  RTP$^{\tilde{\tau}}$

RHSC$^{\sim}$

RHSP$^{Cl_\subseteq \circ \tilde{\sigma}}$  $\longleftrightarrow$  RSCHP$^{HSafe\circ\tilde{\tau}}$

RSC$^{\sim}$

RSP$^{\tilde{\sigma}}$  $\longleftrightarrow$  RTP$^{Safe\circ\tilde{\tau}}$

| R robust | 2r 2-relational | |
| --- | --- | --- |
| H hyperproperties | SCH subset-closed hyperproperties | HS hypersafety |
| T trace properties | S safety properties | |
| C property-full criterion | P property-free criterion based on $\sigma$ and $\tau$ | |

Fig. 4. Hierarchy of trinitarian views of secure compilation criteria preserving classes of hyperproperties and the key to read each acronym. Shorthands 'Ins.' and 'Refl.' stand for Galois Insertion and Reflection. The 🐓 symbol denotes trinities proven in Coq.

be performed by the source language, and against which source-level reasoning can therefore offer no protection.

Both languages and compilation chains now deal with partial programs $P$, contexts $C$ and linking of those two to produce whole programs $C[P]$. In this setting, a whole program $W$ is the combination of a *main expression* to be evaluated and a set of *function definitions fs* (with distinct names) that can refer to their argument (*arg*) symbolically and can be called by the main expression and by other functions ($f(e)$. The set of functions of a whole program is the union of the functions of a partial program and a context; the latter also contains the main expression.

$$e ::= \cdots \mid f(e) \mid \text{out}_S\ n \mid \text{arg} \qquad e ::= \cdots \mid f(e) \mid \text{out}_S\ n \mid \text{arg} \mid \text{out}_T\ n$$
$$i ::= \cdots \mid \text{out}_S\ n \qquad\qquad\qquad i ::= \cdots \mid \text{out}_S\ n \mid \text{out}_T\ n$$
$$fs ::= \langle f_1, e_1 \rangle, \ldots, \langle f_n, e_n \rangle \qquad P ::= \langle fs, e \rangle \qquad C ::= fs \qquad W ::= C[P]$$

The extensions of the typing rules and the operational semantics for whole programs are unsurprising and therefore elided. The trace model also follows closely that of §4.3: it consists of a list of *regular events* (including the new outputs) terminated by a *result event*[18]. A partial program and a context can be linked into a whole program when their functions satisfy the requirements mentioned above.

---

[18]Notice that the languages are strictly terminating.

We define the homomorphic compiler ($\cdot\!\downarrow$) that translates each source construct into its target correspondent. Thus, the compiler never introduces the additional target instruction $\mathbf{out_T\ n}$. Since it is straightforward, the formalisation of the compiler is elided.

**Relating Traces.** In the present model, source and target traces differ only in the fact that the target draws (regular) events from a strictly larger set than the source, i.e., $\Sigma_T \supset \Sigma_S$. A natural relation between source and target traces essentially maps a given target trace $\mathbf{t}$ the source trace that erases from $\mathbf{t}$ those events that exist only at the target level. This is reasonable because only target contexts $C$ (not compiled programs $P\!\downarrow$) can perform the extra target actions as the compiler does not introduce them. Let $\mathbf{t}|_{\Sigma_S}$ indicate trace $\mathbf{t}$ filtered to retain only those elements included in alphabet $\Sigma_S$. We define the trace relation as:

$$\mathbf{s} \sim \mathbf{t} \quad \equiv \quad \mathbf{s} = \mathbf{t}|_{\Sigma_S}$$

In the opposite direction, a source trace $\mathbf{s}$ is related to many target ones, as any target-only events can be inserted at any point in $\mathbf{s}$. The induced mappings for this relation are:

$$\tilde{\tau}(\pi_S) = \{\mathbf{t} \mid \exists \mathbf{s}.\, \mathbf{s} = \mathbf{t}|_{\Sigma_S} \wedge \mathbf{s} \in \pi_S\}$$

$$\tilde{\sigma}(\boldsymbol{\pi}_T) = \{\mathbf{s} \mid \forall \mathbf{t}.\, \mathbf{s} = \mathbf{t}|_{\Sigma_S} \Rightarrow \mathbf{t} \in \boldsymbol{\pi}_T\}$$

That is, the target guarantee of a source property is that the target has the same source-level behavior, sprinkled with arbitrary target-level behavior. Conversely, the source-level obligation of a target property is the aggregate of those source traces all of whose target-level enrichments are in the target property.

Since the languages are very similar, it is simple to prove that our compiler is secure according to the trace relation $\sim$ defined above.

**Theorem 7.1** ($\cdot\!\downarrow$ is Secure 🐾). $\cdot\!\downarrow$ is RTC$^\sim$.

## 7.2 An Instance of Trace-Relating Robust Preservation of Safety Properties

I/O events are not the only instance of events that compilers consider. Especially in the setting of secure compilation, where a compartmentalized partial program interacts with a context, *interaction traces* are often used [? ? ? ?]. Consider a language analogous to that of the previous section, where the context $C$ defines a set of functions $F_c$ and the program defines a different set $F_p$. Interaction traces (generally) record the control flow of calls between these two sets via actions that are *call $f$ $v$* and *ret $v$* [?]. These actions indicate a call to function $f$ with parameter $v$ and a return with return value $v$. In case the context calls a function in $F_p$ (or returns to a function in $F_p$), the action is decorated with a ? (i.e., those actions are *call $f$ $v$?* and *ret $v$?*). Dually, the program calling a function in $F_c$ (or returning to it) generates an action decorated with a ! (i.e., those actions are *call $f$ $v$!* and *ret $v$!*).

? ] consider precisely such a setting. Their languages are simple like those presented here but impure; their source has an ML-like heap and the target has a memory that is indexed by natural numbers and capabilities to protect addresses. Moreover, they define a compiler that preserves safety properties of source programs (i.e., it is RSC$^\sim$ in the sense of Theorem 6.3) by relying on the target capabilities. The interesting point, however, is that they also consider source and target traces to be distinct since the two languages have different values. Concretely, the source has bools and nats and the target only has **nats**, plus in the source, heap addresses are abstract locations $\ell$ while in the target they are **nats**. Thus, to prove RSC$^\sim$, they rely on a cross-language relation on values, which is lifted to trace actions, and then lifted point-wise to traces (analogously to what we have done in Sections 4.3, 4.4 and 7.1). In order to relate addresses, their cross-language

relation is equipped with a partial bijection between source and target addresses, this bijection grows monotonically with every reduction step.

Besides defining a relation on traces (which is an instance of $\sim$), they also define a relation between source and target safety properties that supports concurrent programs.[19] Thus, they really provide an instantiation of $\tau$ that maps all safe source traces to the related target ones. This ensures that no additional target trace is introduced in the target property, and source safety property are mapped to target safety ones by $\tau$. Thus, their compiler is proven to generate code that respects $\tau$, so they really achieve a variation of $\text{RTP}^{Safe \circ \tilde{\tau}}$ from Theorem 6.3. Their proofs are based on standard techniques either for secure compilation (i.e., trace-based backtranslation [? ]) and for correct compilation (i.e., forward/backward simulation [? ]).

## 7.3 An Instance of Trace-Relating Robust Preservation of Hypersafety Properties

? ] study the preservation of hypersafety from the perspective of secure compilation. Again, their result can be interpreted in our setting. They consider reactive systems, where trace alphabets are partitioned in input actions $\alpha$? and output actions $\alpha$!, whose concatenation generate traces $\overline{\alpha?\alpha!}$. We use the same notation as before and indicate such sequences as s and t respectively. The set of target output actions $\alpha$! includes an action $\sqrt{}$ that has no source counterpart (i.e., $\nexists \alpha? \sim \sqrt{}$), and whose output does not depend on internal state (thus it cannot leak secrets).[20] By emitting $\sqrt{}$ whenever undesired inputs are fed to a compiled program (e.g., passing a **nat** when a **bool** is expected), hypersafety is preserved (as $\sqrt{}$ does not leak secrets) [? ].

More formally, they assume a relation on actions $\sim$ that is total on the source actions and injective. From there, they define **TPC**—which here corresponds to an instance of $\tau$—that maps the set of valid source traces to the set of valid target traces (that now mention $\sqrt{}$) as follows:

$$\text{TPC}\,(\pi_S) = \left\{ t \;\middle|\; t \in \bigcup_{n \in \mathbb{N}} \text{int}_n\,(\pi_S) \right\} \quad \text{where } \text{int}_0\,(\pi_S) = \{t \mid \exists s \in \pi_S \wedge s \sim t\}$$

$$\text{int}_{n+1}\,(\pi_S) = \{t \mid t \equiv t_1 \alpha? \sqrt{} t_2 \wedge t_1 t_2 \in \text{int}_n\,(\pi_S) \wedge \text{undesired}\,(\alpha?)\}$$

where undesired $(\alpha?)$ indicates that $\alpha$? is an undesired input (intuitively, this is an information that can be derived from the set of source traces [? ]).

Informally, given a set of source traces $\pi_S$, **TPC** generates all target traces that are related (pointwise) to a source trace (case $\text{int}_0$). Then (case $\text{int}_{n+1}$), it adds all traces (t) with interleavings of undesired input $\alpha$? (third conjunct) followed by $\sqrt{}$ (first conjunct) as long as the interleavings split a trace $t_1 t_2$ that has already been mapped (second conjunct).

**TPC** is an instance of $\tau$ that maps source hypersafety to target hypersafety (and therefore, safety to safety), thus our theory can be instantiated for the preservation of these classes of hyperproperties as well.

## 8 RELATED WORK

We already discussed how our results relate to some existing work in correct compilation [? ? ] and secure compilation [? ? ? ]. We also already mentioned that most of our definitions and results make no assumptions about the structure of traces. One result that partially relies on the structure of traces is Theorem 6.3, that refers to *finite prefix m*, suggesting traces should be some sort of sequences of events (or states), as customary when one wants to refer to safety properties [? ]. Without a notion of finite prefix only RSC$^\sim$ may look different but both $\text{RTP}^{Safe \circ \tilde{\tau}}$ and $\text{RSP}^{\tilde{\sigma}}$ are

---

[19]They call those safety properties monitors since they focus on safety [? ] and indicate s with M and t with M.

[20]Technically, they assume a set of $\sqrt{}$ actions, but for this analogy a single action suffices.

trace agnostic as in general safety properties can be defined as the closed sets of any topology on traces [? ].

Even for reasoning about safety, hypersafety, or arbitrary hyperproperties, traces can therefore be values, sequences of program states, or of input output events, or even the recently proposed *interaction trees* [? ]. In the latter case we believe that the compilation from IMP to ASM proposed by ? ] can be seen as an instance of HC~, for the relation they call "trace equivalence."

**Compilers Where Our Work Could Be Useful.** Our work should be broadly applicable to understanding the guarantees provided by many verified compilers. For instance, ? ] recently proposed a CompCert variant that compiles all the way down to machine code, and it would be interesting to see if the model at the end of §4.1 applies there too. This and many other verified compilers [? ? ? ? ] beyond CakeML [? ] deal with resource exhaustion and it would be interesting to also apply the ideas of §4.2 to them. [CH: There is more related work on resource exhaustion that we should be relating at some point: (1) CompCertTSO [? ] and Gil's PLDI'15 [? ] paper on Integer-Pointer casts. From [? ]: "Therefore, CompCertTSO [? ] models out of memory as the empty set of behaviors, that is no behavior. However, what happens if there were I/O events prior to the out-of-memory error? Discarding I/O events before running out of memory is absurd: the target program should always perform a prefix of the events the source program could have performed. To handle this, CompCert-TSO also observes partial behaviors, possibly before discarding behavior due to running out of memory: 4. A partial execution of the program that has produced a finite sequence of I/O events, e1, …, en, partial. As before, refinement is defined as inclusion of the set of (possibly partial) behaviors of the target program into that of the source program. In this paper, we follow CompCertTSO's approach of handling out of memory. However, unlike CompCertTSO, where only the target language can run out of memory, our source language can also run out of memory due to pointer-to-integer casts, as we explain below." (2) CompCertMC [? ] also deals with memory exhaustion (at least for the stack), but I don't think we ever understood how their top level theorem looks like? From what I still remember they push low-level memory information to the higher levels via an oracle "stackspace" that is produced by the compiler together with the binary, an idea they credit to [? ] (3) Mullen [? ] also deal with a flat finite address space, so they must also have a solution for running out of memory. Andrew's note says they also use oracles.] [JT: Should we talk about this here, or in the related work section? Not sure we can keep this section short if we want to discuss these papers]

? ] devised a correct compiler from an ML language to assembly using a cross-language logical relation to state their CC theorem. They do not have traces, though were one to add them, the logical relation on values would serve as the basis for the trace relation and therefore their result would attain CC~.[CH: Theirs is a *compositional* compiler result though.] [CH: Yale (NI preserving) [? ]]

Switching to more informative traces capturing the interaction between the program and the context is often used as a proof technique for secure compilation [? ? ? ]. Most of these results consider a cross-language relation, so they probably could be proved to attain one of the criteria from Figure 4.

[CH: also compositional compiler correctness works, not just Figure 2]

**Generalizations of Compiler Correctness.** The compiler correctness definition of ? ] was already general enough to account for trace relations, since it considered a translation between the semantics of the source program and that of the compiled program, which he called "decode" in his diagram, reproduced in Figure 5 (left). And even some of the more recent compiler correctness definitions preserve this kind of flexibility [? ]. While CC~ can be seen as an instance of a definition by ? ], we are not aware of any prior work that investigated the preservation of properties when the "decode
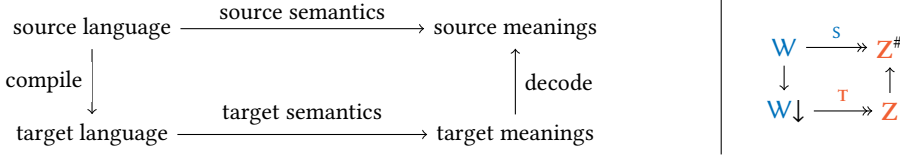
Fig. 5.  ? 's [? ] (left) and ? 's [? ] and ? 's [? ] (right) compiler correctness diagrams

translation" is neither the identity nor a bijection, and source properties need to be re-interpreted as target ones and vice versa.

**Correct Compilation and Galois Connections.** ? ] and ? ] expressed a strong variant of compiler correctness using the diagram of Figure 5 (right). They require that compiled programs *parallel* the computation steps ($\twoheadrightarrow$) of the original source programs, which can be proven showing the existence of a *decompilation* map # that makes the diagram commute, or equivalently, the existence of an adjoint for $\downarrow$ ($W \leq W' \iff W \twoheadrightarrow W'$ for both source and target). The "parallel" intuition can be formalized as an instance of CC~. Take source and target traces to be finite or infinite sequences of program states (maximal trace semantics [? ]), and relate them exactly like ? ] and ? ].

**Translation Validation.** Translation validation is an important alternative to proving that all runs of a compiler are correct as it can be more easily applied to realistic compilers. An interesting work about translation validation of security properties has been recently proposed by ? ]. They can handle many security properties expressible in terms of automata as long as source and target attackers and the observable traces are the same.
Instantiating the definition of any of the presented criteria with a particular program, one has translation validation criteria with the map $\tilde{\tau}$ describing the target property that is (robustly) satisfied once the translation is validated. For example one can consider

$$(tsv\text{~})\ \mathrm{CC}\text{~}(\mathsf{W}{\downarrow}) = \forall \mathsf{t}.\ \mathsf{W}{\downarrow} \rightsquigarrow \mathsf{t} \Rightarrow \exists \mathsf{s} \sim \mathsf{t}.\ \mathsf{W} \rightsquigarrow \mathsf{s}$$

$$(rtsv\text{~})\ \mathrm{RTC}\text{~}(\mathsf{C_T}\,[\,\mathsf{P}{\downarrow}]) = \forall \mathsf{t}.\ \mathsf{C_T}\,[\,\mathsf{P}{\downarrow}] \rightsquigarrow \mathsf{t} \Rightarrow \exists \mathsf{C_S}.\ \exists \mathsf{s} \sim \mathsf{t}.\ \mathsf{C_S}\,[\mathsf{P}] \rightsquigarrow \mathsf{s}$$

While the proof technique proposed by ? ] might be generalized for CC~($\mathsf{W}{\downarrow}$) – as long as *beh*($\mathsf{W}{\downarrow}$) and *beh*($\mathsf{W}$) can be expressed as one of the automata they can handle – they don't work for RTC~($\mathsf{C_T}\,[\,\mathsf{P}{\downarrow}]$) because of the existential in the conclusion.
     ? ] are instead considering translation validation criteria in the spirit of ($rtsv$~), their preliminary work only allows equality as trace relation, but should be subject to a generalization to the trace relating setting similar to the one we presented in this work.

**Proof Techniques.** We believe existing proof techniques (beyond the simulations discussed in Section 4.3.2) that have been devised to prove compiler correctness can also be employed to prove that a compiler attains any of the presented criteria. For example, cross-language binary logical relations can be used to relate two terms of two different languages when they 'behave the same' [? ? ? ]. Additionally, they can also be used when multiple programs 'behave the same' [? ] in a multilanguage semantics setting [? ]. Secure compilation results (which rely on the criteria of Section 6) can be proven using variations of the *backtranslation* proof technique [? ? ? ]. Presenting this proof techniques is beyond the scope of this paper, so we refer the interested reader to the work of ? ].

## 9  CONCLUSION AND FUTURE WORK

We have extended the property preservation view on compiler correctness to arbitrary trace relations, and believe that this will be useful for understanding the guarantees various compilers provide. An open question is whether, given a compiler, there exists a most precise $\sim$ relation for which this compiler is correct. As mentioned in §1, every compiler is $CC^\sim$ for some $\sim$, but under which conditions is there a most precise relation? In practice, more precision may not always be better though, as it may be at odds with compiler efficiency and may not align with more subjective notions of usefulness, leading to tradeoffs in the selection of suitable relations. Finally, another interesting direction for future work is studying whether using the relation to Galois connections allows to more easily compose trace relations for different purposes, say, for a compiler whose target language has undefined behavior, resource exhaustion, and side-channels. In particular, are there ways to obtain complex relations by combining simpler ones in a way that eases the compiler verification burden?

*Composition for Multipass Compilers.* For now, we can already informally argue about the correctness of a multipass compiler, where each step is proved correct for a possibly different trace relation. Concretely, assume $\downarrow_I^S$ is a compilation chain from a source language $S$ to an intermediate language $I$ and $\downarrow_T^I$ from the intermediate language $I$ to a target language $T$.[21] Assume given two relations between traces of these languages: $\sim_{S,I} \subseteq \text{Trace}_S \times \text{Trace}_I$ and $\sim_{I,T} \subseteq \text{Trace}_I \times \textbf{Trace}_T$, such that each compiler is proven to be *CC* w.r.t. the expected trace relation: $\downarrow_I^S \in CC^{\sim_{S,I}}$ and $\downarrow_T^I \in CC^{\sim_{I,T}}$.

Let us consider the source-to-target compiler $\downarrow_T^S$ that is derived of the composition of the two aforementioned compilers, so $\downarrow_T^S = \downarrow_T^I \circ \downarrow_I^S$. In this case, we obtain the expected result: the correctness of the whole compiler $\downarrow_T^S$ is derived from the individual compiler correctness proofs for each step.

$$CC^{(\sim_{I,T} \circ \sim_{S,I})} \equiv \forall W \forall t. \; W \downarrow_T^S \rightsquigarrow t \Rightarrow \exists s \; \sim_{I,T} \circ \sim_{S,I} t. \; W \rightsquigarrow s$$

where $s \sim_{i,t} \circ \sim_{s,i} t \iff \exists i \in \text{Trace}_I. \; s \sim_{S,I} i \wedge i \sim_{I,T} t.$

Generalising this kind of composition to compilers that attain different criteria is unclear. For example, if $\downarrow_I^S$ preserves arbitrary hyperproperties, but $\downarrow_T^I$ preserves 2-relational safety properties, what can we conclude for $\downarrow_T^S$? We leave investigating these interesting matters for future work.

## ACKNOWLEDGMENTS

## A  PROOFS

Proof of Theorem 2.6 (🐾). See Theorems rel_TC_$\tau$TP and rel_TC_$\sigma$TP in `TraceCriterion.v`, where the $TP^{\tilde{\tau}} \iff TP^{\tilde{\sigma}}$ part follows directly from Theorem 2.4.                    □

Proof of Lemma 2.7 (Trace relations $\cong$ Galois connections on trace properties).

? ] show that the existential image is a functor from the category of sets and relations to the category of predicate transformers, mapping a set $X \mapsto 2^X$ and a relation $\sim \subseteq X \times Y \mapsto \tilde{\tau} : 2^X \to 2^Y$.

---

[21]For the intermediate language we use a verbatim, emerald font.

They also show that such a functor is an isomorphism – hence bijective – when one considers only monotonic predicate transformers that have a – unique – upper adjoint. The universal image of $\sim$, $\tilde{\sigma}$, is the unique adjoint of $\tilde{\tau}$ ($\mathcal{K}$), hence $\sim \mapsto \tilde{\tau} \leftrightarrows \tilde{\sigma}$ is itself bijective.                                          □

PROOF OF THEOREM 2.8 (CORRESPONDENCE OF CRITERIA). For a trace relation $\sim$ and the Galois connection $\tilde{\tau} \leftrightarrows \tilde{\sigma}$, the result follows from Theorem 2.6. For a Galois connection $\tau \leftrightarrows \sigma$ and $\hat{\sim}$, use Lemma 2.7 to conclude that the existential and universal images of $\hat{\sim}$ coincide with $\tau$ and $\sigma$, respectively; the goal then follows from Theorem 2.6.                                          □

**Lemma A.1** (Special relations and consequences on the adjoints)**.** Let $X, Y$ be two arbitrary sets and $\sim \subseteq X \times Y$. Assume $\sim$ is a total and surjective map from $Y$ to $X$. Let $\alpha \leftrightarrows \gamma$ be its existential and universal image, i.e.

$$\tilde{\alpha} = \lambda\,\pi_X.\ \{y \mid \exists x \in \pi_X.\ x \sim y\}$$
$$\tilde{\gamma} = \lambda\,\pi_Y.\ \{x \mid \forall y.\ x \sim y \Rightarrow y \in \pi_Y\}$$

Then $\tilde{\gamma} = \lambda\,\pi_Y.\ \{x \mid \exists y \in \pi_Y.\ x \sim y\}$, and $\tilde{\gamma}$ is injective.

PROOF OF LEMMA A.1. See Lemma rel_total_surjective and rel_total_surjective_up_inj in Galois.v
                                          □

PROOF OF THEOREM 4.3 ($\mathcal{K}$). See Theorem correctness in `TypeRelationExampleInput.v`.   □

PROOF FOR LEMMA 4.4 (gensend $(\cdot, \cdot)$ WORKS). We proceed by induction on $\tau$ and then by induction on $\tau'$:

$\tau = \mathsf{N}$ **and** $\tau' = \mathsf{N}$ By canonicity we have that $\mathsf{r} = \langle \mathsf{n}, \mathsf{n}' \rangle$.
      gensend $(\cdot, \cdot)$ translates that into $\mathsf{send\ n; send\ n'}$.
      By Rule Sem-seq, that produces $\mathsf{t} = \mathsf{n; n'}$.
      We need to prove that $\langle \mathsf{n}, \mathsf{n}' \rangle \sim \mathsf{n; n'}$, which holds by Rule Trace-Rel-N-N.
$\tau = \mathsf{N}$ **and** $\tau' = \tau_1 \times \tau_2$ Analogous to the other cases, by IH and Rule Trace-Rel-N-M.
$\tau = \tau_1 \times \tau_2$ **and** $\tau' = \mathsf{N}$ Analogous to the other cases, by IH and Rule Trace-Rel-M-N.
$\tau = \tau_1 \times \tau_1'$ **and** $\tau' = \tau_2 \times \tau_2'$ So by canonicity $\mathsf{r} = \langle \langle \mathsf{r}_1, \mathsf{r}_1' \rangle, \langle \mathsf{r}_2, \mathsf{r}_2' \rangle \rangle$.
      By definition of gensend $(\cdot, \cdot)$:

$$\text{gensend}\,(\mathbf{x}, \tau \times \tau')$$
$$= \text{gensend}\,(\mathbf{x}, \tau).1;\ \text{gensend}\,(\mathbf{x}, \tau').2$$

      By the target reductions we know (gensend $(\mathbf{x}, \tau).1;$ gensend $(\mathbf{x}, \tau').2)[\mathsf{r}/\mathbf{x}] \rightsquigarrow \mathsf{is}_1; \mathsf{is}_2$, so by IH we have $\langle \mathsf{r}_1, \mathsf{r}_1' \rangle \sim \mathsf{is}_1$ and $\langle \mathsf{r}_2, \mathsf{r}_2' \rangle \sim \mathsf{is}_2$.
      We need to prove that $\langle \langle \mathsf{r}_1, \mathsf{r}_1' \rangle, \langle \mathsf{r}_2, \mathsf{r}_2' \rangle \rangle \sim \mathsf{is}_1; \mathsf{is}_2$, which holds by Rule Trace-Rel-M-M, for $\mathsf{i} = \langle \mathsf{r}_1, \mathsf{r}_1' \rangle$ and $\mathsf{i}' = \langle \mathsf{r}_2, \mathsf{r}_2' \rangle$.                                          □

PROOF OF THEOREM 4.5. Trivial induction on the typing derivation of $\mathsf{e}$, the only interesting case is the compilation of $\mathsf{send\ e}$ in the inductive cases.

**Inductive.** $\mathsf{e} = \mathsf{send\ e}$ By IH we have that if $(\vdash \mathsf{e} : \tau \times \tau')\!\downarrow \rightsquigarrow \mathsf{t}$ then $\exists \mathsf{s} \sim \mathsf{t}$ and $\mathsf{e} \rightsquigarrow \mathsf{t}$.
      By definition of $(\cdot)\!\downarrow$ and of $\rightsquigarrow$ we need to prove that if

$$\mathsf{let\ x} = (\vdash \mathsf{e} : \tau \times \tau')\!\downarrow\ \mathsf{in\ gensend}\,(\mathbf{x}, \tau \times \tau') \rightsquigarrow \mathsf{t}$$

      Then $\mathsf{send\ e} \rightsquigarrow \mathsf{s}$ and $\mathsf{s} \sim \mathsf{t}$.
      The reductions proceed as follows in the target:

$$\frac{(\vdash e : \tau \times \tau')\!\downarrow \;\rightsquigarrow\; \langle \mathsf{is}, (\vdash r : \tau \times \tau')\!\downarrow\rangle \qquad \mathsf{gensend}\,(\mathbf{x}, \tau \times \tau')[(\vdash r : \tau \times \tau')\!\downarrow/\mathbf{x}] \;\rightsquigarrow\; \langle \mathsf{is}', \mathbf{r}'\rangle}{\mathbf{let\; x}= (\vdash e : \tau \times \tau')\!\downarrow \mathbf{\;in\;} \mathsf{gensend}\,(\mathbf{x}, \tau \times \tau') \;\rightsquigarrow\; \langle \mathsf{is} \cdot \mathsf{is}', \mathbf{r}'\rangle}$$

In the source we have $\quad \dfrac{e \;\rightsquigarrow\; \langle \mathsf{is}, \mathsf{r}\rangle}{\mathsf{send}\; e \;\rightsquigarrow\; \langle \mathsf{is} \cdot \mathsf{r}, \mathsf{r}\rangle}$

By IH we have that $\mathsf{is} \sim \mathbf{is}$.

By Rule Trace-Rel-Single, to prove that $\mathsf{is}; \mathsf{r} \sim \mathbf{is}; \mathbf{is}'$ we need to prove that $\mathsf{is} \sim \mathbf{is}'$.

By Lemma 4.4 (gensend $(\cdot, \cdot)$ works) we have that $\mathsf{r} \sim \mathbf{is}'$, so this case holds.          □

---

PROOF OF THEOREM 5.1. First of all we show that $\phi^{\#}$ is an *uco*, the proof for $\rho^{\#}$ is the same.

**Monotonicity.** $\phi^{\#}$ is composition of monotonic functions, hence it is itself monotonic.

**Idempotence.** We have to show that for $\pi_{\mathrm{T}}$, $\phi^{\#}(\phi^{\#}(\pi_{\mathrm{T}})) = \phi^{\#}(\pi_{\mathrm{T}})$, that unfolding the definition means

$$g^{\circ} \circ \phi \circ f^{\circ} \circ g^{\circ} \circ \phi \circ f^{\circ}(\pi_{\mathrm{T}}) = g^{\circ} \circ \phi \circ f^{\circ}(\pi_{\mathrm{T}})$$

For the inclusion "$\subseteq$",

$$g^{\circ} \circ \phi \circ f^{\circ} \circ g^{\circ} \circ \phi \circ f^{\circ}(\pi_{\mathrm{T}}) \subseteq g^{\circ} \circ \phi \circ \phi \circ f^{\circ}(\pi_{\mathrm{T}}) = g^{\circ} \circ \phi \circ f^{\circ}(\pi_{\mathrm{T}})$$

the inclusion holds because $f^{\circ} \circ g^{\circ}(x) \subseteq x$ and the equality comes from idempotency of $\phi$.

For the inclusion "$\supseteq$",

$$g^{\circ} \circ \phi \circ f^{\circ} \circ g^{\circ} \circ \phi \circ f^{\circ}(\pi_{\mathrm{T}}) \supseteq g^{\circ} \circ \phi \circ f^{\circ} \circ g^{\circ} \circ f^{\circ}(\pi_{\mathrm{T}}) = g^{\circ} \circ \phi \circ f^{\circ}(\pi_{\mathrm{T}})$$

the inclusion comes from $\phi(f^{\circ}(\pi_{\mathrm{T}})) \supseteq f^{\circ}(\pi_{\mathrm{T}})$ by extensiveness of $\phi$, and the equality from $f^{\circ} \circ g^{\circ} \circ f^{\circ} = f^{\circ}$.

**Extensiveness.** We have to show that $\pi^{\#}(\pi_{\mathrm{T}}) \supseteq \pi_{\mathrm{T}}$.

$$\pi^{\#}(\pi_{\mathrm{T}}) = g^{\circ} \circ \phi \circ f^{\circ}(\pi_{\mathrm{T}}) \supseteq g^{\circ} \circ f^{\circ}(\pi_{\mathrm{T}}) \supseteq \pi_{\mathrm{T}}$$

The first inclusion is due to extensiveness of $\phi$, the second by $g^{\circ}$ being the upper adjoint of $f^{\circ}$.

For the statement of the theorem to hold, assume $\mathsf{W} \models ANI^{\rho}_{\phi}$ and $\mathsf{W}\!\downarrow\!\rightsquigarrow t_1, t_2$ with $\phi^{\#}(t_1^{\circ}) = \phi^{\#}(t_2^{\circ})$, we have to show that $\rho^{\#}(t_1^{\bullet}) = \rho^{\#}(t_1^{\bullet})$.

By CC$^{\sim}$ there exists $s_1 \sim t_1$ and $s_2 \sim t_2$ such that $\mathsf{W} \rightsquigarrow s_1, s_2$. As a preliminary, apply Lemma A.1 to the relations $\overset{\circ}{\sim} \circ \, swap$ and deduce $g^{\circ}$ is injective. Notice also that by functionality and totality, of $\overset{\circ}{\sim}$ and of $\overset{\bullet}{\sim}$, $f^{\circ}(t_1^{\circ}) = \{s_1^{\circ}\}$ and $f^{\bullet}(t_1^{\bullet}) = \{s_1^{\bullet}\}$ and a similar fact holds for $s_2$ and $t_2$.

$$
\begin{aligned}
\phi^{\#}(t_1^{\circ}) = \phi^{\#}(t_2^{\circ}) \Rightarrow &\qquad\qquad [\text{ definition of } \phi^{\#}] \\
g^{\circ} \circ \phi \circ f^{\circ}(t_1^{\circ}) = g^{\circ} \circ \phi \circ f^{\circ}(t_2^{\circ}) \Rightarrow &\qquad\qquad [g^{\circ} \text{ injective}] \\
\phi \circ f^{\circ}(t_1^{\circ}) = \phi \circ f^{\circ}(t_2^{\circ}) \Rightarrow &\qquad\qquad [f^{\circ}(t_i^{\circ}) = s_i^{\circ}\; i = 1, 2] \\
\phi(s_1^{\circ}) = \phi(s_2^{\circ}) \Rightarrow &\qquad\qquad [\mathsf{W} \models ANI^{\rho}_{\phi}] \\
\rho(s_1^{\bullet}) = \rho(s_2^{\bullet}) \Rightarrow &\qquad\qquad [s_i^{\bullet} = f^{\bullet}(t_i^{\bullet})\; i = 1, 2] \\
\rho \circ f^{\bullet}(t_1^{\bullet}) = \rho \circ f^{\bullet}(t_2^{\bullet}) \Rightarrow &\qquad\qquad [\text{ functionality of } g^{\bullet}] \\
g^{\bullet} \circ \rho \circ f^{\bullet}(t_1^{\bullet}) = g^{\bullet} \rho \circ f^{\bullet}(t_2^{\bullet}) \Rightarrow &\qquad\qquad [\text{ definition of } \rho^{\#}] \\
\rho^{\#}(t_1^{\bullet}) = \rho^{\#}(t_2^{\bullet}), &
\end{aligned}
$$

so that $\mathsf{W}\!\downarrow\, \models ANI^{\rho^{\#}}_{\phi^{\#}}$.

We now show that if $\overset{\bullet}{\sim}$ is surjective, i.e., $g^{\bullet}$ injective, $ANI^{\rho^{\#}}_{\phi^{\#}} \subseteq Cl_{\subseteq} \circ \tilde{\tau}(ANI^{\rho}_{\phi})$.

Let $\pi_T \in ANI_{\phi^\#}^{\rho^\#}$, we show that $\pi_T \subseteq \tilde\tau(\pi_S)$ for some $\pi_S \in ANI_\phi^\rho$.

The source property $\pi_S = \{s \mid \exists t \in \pi_T.\ s \sim t\} = f(\pi_T)$ is such that $\pi_T \subseteq \tilde\tau(\pi_S)$. We only need to show $\pi_S \in ANI_\phi^\rho$. Let $s_1, s_2 \in \pi_S$,

$$\phi(s_1^\circ) = \phi(s_2^\circ) \Rightarrow \qquad\qquad [\text{by } f^\circ(t^\circ) = s^\circ \text{ for some } t \in \pi_T]$$
$$\phi(f^\circ(t_1^\circ)) = \phi(f^\circ(t_2^\circ)) \Rightarrow \qquad\qquad [g^\circ \text{ is a function}]$$
$$g^\circ(\phi(f^\circ(t_1^\circ))) = g^\circ(\phi(f^\circ(t_2^\circ))) \Rightarrow \qquad\qquad [\text{by definition of } \phi^\#]$$
$$\phi^\#(t_1^\circ) = \phi^\#(t_2^\circ) \Rightarrow \qquad\qquad [\pi_T \in ANI_{\phi^\#}^{\rho^\#}]$$
$$\rho^\#(t_1^\bullet) = \rho^\#(t_2^\bullet) \Rightarrow \qquad\qquad [\text{definition of } \rho^\#]$$
$$g^\bullet(\rho(f^\bullet(t_1^\bullet))) = g^\bullet(\rho(f^\bullet(t_2^\bullet))) \Rightarrow \qquad\qquad [\text{by injectivity of } g^\bullet]$$
$$\rho(f^\bullet(t_1^\bullet)) = \rho(f^\bullet(t_2^\bullet)) \Rightarrow \qquad\qquad [f^\bullet(t_i) = s_i^\bullet, i = 1, 2]$$
$$\rho(s_1^\bullet) = \rho(s_2^\bullet),$$

that shows $\pi_S \in ANI_\phi^\rho$ and concludes the proof.

□

PROOF OF THEOREM 5.2. Assume $W \models ANI_\phi^\rho$ and $W\downarrow \leadsto t_1, t_2$ with $\phi^\#(t_1^\circ) = \phi^\#(t_2^\circ)$. We have to show that $\rho^\#(t_1^\bullet) = \rho^\#(t_1^\bullet)$, for an arbitrary $\rho^\#$ that satisfies the condition

$$H \equiv \forall s\ t.\ s^\bullet \overset{\cdot}{\sim} t^\bullet \Rightarrow \rho^\#(\tilde\tau^\bullet(\rho(s^\bullet))) = \rho^\#(t^\bullet).$$

By CC$^\sim$ there exists $s_1 \sim t_1$ and $s_2 \sim t_2$ such that $W \leadsto s_1, s_2$. As a preliminary, recall that Lemma A.1 ensures $g^\circ$ is injective. Morevoer notice that by functionality and totality, of $\overset{\cdot}{\sim}$, $f^\circ(t_1^\circ) = \{s_1^\circ\}$ and $f^\circ(t_2^\circ) = \{s_2^\circ\}$.

$$\phi^\#(t_1^\circ) = \phi^\#(t_2^\circ) \Rightarrow \qquad\qquad [\text{ definition of } \phi^\#]$$
$$g^\circ \circ \phi \circ f^\circ(t_1^\circ) = g^\circ \circ \phi \circ f^\circ(t_2^\circ) \Rightarrow \qquad\qquad [g^\circ \text{ injective}]$$
$$\phi \circ f^\circ(t_1^\circ) = \phi \circ f^\circ(t_2^\circ) \Rightarrow \qquad\qquad [f^\circ(t_i^\circ) = s_i^\circ\ i = 1, 2]$$
$$\phi(s_1^\circ) = \phi(s_2^\circ) \Rightarrow \qquad\qquad [W \models ANI_\phi^\rho]$$
$$\rho(s_1^\bullet) = \rho(s_2^\bullet) \Rightarrow \qquad\qquad [\text{ functionality of } \tilde\tau^\bullet]$$
$$\tilde\tau^\bullet(\rho(s_1^\bullet)) = \tilde\tau^\bullet(\rho(s_2^\bullet)) \Rightarrow \qquad\qquad [\text{ by functionality of } \rho^\#]$$
$$\rho^\#(\tilde\tau^\bullet(\rho(s_1^\bullet))) = \rho^\#(\tilde\tau^\bullet(\rho(s_2^\bullet))) \Rightarrow \qquad\qquad [\text{ by condition } H]$$
$$\rho^\#(t_1^\bullet) = \rho^\#(t_2^\bullet)$$

so that $W\downarrow \models ANI_{\phi^\#}^{\rho^\#}$.

□

PROOF OF THEOREM 5.3. Assume $W \models ANI_{\phi^\#}^{\rho^\#}$ and $W\downarrow \leadsto t_1, t_2$ with $\phi(t_1^\circ) = \phi(t_2^\circ)$ and $\phi$ satisfying the condition $H \equiv \forall s\ t.\ s^\circ \overset{\cdot}{\sim} t^\circ \Rightarrow \phi(t^\circ) = \phi(\tilde\tau^\circ(s^\circ))$. We have to show that $\rho(t_1^\bullet) = \rho(t_2^\bullet)$. By CC$^\sim$ there exists $s_1 \sim t_1$ and $s_2 \sim t_2$ such that $W \leadsto s_1, s_2$. As a preliminary, recall that Lemma A.1 ensures $\tilde\sigma^\bullet$ is injective. Morevoer notice that by functionality and totality, of $\overset{\cdot}{\sim}$, $\tilde\tau^\bullet(s_1^\bullet) = \{t_1^\bullet\}$ and $\tilde\tau^\bullet(s_2^\bullet) = \{t_2^\bullet\}$.

$$\phi(t_1^\circ) = \phi(t_2^\circ) \Rightarrow \qquad\qquad [\text{ by } H]$$
$$\phi(\tilde\tau^\circ(s_1^\circ)) = \phi(\tilde\tau^\circ(s_2^\circ)) \Rightarrow \qquad\qquad [\text{ functionality of } \tilde\sigma^\circ]$$

$$\tilde{\sigma}^{\circ}(\phi(\tilde{\tau}^{\circ}(s_1^{\circ}))) = \tilde{\sigma}^{\circ}(\phi(\tilde{\tau}^{\circ}(s_1^{\circ}))) \Rightarrow \qquad [\text{ definition of } \phi^{\#}]$$

$$\phi^{\#}(s_1^{\circ}) = \phi^{\#}(s_2^{\circ}) \Rightarrow \qquad [W \models ANI_{\phi^{\#}}^{\phi^{\#}}]$$

$$\rho^{\#}(s_1^{\bullet}) = \rho^{\#}(s_2^{\bullet}) \Rightarrow \qquad [\text{ by definition of } \rho^{\#}]$$

$$\tilde{\sigma}^{\bullet}(\rho(\tilde{\tau}^{\bullet}(s_1^{\bullet}))) = \tilde{\sigma}^{\bullet}(\rho(\tilde{\tau}^{\bullet}(s_2^{\bullet}))) \Rightarrow \qquad [\text{ injectivity of } \tilde{\sigma}^{\bullet}]$$

$$\rho(\tilde{\tau}^{\bullet}(s_1^{\bullet})) = \rho(\tilde{\tau}^{\bullet}(s_2^{\bullet})) \Rightarrow \qquad [\tilde{\tau}^{\bullet}(s_i^{\bullet}) = \{t_i^{\bullet}\}\ i = 1, 2]$$

$$\rho(t_1^{\bullet}) = \rho(t_2^{\bullet})$$

so that $W{\downarrow} \models ANI_{\phi}^{\rho}$.

□

PROOF OF THEOREM 6.1 (🐔). Theorems rel_RTC_$\tau$RTP and rel_RTC_$\sigma$RTP in
`RobustTraceCriterion.v`. □

PROOF OF THEOREM 6.3 (🐔). Theorems tilde_RSC_$\sigma$RSP and tilde_RSC_Cl_$\tau$RTP in
`RobustSafetyCriterion.v`. □

PROOF OF THEOREM 6.5 (🐔). Lemmas $\sigma$RHP_rel_RHC and rel_RHC_$\sigma$RHP and Theorem
rel_RHC_$\tau$RHP in `RobustHyperCriterion.v`. □

PROOF OF THEOREM 7.1 (🐔). (See theorem extra_target_RTCt in `MoreTargetEventsExample.v`,
mechanizing a slightly simplified model.) By definition of RTC$^{\sim}$ we need to find a source context
and source trace given a source program, target context and target trace related by compilation and
program semantics: This instantiation is simple since the trace relation is a *function* from target
traces to source traces, and it is easy to clean target contexts to produce equivalent source context
without target-only events. The proof is a trivial instance of *precise, context-based backtranslation*
[? ? ? ? ], aided by a few straightforward lemmas and where the case of function calls is guaranteed
to terminate by the language. □

PROOF OF THEOREM 3.6 (🐔). Theorems tilde_SC_$\sigma$SP and tilde_SC_Cl_$\tau$TP in
`SafetyCriterion.v`. □

PROOF OF THEOREM 3.7. For the implication from left to right, assume $W \models H$. By CC$^{=}$ have
$\mathbf{beh}(W{\downarrow}) = \mathrm{beh}(W)$, so that $W{\downarrow} \models H$ as well.
For the implication from right to left, instantiate HP with the hyperproperty $\{\mathrm{beh}(W)\}$, for a given
$W$, and deduce that $W{\downarrow} \models \{\mathrm{beh}(W)\}$ i.e., $\mathbf{beh}(W{\downarrow}) = \mathrm{beh}(W)$. □

PROOF OF THEOREM 3.8 (🐔). Theorems rel_HC_$\tau$HP, rel_HC_$\sigma$HP and $\sigma$HP_rel_HC in
`HyperCriterion.v`. □