# Abstract Symbolic Automata

## Mixed syntactic/semantic similarity analysis of executables

Mila Dalla Preda[1]     Roberto Giacobazzi[1,3]     Arun Lakhotia[2]     Isabella Mastroeni[1]

[1]University of Verona     [2]University of Louisiana     [3] Irdeto Canada

mila.dallapreda@univr.it, roberto.giacobazzi@univr.it, arun@louisiana.edu, isabella.mastroeni@univr.it

## Abstract

We introduce a model for mixed syntactic/semantic approximation of programs based on symbolic finite automata (SFA). The edges of SFA are labeled by predicates whose semantics specifies the denotations that are allowed by the edge. We introduce the notion of abstract symbolic finite automaton (ASFA) where approximation is made by abstract interpretation of symbolic finite automata, acting both at syntactic (predicate) and semantic (denotation) level. We investigate in the details how the syntactic and semantic abstractions of SFA relate to each other and contribute to the determination of the recognized language. Then we introduce a family of transformations for simplifying ASFA. We apply this model to prove properties of commonly used tools for similarity analysis of binary executables. Following the structure of their control flow graphs, disassembled binary executables are represented as (concrete) SFA, where states are program points and predicates represent the (possibly infinite) I/O semantics of each basic block in a constraint form. Known tools for binary code analysis are viewed as specific choices of symbolic and semantic abstractions in our framework, making symbolic finite automata and their abstract interpretations a unifying model for comparing and reasoning about soundness and completeness of analyses of low-level code.

***Categories and Subject Descriptors***   D.3.1 [*Programming Languages*]: Formal Definitions and Theory;   F.3.2 [*Logics and Meanings of Programs*]: Semantics of Programming Languages—program analysis

***General Terms***   Languages.

***Keywords***   Symbolic automata, abstract interpretation.

## 1.  Introduction

**The problem.** Similarity analysis is a key component in mining and understanding huge software enclaves, including code, e,g., coming from malware repositories, specifications, analyses and other heterogeneous meta-data. This is particularly relevant when dealing with binary executables, which, besides representing a large portion of existing malware, also represent a highly malleable often hard to analyze carrier. This is due to its unstructured nature,

allowing self-modification, overlapping instructions, and untyped computations where data and code coexist without any predefined (static) boundary.

In order to mine both semantic meanings and syntactic patterns from programs, existing tools for similarity analysis of binary executables always employ mixed syntactic/symbolic and semantic representations of programs. At syntactic level properties concerning the control flow graph, such as in *BinHunt* [15] and *BinDiff* [12, 25], or feature vectors concerning sequences of instructions, are used together with graph-isomorphism, sequence comparison algorithms, and hash functions for extracting structural similarities in code. At semantic level, more advanced semantic properties such as those extracted from symbolic executions, dynamic analysis and emulation, such as those used in *BinJuice* [18] and *BinHunt* [15], are employed for bypassing semantic preserving code transformations for code obfuscation, e.g., for similarity analysis in malware detection. The use of mixed syntactic/semantic representation of code in similarity analysis is becoming a good practice because pure semantic similarity is too complex and often undecidable while pure syntactic similarities is too imprecise and prone to false negatives due to code obfuscation techniques. This is precisely what happens in most known tools and methods for dissecting and comparing programs in order to extract semantic similarities from syntactically different code. However, none of these tools have a formal semantic model in which relative precision and soundness can be formally proved. This paper is intended to fill this gap.

**Our contribution.** We attack this problem by observing that most known methods employed in similarity analysis of disassembled binaries can be seen as peculiar abstract interpretations of *symbolic finite state automata* (SFA). Symbolic finite automata, introduced in [23] and further developed in [8, 9], provide the ideal formal setting in order to treat within the same model the abstraction of both the syntactic structure of programs and their intended semantics.

SFA have been introduced as an extension of traditional finite state automata for modeling languages with a potential infinite alphabet. Transitions in SFA are therefore modeled as constraints interpreted in a given Boolean algebra, providing the semantic interpretation of constraints, and therefore the (potentially infinite) structural components of the language recognized (see [9, 23]).

Our main contribution is the introduction of the notion of *abstract symbolic finite automaton*, where approximation is made by abstract interpretation of standard SFA. Abstract interpretation here acts both at syntactic (predicate), topological (graph), and semantic (denotation) level. We investigate in details how the syntactic, topological, and semantic abstractions of SFA relate to each other and interfere when automata, at different levels of abstractions, are compared with respect to their recognized language.

The abstraction respectively on syntactic predicates and semantic structures corresponds precisely to the abstract interpretation of the underlying Boolean algebra of a concrete SFA $M$, resulting in a

different SFA $A$ whose language recognized is an over approximation of the language of $M$. The key aspect here is to maintain a relative compatibility between syntactic abstractions on predicates and constraint formulae and the abstractions of their semantics. This intuitively means that the approximate predicates and their interpretation provide, one over the others, coherent partitions of objects (respectively interpretations and predicates).

Topological abstraction means instead changing the graph structure of SFA, yet keeping correctness, namely providing an over approximation of the recognized language of $M$. This is achieved by generalizing a minimization algorithm proposed in [9] with respect to a family of equivalence relations on SFA states. The result is a simplification of $M$ which is still correct in the sense of abstract interpretation with respect to $M$.

Abstract SFA provide a general enough model for representing syntactic and semantic properties of arbitrary programming languages. We apply our model in the attempt to formalize and prove properties of two commonly used tools for similarity analysis of binary executables, notably *BinJuice* [18] and *BinDiff*. Following the structure of their control flow graphs, disassembled binary executables are represented as (concrete) SFA, where states are program points between basic blocks and predicates represent the (possibly infinite) I/O semantics of each basic block in a constraint form. Tools for binary-level similarity analysis are then formalized as abstract interpretations of these concrete SFA. By studying the properties of the corresponding abstractions we can provide a first unifying model for formally proving properties for these tools. Moreover, our model suggests potential refinements of similarity analyses for disassembled binaries such as the possibility of extracting minimal SFA from binaries as canonical signatures for code fragments.

## 2. Preliminaries

***Mathematical Notation.*** Given two sets $S$ and $T$, we denote with $\wp(S)$ the powerset of $S$, $\wp^{\mathrm{re}}(S)$ the set of recursive enumerable (r.e.) subsets of $S$, with $S \smallsetminus T$ the set-difference between $S$ and $T$, with $S \subset T$ strict inclusion and with $S \subseteq T$ inclusion. $S^*$ denotes the set of all finite sequences of elements in $S$. A set $L$ with ordering relation $\leq$ is a poset and it is denoted as $\langle L, \leq \rangle$. A poset $\langle L, \leq \rangle$ is a lattice if $\forall x.y \in L$ we have that $x \vee y$ and $x \wedge y$ belong to $L$. A lattice $\langle L, \leq \rangle$ is complete when for every $X \subseteq L$ we have that $\bigvee X, \bigwedge X \in L$. As usual a complete lattice $L$, with ordering $\leq$, least upper bound (lub) $\vee$, greatest lower bound (glb) $\wedge$, greatest element (top) $\top$, and least element (bottom) $\bot$ is denoted by $\langle L, \leq, \vee, \wedge, \top, \bot \rangle$. Given $f : S \longrightarrow T$ and $g : T \longrightarrow Q$ we denote with $g \circ f : S \longrightarrow Q$ their composition, i.e., $g \circ f = \lambda x.g(f(x))$. $f : L \longrightarrow D$ on complete lattices is *additive (co-additive)* if for any $Y \subseteq L, f(\vee_L Y) = \vee_D f(Y) \ (f(\wedge_L Y) = \wedge_D f(Y))$. Continuity holds when $f$ preserves *lubs*'s of chains. Co-continuity is dually defined. For a continuous function $f$: $lfp(f) = \bigwedge \{ x \mid x = f(x) \} = \bigvee_{n \in \mathbb{N}} f^n(\bot)$ where $f^0(\bot) = \bot$ and $f^{n+1}(\bot) = f(f^n(\bot))$.

***Abstract Interpretation.*** Abstract domains can be equivalently formalized either as Galois connections or closure operators on a given concrete domain which is a complete lattice $C$ (cf. [4]). Let $C$ and $A$ be complete lattices, a pair of monotone functions $\alpha : C \longrightarrow A$ and $\gamma : A \longrightarrow C$ forms a *Galois connection (GC)* between $C$ and $A$ if for every $x \in C$ and $y \in A$ we have $\alpha(x) \leq_A y \Leftrightarrow x \leq_C \gamma(y)$. $\alpha$ (resp. $\gamma$) is the *left-adjoint* (resp. *right-adjoint*) to $\gamma$ (resp. $\alpha$) and it is additive (resp. co-additive). If $\langle \alpha, \gamma \rangle$ is a GC between $C$ and $A$ then $\gamma \circ \alpha \in uco(C)$. If $\rho \in uco(C)$ then $\langle \rho, \mathtt{id} \rangle$ is a CG between $C$ and $\rho(C)$. Given an additive (resp. co-additive) function $\alpha$ (resp. $\gamma$) we have a GC $\langle \alpha, \alpha^+ \rangle$ (resp. $\langle \gamma^-, \gamma \rangle$) by considering its right (resp. left) adjoint

$\alpha^+ = \lambda x. \bigvee \{ y \mid \alpha(y) \leq x \}$ (resp. $\gamma^- = \lambda x. \bigwedge \{ y \mid x \leq \gamma(y) \}$). An *upper closure operator* (or simply a *closure*) on a poset $\langle L, \leq \rangle$ is an operator $\rho : L \longrightarrow L$ which is monotone, idempotent, and extensive (i.e., $x \leq \rho(x)$). We denote with $uco(L)$ the set of all closure operators on the poset $L$. If $C$ is a complete lattice, then $\langle uco(C), \sqsubseteq, \sqcup, \sqcap, \lambda x.\ C, \mathtt{id} \rangle$ forms a complete lattice [24], which is the set of all possible abstractions of $C$, where the bottom is $\mathtt{id} = \lambda x.x$ and for every $\rho, \eta \in uco(C)$, $\rho$ is *more concrete than* $\eta$ iff $\rho \sqsubseteq \eta$ iff $\forall y \in C.\ \rho(y) \leq \eta(y)$ iff $\eta(C) \subseteq \rho(C)$, $(\sqcap_{i \in I} \rho_i)(x) = \wedge_{i \in I} \rho_i(x)$; $(\sqcup_{i \in I} \rho_i)(x) = x$ iff $\forall i \in I.\ \rho_i(x) = x$. $\rho \in uco(C)$ is disjunctive when $\rho(C)$ is a join-sublattice of $C$ which holds iff $\rho$ is additive (cf. [4]). $\rho \in uco(\wp(C))$ is *partitioning* (or induces a partition) if it is additive and $\{ \rho(\{c\}) \}_{c \in C}$ is a partition of $C$ [17]. If $\rho \in uco(\wp(C))$ then the most abstract partitioning closure containing $\rho$:

$$ \Pi(\rho) \stackrel{\text{def}}{=} \bigsqcup \left\{ \beta \in uco(\wp(C)) \mid \beta \sqsubseteq \rho \wedge \beta \text{ is partitioning} \right\}. $$

The key aspect of partitioning closures is that they preserve the structure of Boolean algebras.

If $f : C \longrightarrow C$ is a continuous function and $\rho \in uco(C)$ is an abstraction, then $f$ always has a *best correct approximation* in $\rho(C)$ which is $f^\rho \stackrel{\text{def}}{=} \rho \circ f \circ \rho$. Any approximation $f^\sharp : \rho(C) \longrightarrow \rho(C)$ of $f$ in $\rho(C)$ is *sound* if $f^\rho \sqsubseteq f^\sharp$. In this case we have the fixpoint soundness $\rho(lfp\, f) \leq lfp(f^\rho) \leq lfp(f^\sharp)$(cf. [3]). $f^\sharp$ is *complete* when $\rho \circ f = f^\sharp \circ \rho$ which holds iff $\rho \circ f = \rho \circ f \circ \rho$ (cf. [16]). Therefore the possibility of defining a complete approximation $f^\sharp$ of $f$ on some abstract domain $\rho$ only depends on $f$ and $\rho$. In this case we have: $\rho(lfp\, f) = lfp(f^\rho) = lfp(f^\sharp)$. In the following, for any semantics $\llbracket \cdot \rrbracket : \mathcal{S} \longrightarrow \mathcal{D}$ mapping syntactic objects in $\mathcal{S}$ into denotations in $\mathcal{D}$ such that $\llbracket \cdot \rrbracket$ is an element in the set of fixpoint semantics $\mathfrak{S} \subseteq \mathcal{S} \longrightarrow \mathcal{D}$ inductively defined as follows

$$ \mathfrak{S} \quad ::= \quad f : \mathcal{S} \longrightarrow \mathcal{D} \ \mid \ lfp(\mathfrak{S}) \ \mid \ \mathfrak{S} \circ \mathfrak{S} $$

and if $\rho \in uco(\mathcal{D})$, we denote by $\llbracket \cdot \rrbracket^\rho \in \mathfrak{S}^\rho \subseteq \mathcal{S} \longrightarrow \rho(\mathcal{D})$ the corresponding best correct approximation which is defined inductively on the structure of $\mathfrak{S}$ as follows:

$$ \mathfrak{S}^\rho \quad ::= \quad \rho \circ f \circ \rho \ \mid \ lfp(\mathfrak{S}^\rho) \ \mid \ \mathfrak{S}^\rho \circ \mathfrak{S}^\rho $$

It is known that $\llbracket \cdot \rrbracket^\rho$ is sound and, whenever $\rho$ is complete for the basic semantic operators $f$ defining $\llbracket \cdot \rrbracket \in \mathfrak{S}$, then $\llbracket \cdot \rrbracket^\rho$ is complete, i.e. for any $s \in \mathcal{S}$: $\rho(\llbracket s \rrbracket) = \llbracket s \rrbracket^\rho$ (cf. [4, 16]).

***Symbolic Finite Automata.*** Symbolic automata and finite state transducers have been introduced to deal with specifications involving a potentially infinite alphabet of symbols [8, 9, 23]. We follow [9] in specifying symbolic automata in terms of effective Boolean algebra. Consider an effective Boolean algebra $\mathcal{A} = \langle \mathfrak{D}_\mathcal{A}, \Psi_\mathcal{A}, \llbracket \cdot \rrbracket, \bot, \top, \wedge, \vee, \neg \rangle$, with domain elements in a r.e. set $\mathfrak{D}_\mathcal{A}$, a r.e. set of predicates $\Psi_\mathcal{A}$ closed under boolean connectives $\wedge, \vee$ and $\neg$. The semantic function $\llbracket \cdot \rrbracket : \Psi_\mathcal{A} \longrightarrow \wp(\mathfrak{D}_\mathcal{A})$ is a partial recursive function such that $\llbracket \bot \rrbracket = \varnothing$, $\llbracket \top \rrbracket = \mathfrak{D}_\mathcal{A}$, and $\forall \varphi, \phi \in \Psi_\mathcal{A}$ we have that $\llbracket \varphi \vee \phi \rrbracket = \llbracket \varphi \rrbracket \cup \llbracket \phi \rrbracket$, $\llbracket \varphi \wedge \phi \rrbracket = \llbracket \varphi \rrbracket \cap \llbracket \phi \rrbracket$, and $\llbracket \neg \varphi \rrbracket = \mathfrak{D}_\mathcal{A} \smallsetminus \llbracket \varphi \rrbracket$. In the following we abuse notation by denoting with $\llbracket \cdot \rrbracket$ also its additive lift to $\wp(\Psi_\mathcal{A})$, i.e., for any $\Phi \in \wp(\Psi_\mathcal{A})$: $\llbracket \Phi \rrbracket = \{ \llbracket \varphi \rrbracket \mid \varphi \in \Phi \}$. For $\varphi \in \Psi_\mathcal{A}$ we write $IsSat(\varphi)$ when $\llbracket \varphi \rrbracket \neq \varnothing$ and say that $\varphi$ is *satisfiable*. $\mathcal{A}$ is decidable if $IsSat$ is decidable.

DEFINITION 2.1. *A symbolic automaton (SFA) is $\langle \mathcal{A}, Q, q_0, F, \Delta \rangle$ where $\mathcal{A}$ is an effective Boolean algebra, $Q$ is a finite set of states, $q_0 \in Q$ is the initial state, $F \subseteq Q$ is the set of final states and $\Delta \subseteq Q \times \Psi_\mathcal{A} \times Q$ is a finite set of transitions.*

A transition in $M = \langle \mathcal{A}, Q, q_0, F, \Delta \rangle$ labeled $\varphi$ from state $p$ to state $q$, $(p, \varphi, q) \in \Delta$ is often denoted $p \xrightarrow{\varphi} q$. $\varphi$ is called the *guard*

of the transition. An $a$-move of a SFA $M$ is a transition $p \xrightarrow{\varphi} q$ such that $a \in \llbracket \varphi \rrbracket$, also denoted $p \xrightarrow{a} q$. The language recognized by a state $q \in Q$ in $M$ is defined as:

$$\mathscr{L}_q(M) = \left\{ a_1, \ldots, a_n \in \mathfrak{D}_{\mathcal{A}} \;\middle|\; \begin{array}{l} \forall 1 \leq i \leq n.\ p_{i-1} \xrightarrow{a_i} p_i \\ p_0 = q,\ p_n \in F \end{array} \right\}$$

in this case, $\mathscr{L}(M) = \mathscr{L}_{q_0}(M)$. We assume *complete SFA*, namely where all states hold an out-going $a$-move, for any character $a \in \mathfrak{D}$. This can be simply achieved by adding a shaft-state $q_\perp \in Q$ such that $q_\perp \xrightarrow{\top} q_\perp \in \Delta$ and for all states $q$ lacking an out-going $a$-move, for $a \in \mathfrak{D}$, then $q \xrightarrow{\neg \beta} q_\perp \in \Delta$ with $\beta = \bigvee \left\{ \varphi \;\middle|\; q \xrightarrow{\varphi} p \wedge p \in Q \right\}$.

The following terminology holds for SFA: $M$ is *deterministic* whenever $p \xrightarrow{\varphi} q, p \xrightarrow{\beta} q' \in \Delta$: if $IsSat(\varphi \wedge \beta)$ then $q = q'$. $M$ is *clean* if for all $p \xrightarrow{\varphi} q \in \Delta$: $p$ is reachable from $q_0$ and $IsSat(\varphi)$. $M$ is *normalized* if for all $p, q \in Q$: there is at most one move from $p$ to $q$. $M$ is *minimal* if $M$ is deterministic, clean, normalized and for all $p, q \in Q$:

$$p = q \iff \mathscr{L}_q(M) = \mathscr{L}_p(M)$$

Given a SFA $M = \langle \mathcal{A}, Q, q_0, F, \Delta \rangle$ and $\equiv \subseteq Q \times Q$, we define the *quotient SFA* $M_{/\equiv} \stackrel{\text{def}}{=} \langle \mathcal{A}, Q', q_0', F', \Delta' \rangle$ as follows: $Q' = \left\{ [q]_\equiv \;\middle|\; q \in Q \right\}$, $\Delta' \subseteq Q' \times \Psi_{\mathcal{A}} \times Q'$ is such that $\Delta' = \left\{ ([q]_\equiv, \Phi, [q']_\equiv) \;\middle|\; (p, \Phi, q') \in \Delta, p \in [q]_\equiv \right\}$, $q_0' = [q_0]_\equiv$, and $F' = \left\{ [q]_\equiv \;\middle|\; q \in F \right\}$.

## 3. Abstracting Symbolic Automata

Approximating symbolic automata means building different automata recognizing an upper approximation of the original recognized language. This can be achieved by abstract interpretation of the underlying effective Boolean algebra $\mathcal{A}$ and by approximating the automaton's structure. When acting on the Boolean algebra we may either approximate the domain of denotations $\mathfrak{D}_{\mathcal{A}}$ where formulae and predicates are interpreted, or approximate the predicates in $\Psi_{\mathcal{A}}$ where formulae are built. In both cases we need to obtain as result an abstract effective Boolean algebra.

### 3.1 Abstract effective Boolean algebras

The duality of syntax and semantics is perfectly encoded in SFA by the underlying algebraic structure of effective Boolean algebras. They represent the universe of predicates and formulae (later called *syntax*) as well as the domain for their interpretation and semantics, providing the structure for expressing the language recognized by the given SFA. The abstraction of syntactic and semantic structures applies on sets of predicates and semantic structures representing, as usual in abstract interpretation, properties respectively of predicates and semantics. In the following $\mathcal{A} = \langle \mathfrak{D}_{\mathcal{A}}, \Psi_{\mathcal{A}}, \llbracket \cdot \rrbracket, \perp, \top, \wedge, \vee, \neg \rangle$ is an effective Boolean Algebra.

DEFINITION 3.1 (Semantic abstraction). *Let $\mathcal{A}$ be an effective Boolean Algebra and $\rho \in uco(\wp(\mathfrak{D}_{\mathcal{A}}))$ be a partitioning abstraction of its domain of denotations. The semantic abstraction of $\mathcal{A}$ w.r.t. $\rho$, denoted $\langle\!\langle \rho \rangle\!\rangle$-abstraction, is the effective Boolean algebra*

$$\mathcal{A}^\rho = \langle \mathfrak{D}_{\mathcal{A}}^\rho, \Psi_{\mathcal{A}}, \llbracket \cdot \rrbracket^\rho, \rho(\perp), \top, \wedge, \vee^\rho, \neg^\rho \rangle$$

*where:*

$$\mathfrak{D}_{\mathcal{A}}^\rho = \bigcup \left\{ \rho(d) \;\middle|\; d \in \mathfrak{D}_{\mathcal{A}} \right\}$$

$\llbracket \cdot \rrbracket^\rho : \Psi_{\mathcal{A}} \longrightarrow \wp(\mathfrak{D}_{\mathcal{A}}^\rho)$ *such that*
$$\llbracket \varphi \rrbracket^\rho = \rho(\llbracket \varphi \rrbracket) = \bigcup \left\{ \rho(d) \;\middle|\; d \in \llbracket \varphi \rrbracket \right\}$$

$\varphi_1, \varphi_2 \in \Psi_{\mathcal{A}} : \llbracket \varphi_1 \vee^\rho \varphi_2 \rrbracket^\rho = \llbracket \varphi_1 \rrbracket^\rho \cup \llbracket \varphi_2 \rrbracket^\rho$

$\varphi \in \Psi_{\mathcal{A}} : \llbracket \neg^\rho \varphi \rrbracket^\rho = \mathfrak{D}_{\mathcal{A}}^\rho \smallsetminus \llbracket \varphi \rrbracket^\rho$

Before abstracting predicates, i.e., syntax, we have to guarantee the effectiveness of symbolic computation in the SFA. Next lemma proves that if $S$ is a set, whenever $\eta \in uco(\wp(S))$ is additive $\eta$ maps any r.e. subset $X$ of $S$ into a r.e. (abstract) subset $\eta(X)$ of $S$.

LEMMA 3.2. *If $X \subseteq S$ is r.e. and $\eta \in uco(\wp(S))$ is additive, then $\eta(X)$ is r.e., namely $\eta \in uco(\wp^{\text{re}}(S))$.*

By Lemma 3.2, because $\eta$ is a recursive function, and by Kleene's characterization of recursive enumerable sets, the range of $\eta$ over r.e. sets is itself r.e. (see [21]).

THEOREM 3.3. *If $S$ is a set and $\eta \in uco(\wp^{\text{re}}(S))$ is additive then $\left\{ \eta(X) \;\middle|\; X \subseteq S \wedge X \text{ is r.e.} \right\}$ is r.e.*

DEFINITION 3.4 (Syntactic abstraction). *Let $\mathcal{A}$ be an effective Boolean Algebra and let $\eta \in uco(\wp^{\text{re}}(\Psi_{\mathcal{A}}))$ be an additive abstraction of predicates. The syntactic abstraction of $\mathcal{A}$ w.r.t. $\eta$, denoted $\langle \eta \rangle$-abstraction, is the effective Boolean algebra*

$$\mathcal{A}_\eta = \langle \mathfrak{D}_{\mathcal{A}}, \eta(\wp^{\text{re}}(\Psi_{\mathcal{A}})), \llbracket \cdot \rrbracket, \perp, \top, \wedge, \vee, \neg \rangle$$

*where $\llbracket \cdot \rrbracket : \eta(\wp^{\text{re}}(\Psi_{\mathcal{A}})) \longrightarrow \wp(\mathfrak{D}_{\mathcal{A}})$ is defined as in SFA.*

If we have both a $\langle\!\langle \rho \rangle\!\rangle$-abstraction and a $\langle \eta \rangle$-abstraction of an effective Boolean algebra $\mathcal{A}$, then we define the combined abstraction $\langle\!\langle \rho \rangle\!\rangle\langle \eta \rangle$-abstraction of $\mathcal{A}$ by combining them as follows. Let $\rho \in uco(\wp(\mathfrak{D}_{\mathcal{A}}))$ and $\eta \in uco(\wp^{\text{re}}(\Psi_{\mathcal{A}}))$. The abstraction of $\mathcal{A}$ w.r.t. $\rho$ and $\eta$ is the effective Boolean algebra

$$\mathcal{A}_\eta^\rho = \langle \mathfrak{D}_{\mathcal{A}}^\rho, \eta(\wp^{\text{re}}(\Psi_{\mathcal{A}})), \llbracket \cdot \rrbracket^\rho, \rho(\perp), \top, \wedge, \vee^\rho, \neg^\rho \rangle$$

It is clear that $\mathcal{A}_\eta = \mathcal{A}_\eta^{\text{id}}$ and $\mathcal{A}^\rho = \mathcal{A}_{\text{id}}^\rho$. In the following of the paper we assume that $\langle\!\langle \rho \rangle\!\rangle$- and $\langle \eta \rangle$-abstractions satisfy the hypothesis in Definition 3.1 and 3.4 respectively.

THEOREM 3.5. *If $\mathcal{A}$ is decidable then for any $\rho \in uco(\wp(\mathfrak{D}_{\mathcal{A}}))$ and $\eta \in uco(\wp^{\text{re}}(\Psi_{\mathcal{A}}))$, $\mathcal{A}_\eta^\rho$ is decidable.*

Note that, in the definition of symbolic automata there is a strong relation in the underlying effective Boolean algebra $\mathcal{A}$ between the domain of denotations $\mathfrak{D}_{\mathcal{A}}$ and the set of predicates $\Psi_{\mathcal{A}}$ used to symbolically represent them. This means that, if we abstract the domain of denotations by considering $\rho \in uco(\wp(\mathfrak{D}_{\mathcal{A}}))$, leaving unchanged $\Psi_{\mathcal{A}}$ we are implicitly changing the interpretation of predicates in $\mathfrak{D}_{\mathcal{A}}$. On the other hand, if we abstract the predicates by considering $\eta \in uco(\wp^{\text{re}}(\Psi_{\mathcal{A}}))$ we explicitly describe how symbols are abstracted and the semantics is simply the collection of all the semantics denoting the same abstracted predicate. This leads to the following notion of *compatible abstractions*.

### 3.2 Compatible syntactic and semantic abstractions

Let us consider a $\langle\!\langle \rho \rangle\!\rangle$-abstraction of $\mathcal{A}$, we aim at characterizing the syntactic abstractions that produce abstract predicates which may have semantics in $\mathfrak{D}_{\mathcal{A}}^\rho$. This is captured by the notion of $\langle\!\langle \rho \rangle\!\rangle$-compatibility of a syntactic abstraction. Any semantic $\langle\!\langle \rho \rangle\!\rangle$-abstraction naturally induces a corresponding syntactic $\langle \Omega(\rho) \rangle$-abstraction with $\Omega(\rho) \in uco(\wp(\Psi_{\mathcal{A}}))$ defined as follows:

$$\Omega(\rho) \stackrel{\text{def}}{=} \lambda \Phi. \bigcup \left\{ \Phi' \;\middle|\; \llbracket \Phi' \rrbracket \subseteq \llbracket \Phi \rrbracket^\rho \right\}$$

Analogously, any syntactic $\langle \eta \rangle$-abstraction naturally induces a corresponding semantic $\langle\!\langle \mho(\eta) \rangle\!\rangle$-abstraction with $\mho(\eta) \in uco(\wp(\mathfrak{D}_{\mathcal{A}}^\rho))$. In order to characterize when and how a syntactic abstraction induces a semantic abstraction, we need to characterize the syntactic abstraction that precisely corresponds to the semantics $\llbracket \cdot \rrbracket$, namely the abstraction collecting all the predicates having the same semantics $\llbracket \cdot \rrbracket$. This is precisely $\Omega(\text{id})$, which can be rewritten as $\lambda \Phi. \llbracket \llbracket \Phi \rrbracket \rrbracket^+$. Here $\llbracket \cdot \rrbracket^+$ is the *adjoint semantic function* defined as follows:

$$\llbracket \cdot \rrbracket^+ \stackrel{\text{def}}{=} \lambda X \in \wp(\mathfrak{D}_{\mathcal{A}}). \bigcup \left\{ \Phi \;\middle|\; \llbracket \Phi \rrbracket \subseteq X \right\}$$

Then we can define the induced $\langle\!|\mho(\eta)|\!\rangle$-abstraction:

$$\mho(\eta) \overset{\text{def}}{=} \lambda X. \bigcup \left\{ Y \mid \; [\![Y]\!]^+ \subseteq \eta([\![X]\!]^+) \; \right\}$$

Observe that, when $[\![\cdot]\!] : \Psi_{\mathcal{A}} \to \wp(\mathfrak{D}_{\mathcal{A}})$ is surjective, namely when there exists at least one predicate for each possible semantics in $\wp(\mathfrak{D}_{\mathcal{A}})$ we have that $\mho(\mathtt{id}) = \mathtt{id}$. Indeed, $\mathtt{id} \in uco(\wp^{\text{re}}(\Psi_{\mathcal{A}}))$ considers every single predicate and we have a predicate for each semantic object so in this case we have no effects on the semantics and $\mho(\mathtt{id})$ return precisely the identity on the semantics.

Compatibility of a $\langle\eta\rangle$-abstraction w.r.t. $\langle\!|\rho|\!\rangle$-abstraction can therefore be defined in terms of relative abstraction of $\eta$ and $\Omega(\rho)$, or analogously, in terms of relative abstraction of $\rho$ and $\mho(\eta)$.

DEFINITION 3.6 (Semantic compatibility). *Given a $\langle\!|\rho|\!\rangle$-abstracted effective Boolean algebra $\mathcal{A}^\rho$ and a syntactic abstraction $\eta \in uco(\wp^{\text{re}}(\Psi_{\mathcal{A}}))$, $\eta$ is $\langle\!|\rho|\!\rangle$-compatible if:*

$$\eta \sqsubseteq \Omega(\rho) \tag{1}$$

Intuitively we have semantic compatibility when the syntactic abstraction is more concrete than the semantic abstraction, when they are compared on the domain of abstractions of predicates. Indeed, semantic compatibility means that the way a syntactic abstraction $\eta$ partitions the set of predicates of $\mathcal{A}$ is a refinement of the partition induced by the syntactic abstraction $\Omega(\rho)$ that corresponds to the semantic abstraction $\rho$. We can say that when we have semantic compatibility the abstraction of the syntax distinguishes programs with the same abstract semantics, namely the abstract program provides an under-approximation of the abstract program behavior.

THEOREM 3.7. *Let $\rho \in uco(\wp(\mathfrak{D}_{\mathcal{A}}))$, then $\Omega(\rho)$ is the most abstract syntactic abstraction $\langle\!|\rho|\!\rangle$-compatible.*

Note that $\mathcal{A}^{\Omega(\rho)}$ may not be an effective Boolean algebra because $\Omega(\rho)(\wp^{\text{re}}(\Psi_{\mathcal{A}}))$ may not be a r.e. set.

EXAMPLE 3.8. *Consider the domains depicted in Fig. 1 (the missing point labels are the set union of smaller elements). The first three domains on the left represent possible syntactic abstractions of $\wp(\Psi_{\mathcal{A}})$, where*

$$\Psi_{\mathcal{A}} \overset{\text{def}}{=} \{x + y > 3, x \geq 3, y \geq 0, x + y > 3 \wedge x \geq 3 \wedge y \geq 0\}.$$

*The last domain on the right represents possible semantic abstractions of $\wp(\mathfrak{D}_{\mathcal{A}})$, where*

$$\mathfrak{D}_{\mathcal{A}} \overset{\text{def}}{=} \{[\![x + y > 3]\!], [\![x \geq 3]\!], [\![y \geq 0]\!]\}.$$

*Consider for instance the semantic abstraction $\rho$ of $\wp(\mathfrak{D}_{\mathcal{A}})$, depicted with circles on the last domain on the right. The corresponding syntactic abstraction $\Omega(\rho)$ is depicted on the three syntactic domain on the left. Considering the closures depicted on the first domain on the left we observe that the closure $\eta_1 \in \wp(\Psi_{\mathcal{A}})$ is $\langle\!|\rho|\!\rangle$-compatible being more concrete that $\Omega(\rho)$. This means that the syntactic abstraction can distinguish predicates with the same abstract semantics. In particular, while $\rho([\![x + y > 3]\!]) = \rho([\![x \geq 3]\!])$ we have that $\eta_1(x + y > 3) = \{x + y > 3, x \geq 3\}$ while $\eta_1(x \geq 3) = \{x \geq 3\}$.*

Now consider a $\langle\eta\rangle$-compatible abstraction of $\mathcal{A}$. We introduce the notion of $\langle\eta\rangle$-compatibility of a semantic abstraction.

DEFINITION 3.9 (Syntactic compatibility). *A semantic abstraction $\rho \in uco(\wp(\mathfrak{D}_{\mathcal{A}}))$ is $\langle\eta\rangle$-compatible for a syntactic $\langle\eta\rangle$-abstraction $\mathcal{A}_\eta$ if:*

$$\eta \sqsupseteq \Omega(\rho) \tag{2}$$

Intuitively we have syntactic compatibility when the syntactic abstraction is more abstract than the semantic abstraction when they are compared on the domain of abstractions of predicates.
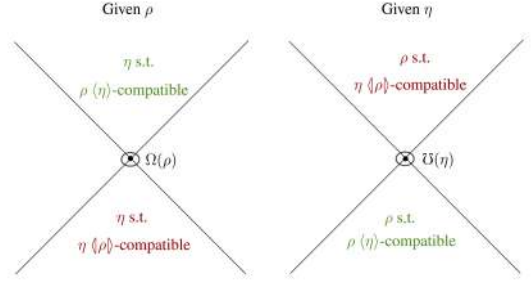


**Figure 2.** Relation between compatibilities.

Indeed, syntactic compatibility means that the semantic abstraction $\rho$ corresponds to a syntactic abstraction $\Omega(\rho)$ and that the partition on the set of predicates of $\mathcal{A}$ induced by $\Omega(\rho)$ is a refinement of the partition induced by $\eta$. In other words, when we have syntactic compatibility the abstraction $\eta$ of the syntax collapses programs with different abstract semantics $\rho$, hence capturing behaviors that, according to $\rho$, are not related with the program to analyze, yet providing an over-approximation of the abstract program behavior.

THEOREM 3.10. *Let $\eta \in uco(\wp^{\text{re}}(\Psi_{\mathcal{A}}))$, then $\mho(\eta)$ is the most concrete semantic abstraction $\langle\eta\rangle$-compatible.*

EXAMPLE 3.11. *Consider again the example in Fig. 1 introduced in Example 3.8. Consider in this case the syntactic abstraction $\eta_3$ depicted on the third domain . We observe that $\rho$ is $\langle\eta_3\rangle$-compatible since $\eta_3$ is more abstract than $\Omega(\rho)$. This means that $\eta_3$ induces a further semantic abstraction collapsing elements with different $\rho$ abstract semantics. In particular, $\rho([\![x + y > 3 \wedge x \geq 3 \wedge x \geq 3]\!]) \neq \rho([\![x \geq 3]\!])$ while $\eta_3(x + y > 3 \wedge x \geq 3 \wedge y \geq 0) = \eta_3(x \geq 3) = \top$. In this example we can also observe a syntactic abstraction $\eta_2$ (depicted on the second domain) which fails both the compatibilities since it not comparable with $\Omega(\rho)$.*

Finally, we show when a syntactic abstraction does induce an abstraction of the semantic denotations and vice versa.

LEMMA 3.12. *Let $\eta \in uco(\wp^{\text{re}}(\Psi_{\mathcal{A}}))$:*

1. $\eta \sqsupseteq \Omega(\mathtt{id})$    *iff*    $\forall \Phi \in \wp^{\text{re}}(\Psi_{\mathcal{A}}). [\![[\![\eta(\Phi)]\!]]\!]^+ = \eta(\Phi)$
2. $\eta \sqsubseteq \Omega(\mathtt{id})$    *iff*    $\forall \Phi \in \wp^{\text{re}}(\Psi_{\mathcal{A}}). \eta([\![[\![\Phi]\!]]\!]^+) = [\![[\![\Phi]\!]]\!]^+$

THEOREM 3.13. *Let $\eta \in uco(\wp^{\text{re}}(\Psi_{\mathcal{A}}))$, then*

$$\eta \sqsubseteq \Omega(\mathtt{id}) \Rightarrow \mho(\eta) = \mathtt{id}$$

This result tells us that when we have a syntactic abstraction distinguishing predicates with the same semantics, then we cannot abstract the semantics.

We prove that we can characterize compatibilities both in the domain of semantic abstractions and in the domain of syntactic abstractions.

THEOREM 3.14. *Let $\eta \in uco(\wp^{\text{re}}(\Psi_{\mathcal{A}}))$ be such that $\eta \sqsupseteq \Omega(\mathtt{id})$, and $\rho \in uco(\wp(\mathfrak{D}_{\mathcal{A}}))$:*

$$\Omega(\rho) \sqsubseteq \eta \qquad iff \qquad \rho \sqsubseteq \mho(\eta)$$

In Fig. 2 we can see the relation between the two compatibilities. In particular we observe that the two transformers, form syntax to semantics and viceversa, show a relation similar to an adjunction, as observed in the following result.

PROPOSITION 3.15. *Let $\eta \in uco(\wp^{\text{re}}(\Psi_{\mathcal{A}}))$ and $\rho \in uco(\wp(\mathfrak{D}_{\mathcal{A}}))$ the following conditions holds:*

$$(1) \; \mho(\Omega(\rho)) \sqsupseteq \rho \qquad\qquad (2) \; \Omega(\mho(\eta)) \sqsubseteq \eta.$$
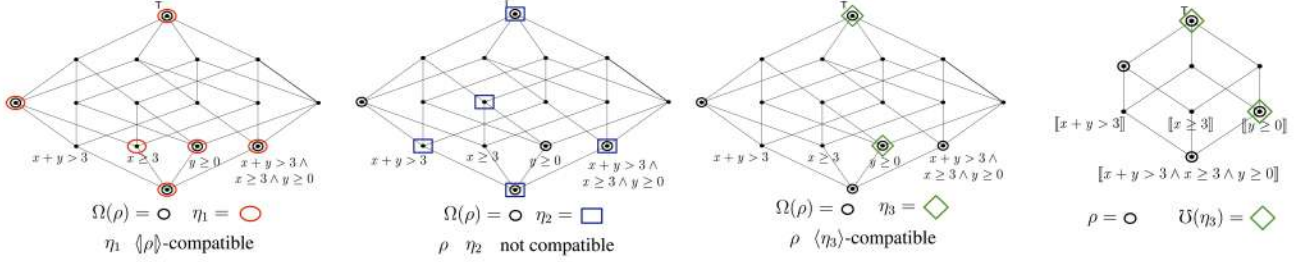
**Figure 1.** Compatible abstractions.

EXAMPLE 3.16. *Consider again the example in Fig. 1. For $\eta_3$ which satisfies the hypotheses of Th. 3.14, we have a corresponding semantic abstraction $\mho(\eta_3)$ (depicted on the right) which is indeed more abstract than $\rho$.*

As a corollary of the previous results we show when a $\langle\!|\rho|\!\rangle\langle\eta\rangle$-abstraction of $\mathcal{A}$ satisfies both the compatibilities. The computational cost of making analyses compatible is still to be explored.

PROPOSITION 3.17. *Let $\eta \in uco(\wp^{\mathrm{re}}(\Psi_{\mathcal{A}}))$ and $\rho \in uco(\wp(\mathfrak{D}_{\mathcal{A}}))$ such that $\Omega(\rho) \in uco(\wp^{\mathrm{re}}(\Psi_{\mathcal{A}}))$, the following facts are equivalent:*

1. *$\eta$ is $\langle\!|\rho|\!\rangle$-compatible and $\rho$ is $\langle\eta\rangle$-compatible;*
2. *$\eta = \Omega(\rho)$;*
3. *$\rho = \mho(\eta)$.*

### 3.3 Abstracting symbolic automata

Consider a SFA $M = \langle \mathcal{A}, Q, q_0, F, \Delta \rangle$ and the $\langle\!|\rho|\!\rangle\langle\eta\rangle$-abstraction of the effective Boolean algebra $\mathcal{A}$, denoted as $\mathcal{A}_\eta^\rho$. We define the symbolic finite automaton corresponding to $M$ on the abstract effective Boolean algebra $\mathcal{A}_\eta^\rho$ as $M_\eta^\rho \stackrel{\text{def}}{=} \langle \mathcal{A}_\eta^\rho, Q, q_0, F, \Delta_\eta \rangle$ where:

$$\Delta_\eta \stackrel{\text{def}}{=} \big\{ (q, \eta(\varphi), q') \,\big|\, (q, \varphi, q') \in \Delta \big\}$$

Note that $M_\eta = M_{\mathrm{id}}^\eta$ and $M^\rho = M_{\mathrm{id}}^\rho$. In the following we prove that when abstracting the underling effective Boolean algebra of an SFA we over-approximate the recognized language, providing a sound approximation in the sense of abstract interpretation.

THEOREM 3.18. *Given a SFA $M = \langle \mathcal{A}, Q, q_0, F, \Delta \rangle$, two closures $\eta \in uco(\wp^{\mathrm{re}}(\Psi_{\mathcal{A}}))$ and $\rho \in uco(\wp(\mathfrak{D}_{\mathcal{A}}))$, the abstract effective Boolean algebra $\mathcal{A}_\eta^\rho$ and the corresponding SFA $M_\eta^\rho = \langle \mathcal{A}_\eta^\rho, Q, q_0, F, \Delta_\eta \rangle$. Then: $\mathscr{L}(M) \subseteq \mathscr{L}(M_\eta^\rho)$.*

For this reason is the following we abuse terminology and refer to the SFA whose underlying Boolean algebra is an $\langle\!|\rho|\!\rangle\langle\eta\rangle$-abstraction of a Boolean algebra $\mathcal{A}$ as an $\langle\!|\rho|\!\rangle\langle\eta\rangle$-abstract SFA. Moreover, we can observe that given two abstract Boolean algebra $\mathcal{A}_{\eta_1}^{\rho_1}$ and $\mathcal{A}_{\eta_2}^{\rho_2}$ and an SFA $M$ on $\mathcal{A}$, then the relation between the languages recognized respectively by $M_{\eta_1}^{\rho_1}$ and by $M_{\eta_2}^{\rho_2}$ corresponds to the relation existing between the best correct approximation of the semantics $[\![\cdot]\!]$ with respect to the pair of abstractions $\rho_1, \eta_1$ and $\rho_2, \eta_2$. This is formally stated in the following Proposition.

PROPOSITION 3.19. *Consider a SFA $M = \langle \mathcal{A}, Q, q_0, F, \Delta \rangle$, the closures $\eta_1, \eta_2 \in uco(\wp^{\mathrm{re}}(\Psi_{\mathcal{A}}))$ and $\rho_1, \rho_2 \in uco(\wp(\mathfrak{D}_{\mathcal{A}}))$, then:*

$$\mathscr{L}(M_{\eta_1}^{\rho_1}) \subseteq \mathscr{L}(M_{\eta_2}^{\rho_2}) \quad \Leftrightarrow \quad \rho_1 \circ [\![\cdot]\!] \circ \eta_1 \sqsubseteq \rho_2 \circ [\![\cdot]\!] \circ \eta_2$$
$$\Leftarrow \quad \rho_1 \sqsubseteq \rho_2 \wedge \eta_1 \sqsubseteq \eta_2$$

## 4. Minterms

A notion which plays a central role in our transformations of SFA is the notion of *minterm*. This notion has been introduced in [9] for

```
1. MINTERMS_A(Φ) ≝
2.     tree := new Tree(⊤_A, null, null);
3.     foreach φ in Φ tree.Refine(φ);
4.     return Leaves(tree);
         //The minterms are the leaf predicates
5. class Tree
6.     Predicate ψ; Tree left; Tree right;
7.     Refine(φ) ≝
8.     if (IsSat_A(ψ ∧ φ) and IsSat_A(ψ ∧ ¬φ))
9.         if (left = null)        // If the tree is a leaf then split ψ
10.            left := new Tree(ψ ∧ φ, null, null);
11.            right := new Tree(ψ ∧ ¬φ, null, null);
12.        else left.Refine(φ); right.Refine(φ);
```

**Figure 3.** Minterm generation algorithm.

providing a minimal and univocal representation of the predicates in a given set of predicates, e.g., the guards of a given program. In this context we observe some peculiar properties of minterms which make them powerful tools for reasoning on semantics in a syntactic way. A minterm is a minimal satisfiable boolean combination of all predicates occurring in a given SFA. Minterms can be generated from a set of predicates by the algorithm proposed in [9] and reported in Fig. 3. As observed in [9] the set of minterms of an SFA may be expensive to compute, indeed in the worst case the complexity of the algorithm that computes the minterms is exponential in the number of guards of the SFA.

### 4.1 Basic properties of Minterms

The minterm generation for a formula $\varphi$ produces a tree $T_\varphi$ that satisfies the following basic properties.

PROPOSITION 4.1. *Let* tree *be the tree built during the minterm generation, starting from a set $\Phi \in \wp^{\mathrm{re}}(\Psi_{\mathcal{A}})$ of predicates. Given $\varphi \in \Psi_{\mathcal{A}}$, let us denote by $T_\varphi$ the subtree of* tree *having $\varphi$ as root. Then the following properties hold:*

1. *Let $Leaves(T_\varphi) = \{\varphi_1, \ldots, \varphi_k\}$, then $\varphi \Leftrightarrow \bigvee_{i \in \{1..k\}} \varphi_i$;*
2. *Any $\varphi \in$ MINTERMS$(\Phi)$ $\varphi$ satisfiable implies that for all $\varphi' \in$ MINTERMS$(\Phi) \smallsetminus \{\varphi\}$ is not satisfiable.*
3. *For all $\varphi_1, \varphi_2 \in \Phi$ we have that $\varphi_1 \wedge \varphi_2$ is satisfiable iff $Leaves(T_{\varphi_1}) \cap Leaves(T_{\varphi_2}) \neq \varnothing$;*
4. *For any $\varphi_1, \varphi_2 \in \Phi$ we have that $\varphi_1 \Rightarrow \varphi_2$ is satisfiable with $\varphi_1$ satisfiable iff $Leaves(T_{\varphi_1})_{\mathrm{SAT}} \subseteq Leaves(T_{\varphi_2})$[1];*

The following proposition shows that the semantics of minterms is a partition of the domain $\mathfrak{D}_{\mathcal{A}}$ of denotations.

---

[1] where $Leaves(T_{\varphi_1})_{\mathrm{SAT}} \stackrel{\text{def}}{=} \big\{ \varphi \in Leaves(T_{\varphi_1}) \,\big|\, \varphi \text{ is satisfiable} \big\}$.

PROPOSITION 4.2. *Let* $\mathcal{A} = \langle \mathfrak{D}_{\mathcal{A}}, \Psi_{\mathcal{A}}, [\![\cdot]\!], \bot, \top, \wedge, \vee, \neg \rangle$ *be an effective Boolean algebra, then* $\{\, [\![\varphi]\!] \mid \varphi \in \mathrm{MINTERMS}(\Psi_{\mathcal{A}}) \,\}$ *is a partition of* $\mathfrak{D}_{\mathcal{A}}$.

## 4.2 Approximated Minterms

Minterms change their structure when the underlying Boolean algebra is approximated by abstract interpretation. We consider an effective Boolean algebra $\mathcal{A} = \langle \mathfrak{D}_{\mathcal{A}}, \Psi_{\mathcal{A}}, [\![\cdot]\!], \bot, \top, \wedge, \vee, \neg \rangle$, where the semantic function $[\![\cdot]\!] : \Psi_{\mathcal{A}} \to \wp(\mathfrak{D}_{\mathcal{A}})$ is surjective. Consider a subset $\Psi \subseteq \Psi_{\mathcal{A}}$ of such predicates, for example the set of predicates that label a given SFA. We define the syntactic abstraction $\eta_{\Psi} \in uco(\wp^{\mathrm{re}}(\Psi_{\mathcal{A}}))$ as that abstraction of predicates that observes precisely only the predicates in $\Psi$ and abstract in $\top$ any other predicate. Let $\varphi \in \Psi_{\mathcal{A}}$, then $\eta_{\Psi}$ is formally defined as additive lift of:

$$\eta_{\Psi}(\{\varphi\}) \stackrel{\mathrm{def}}{=} \begin{cases} \{\varphi\} & \text{if } \varphi \in \Psi \\ \top & \text{otherwise} \end{cases}$$

Note that the fixpoints of $\eta_{\Psi}$ is $\eta_{\Psi}(\wp^{\mathrm{re}}(\Psi_{\mathcal{A}})) = \wp(\Psi) \cup \{\top\}$. Of course $\eta_{\Psi}$ corresponds to an abstraction $\mho(\eta_{\Psi})$ on the semantics that precisely observes only the semantics of the predicates in $\Psi$, as stated by the following result.

LEMMA 4.3. *Let* $\mathcal{A} = \langle \mathfrak{D}_{\mathcal{A}}, \Psi_{\mathcal{A}}, [\![\cdot]\!], \bot, \top, \wedge, \vee, \neg \rangle$ *be an effective Boolean algebra and consider* $\eta_{\Psi} \in uco(\wp^{\mathrm{re}}(\Psi_{\mathcal{A}}))$ *which is* $\langle \mathrm{id} \rangle$*-compatible, then:*

$$\mho(\eta_{\Psi})(\wp(\mathfrak{D}_{\mathcal{A}})) = \{\, [\![\Phi]\!] \mid \Phi \in \wp^{\mathrm{re}}(\Psi_{\mathcal{A}}) \,\}$$

Observe that $\eta_{\Psi}$ is $\langle \mathrm{id} \rangle$-compatible if whenever there is a predicate in $\Psi$ then $\Psi$ contains also all the predicates with the same semantics. The closure $\mho(\eta_{\Psi}) \in uco(\wp(\mathfrak{D}_{\mathcal{A}}))$ may not be partitioning in general, so we consider $\Pi(\mho(\eta_{\Psi}))$ and we observe that the equivalence classes of the partition induced by $\Pi(\mho(\eta_{\Psi}))$ on $\mathfrak{D}_{\mathcal{A}}$ are precisely the semantics of the minterms of $\Psi$.

PROPOSITION 4.4. *Let* $\mathcal{A} = \langle \mathfrak{D}_{\mathcal{A}}, \Psi_{\mathcal{A}}, [\![\cdot]\!], \bot, \top, \wedge, \vee, \neg \rangle$ *be an effective Boolean algebra, and consider* $\Psi \subseteq \Psi_{\mathcal{A}}$ *such that the abstraction* $\eta_{\Psi} \in uco(\wp^{\mathrm{re}}(\Psi_{\mathcal{A}}))$ *is* $\langle \mathrm{id} \rangle$*-compatible, then:*

$$\{\, [\![\varphi]\!] \mid \varphi \in \mathrm{MINTERMS}(\Psi) \,\} = \{\, \Pi(\mho(\eta_{\Psi}))(d) \mid d \in \mathfrak{D}_{\mathcal{A}} \,\}$$

It is now interesting to observe what happens when we consider a generic syntactic abstraction $\eta \in uco(\wp^{\mathrm{re}}(\Psi_{\mathcal{A}}))$ such that $\eta_{\Psi} \sqsubseteq \eta$, namely that further abstracts the set of predicates $\Psi$ that we are considering. In this case, the semantics of the minterms of the approximated predicates $\eta(\Psi)$ are precisely given by the abstraction $\Pi(\mho(\eta))$ of the semantics of the minterms of $\Psi$.

THEOREM 4.5. *Let* $\mathcal{A} = \langle \mathfrak{D}_{\mathcal{A}}, \Psi_{\mathcal{A}}, [\![\cdot]\!], \bot, \top, \wedge, \vee, \neg \rangle$ *be an effective Boolean algebra, and consider* $\Psi \subseteq \Psi_{\mathcal{A}}$ *such that the abstraction* $\eta_{\Psi} \in uco(\wp^{\mathrm{re}}(\Psi_{\mathcal{A}}))$ *is* $\langle \mathrm{id} \rangle$*-compatible, and an abstraction* $\eta \in uco(\wp^{\mathrm{re}}(\Psi_{\mathcal{A}}))$ *such that* $\eta_{\Psi} \sqsubseteq \eta$. *Then:*

$$\{\, [\![\varphi]\!] \mid \varphi \in \mathrm{MINTERMS}(\eta(\Psi)) \,\} = \\ \{\, \Pi(\mho(\eta))([\![\varphi]\!]) \mid \varphi \in \mathrm{MINTERMS}(\Psi) \,\}$$

This means that the semantics of the minterms of a set of abstract predicates is precisely the abstraction of the semantics of the original predicates.

EXAMPLE 4.6. *Let* $\mathcal{A} = \langle \mathfrak{D}_{\mathcal{A}}, \Psi_{\mathcal{A}}, [\![\cdot]\!], \bot, \top, \wedge, \vee, \neg \rangle$ *be an effective Boolean algebra where* $\Psi_{\mathcal{A}} = \{\, x \in N \mid N \subseteq \mathbb{Z} \,\}$, *and the semantic function* $[\![\cdot]\!] : \Psi_{\mathcal{A}} \to \wp(\mathbb{Z})$ *is naturally defined as* $[\![x \in N]\!] = N$.
*Let us consider the following subset of* $\Psi_{\mathcal{A}}$:

$$\Psi = \{\, x \in \{4, 6\}, x \in \{5, 6\}, x \in \{-5\}, x \in \{-8\} \,\}$$

*the corresponding set of minterms is* $\mathrm{MINTERMS}(\Psi)$:

$$\begin{aligned} &\{\, (x \in \{4, 6\} \wedge x \in \{5, 6\}), \\ &(x \in \{4, 6\} \wedge \neg x \in \{5, 6\}), \\ &(\neg x \in \{4, 6\} \wedge x \in \{5, 6\}), \\ &(\neg x \in \{4, 6\} \wedge \neg x \in \{5, 6\} \wedge x \in \{-5\}), \\ &(\neg x \in \{4, 6\} \wedge \neg x \in \{5, 6\} \wedge \neg x \in \{-5\} \wedge x \in \{-8\}), \\ &(\neg x \in \{4, 6\} \wedge \neg x \in \{5, 6\} \wedge \neg x \in \{-5\} \wedge \neg x \in \{-8\})\} \end{aligned}$$

*observe that:*

$$\begin{aligned} &\{\, [\![\varphi]\!] \mid \varphi \in \mathrm{MINTERMS}(\Psi) \,\} = \\ &\quad \{\, \{4\}, \{6\}, \{5\}, \{-5\}, \{-8\}, \mathbb{Z} \smallsetminus \{4, 6, 5, -5, -8\} \,\} \end{aligned}$$

*the closure* $\eta_{\Psi} \in uco(\wp^{\mathrm{re}}(\Psi_{\mathcal{A}}))$ *is defined as the additive lift of:*

$$\eta_{\Psi}(\{x \in N\}) \stackrel{\mathrm{def}}{=} \begin{cases} \{x \in N\} & \text{if } \{x \in N\} \in \Psi \\ \top & \text{otherwise} \end{cases}$$

*and, as states in Proposition 4.4 we have that:*

$$\begin{aligned} &\{\, \Pi(\mho(\eta_{\Psi}))(d) \mid d \in \mathbb{Z} \,\} = \\ &\quad \{\, \{4\}, \{6\}, \{5\}, \{-5\}, \{-8\}, \mathbb{Z} \smallsetminus \{4, 6, 5, -5, -8\} \,\} \end{aligned}$$

*Let* $\mathbb{Z}^+ \stackrel{\mathrm{def}}{=} \{\, v \mid v \geq 0 \,\}$ *and* $\mathbb{Z}^- \stackrel{\mathrm{def}}{=} \{\, v \mid v < 0 \,\}$ *and let the closure* $\eta_{Sign} \in uco(\wp^{\mathrm{re}}(\Psi_{\mathcal{A}}))$ *defined as the additive lift of:*

$$\eta_{Sign}(\{x \in N\}) \stackrel{\mathrm{def}}{=} \begin{cases} \{x \in \mathbb{Z}^+\} & \text{if } N \subseteq \mathbb{Z}^+ \\ \{x \in \mathbb{Z}^-\} & \text{if } N \subseteq \mathbb{Z}^- \\ \top & \text{otherwise} \end{cases}$$

*Observe that the* $\mathrm{MINTERMS}(\eta_{Sign}(\Psi))$ *is the set*

$$\{\, \{x \in \mathbb{Z}^+\}, \{x \in \mathbb{Z}^-\} \,\}$$

*and the semantics of the minterms of* $\eta_{Sign}(\Psi)$ *is:*

$$\{\, [\![\varphi]\!] \mid \varphi \in \mathrm{MINTERMS}(\eta_{Sign}(\Psi)) \,\} \quad = \quad \{\mathbb{Z}^+, \mathbb{Z}^-\}$$

*Moreover, as shown in Theorem 4.5:*

$$\{\, \Pi(\mho(\eta_{Sign}))([\![\varphi]\!]) \mid \varphi \in \mathrm{MINTERMS}(\Psi) \,\} = \{\mathbb{Z}^+, \mathbb{Z}^-\}$$

## 5. Topological SFA abstraction

In Section 3 we have seen how an SFA can be abstracted by abstracting its underlying Boolean algebra. This abstraction does not influence directly the topological structure of SFA. When dealing with automata, the natural way of thinking about automata simplification (or abstraction) is the merge of states. In general, we can define a simplification operation on automata that collapses states wrt a given equivalence relation over states. Namely, the equivalence relation establish the criteria that the simplification uses for merging states.

DEFINITION 5.1. *Consider a SFA* $M = \langle \mathcal{A}, Q, q_0, F, \Delta \rangle$ *and an equivalence relation* $R \subseteq Q \times Q$ *over its states. We denote with* $Sim_R(M)$ *the SFA obtained by simplifying* $M$ *wrt* $R$, *namely the SFA computed as the quotient of* $M$ *wrt* $R$, *i.e.,* $Sim_R(M) = M_{/R}$.

Thus, SFA simplification is the operation of quotient made parametric on the equivalence relation used to merge states. It is easy to observe that for every equivalence relation $R$, the SFA $Sim_R(M)$ resulting from SFA simplification recognizes at least the language recognized by $M$. Indeed when we merge states we keep all the transitions of the original SFA and we may add some new *spurious* ones.

PROPOSITION 5.2. *Consider a SFA* $M = \langle \mathcal{A}, Q, q_0, F, \Delta \rangle$. *For any equivalence relation* $R \subseteq Q \times Q$ *we have that* $\mathscr{L}(M) \subseteq \mathscr{L}(Sim_R(M))$.

Given two equivalence relations $R$ and $R'$, we write $R \preceq R'$ when $R$ is a refinement of $R'$. Of course the coarser is the equivalence relation the wider is the language recognized by the corresponding simplified SFA.

**PROPOSITION 5.3.** *Consider a SFA $M = \langle \mathcal{A}, Q, q_0, F, \Delta \rangle$ and two equivalence relations $R, R' \subseteq Q \times Q$ such that $R \preceq R'$. Then $\mathscr{L}(Sim_R(M)) \subseteq \mathscr{L}(Sim_{R'}(M))$.*

Another important property of topological abstractions is that they do not change the set of minterms, since they do not change the predicates. In the following we report a simplification algorithm where the predicates of the SFA to simplify are first rewritten as disjunction of minterms (line 3-7). Thus, whenever the equivalence relation $R$ deals with properties of the languages of strings that reaches or starts from a state, it may be easier to check these properties on minterms instead of checking them on the language of denotations. (Examples will be provided in the following).

---

**Simplify$(M, R)$**

1.   Input: $M = \langle \mathcal{A}, Q, q_0, F, \Delta \rangle, R \subseteq Q \times Q,$
2.   $\qquad \mathcal{A} = \langle \mathfrak{D}_{\mathcal{A}}, \Psi_{\mathcal{A}}, [\![\cdot]\!], \bot, \top, \wedge, \vee, \neg \rangle$
3.   $\mathcal{M}t(M) \stackrel{\text{def}}{=} \text{MINTERMS}\left( \left\{ \psi \,\middle|\, \begin{array}{l} \exists p, q \in Q. \\ p \stackrel{\psi}{\longrightarrow} q \in \Delta \end{array} \right\} \right)$
4.   $M' = \langle \mathcal{A}', Q, q_0, F, \Delta' \rangle$:
5.   $\qquad \mathcal{A}' \stackrel{\text{def}}{=} \langle \mathfrak{D}_{\mathcal{A}}, \mathcal{M}t(M), [\![\cdot]\!], \bot, \top, \wedge, \vee, \neg \rangle$
6.   $\qquad \mu(\psi) \stackrel{\text{def}}{=} \bigvee \{ \varphi \in \mathcal{M}t(M) \mid \varphi \in \text{Leaves}(T_\psi) \}$
7.   $\qquad \Delta' \stackrel{\text{def}}{=} \left\{ p \stackrel{\mu(\psi)}{\longrightarrow} q \,\middle|\, \exists \psi. \, p \stackrel{\psi}{\longrightarrow} q \right\}$
8.   Output: $M'' = M'_{/R}$

---

## 5.1 Examples of SFA Simplifications

***Minimization.*** D'Antoni and Veanes in [9] have extended the standard algorithm of Hopcroft for finite state automata minimization to SFA. This operation is based on the idea of refining an initial partition by checking all the possible moves depending on the considered alphabet symbol. In FSA this is feasible because they have a finite alphabet. In SFA the alphabet is r.e., hence in general infinite. For this reason the algorithm proposed iterates this check on predicates/symbols in a way that makes the number of possible iteration finite: instead of checking transitions for each alphabet symbol, the check is made for each minterm (see [9] for details).

Observe that this SFA minimization algorithm can be seen as a simplification wrt. the equivalence relation that relates all and only the states that are reached exactly by the same language of minterms. Consider a SFA $M = \langle \mathcal{A}, Q, q_0, F, \Delta \rangle$, and for every $q \stackrel{\psi}{\longrightarrow} p \in \Delta$ let $\mu(\psi)$ be the predicate $\psi$ written as a disjunction of minterms (namely as the disjunction of the leaves of the subtree $T_\psi$ with root $\psi$ of the tree generated during the construction of the minterms of the considered SFA). We define the language of strings of minterms that reaches a state $q$ as:

$$\dot{\mathscr{L}}(q) \stackrel{\text{def}}{=} \left\{ \mu(\psi_1) \dots \mu(\psi_{n-1}) \,\middle|\, \begin{array}{l} \exists n \in \mathbb{N} : \exists q_1 \dots q_n \in Q : \\ \forall i \in [1, n[. q_i \stackrel{\mu(\psi_i)}{\longrightarrow} q_{i+1} \in \Delta \\ q_n = q \end{array} \right\}$$

Let $\stackrel{.}{\equiv} \subseteq Q \times Q$ be such that $q \stackrel{.}{\equiv} p$ iff $\dot{\mathscr{L}}(q) = \dot{\mathscr{L}}(p)$. Observe that for the properties of minterms proved in the previous section, we have that checking the language of minterms or checking the language of denotations is equivalent, since minterms provide a minimal and unequivocal representation of predicates. Let $Min(M)$ denote the minimization of $M$.

**PROPOSITION 5.4.** $Min(M) = Sim_{\stackrel{.}{\equiv}}(M)$.

***k-Minimization.*** According to the above formalization of SFA minimization, we can weaken minimization by defining a relation over states that observes the language of stings of minterms that reaches a given state. To this end, given an SFA $M = \langle \mathcal{A}, Q, q_0, F, \Delta \rangle$, and for every $q \stackrel{\psi}{\longrightarrow} p \in \Delta$ let $\mu(\psi)$ be the predicate $\psi$ written as a disjunction of minterms, we define the language $\dot{\mathscr{L}}_k(q)$, which is the language of strings of length $k$ that can reach the state $q$:

$$\dot{\mathscr{L}}_k(q) \stackrel{\text{def}}{=} \left\{ \mu(\psi_1) \dots \mu(\psi_{k-1}) \,\middle|\, \begin{array}{l} \exists q_1 \dots q_k \in Q : \\ \forall i \in [1, .k[. q_i \stackrel{\mu(\psi_i)}{\longrightarrow} q_{i+1} \in \Delta \\ q_k = q \end{array} \right\}$$

Let $\stackrel{.}{\equiv}_k \subseteq Q \times Q$ be such that $q \stackrel{.}{\equiv}_k p$ iff $\dot{\mathscr{L}}_k(q) = \dot{\mathscr{L}}_k(p)$. Let $Min_k(M)$ denote the simplification of $M$ wrt $\stackrel{.}{\equiv}_k$. The following examples illustrate the difference between minimization and $k$-minimization.

**EXAMPLE 5.5.** *Consider the SFA $M$ in Fig. 4 on the left. It is clear that the predicates $x$ odd and $(x + 1)$ even are equivalent, as well as predicates $y$ even and $(y + 1)$ odd. This is captured by the minimization algorithm of D'Antoni and Veanes that correctly collapses state $q_4$ with $q_5$ and $q_7$ with $q_8$. The minimized algorithm $Min(M)$ is shown in Fig. 4 at the top on the right. Observe that the edge between $q_2$ and $\{q_4, q_5\}$, as well as the edge between $\{q_7, q_8\}$ and $q_9$, is labeled by one of the two equivalent predicates. Of course, the SFA $M$ and $Min(M)$ recognize the same language. In order to clarify the difference between minimization and $k$-minimization at the bottom right of Fig. 4 we report the result obtained by applying the simplification algorithm wrt $\stackrel{.}{\equiv}_k$ where $k = 1$ at the SFA $M$. Observe that the simplification algorithm with $k = 1$ merges the state $q_6$ with the states $q_7$ and $q_8$, as shown in the resulting SFA $Min_1(M)$. Indeed, the states $q_6$, $q_7$ and $q_8$ are reached by the same language of strings of length 1 (in this simple case all the denotations with $y$ positive). The edge between $\{q_6, q_7, q_8\}$ and $q_9$ is labeled by* true *since it corresponds to $y$ odd $\vee y$ even. We can observe that the language recognized by $Min_1(M)$ is greater than the one recognized by $M$. Let us consider the pairs $(n_1, n_2)$ with $n_1, n_2 \in \mathbb{Z}$ where the first number denotes the values of $x$ and the second the values of $y$. For example we have that the string of pairs $(1, 2)(2, 4)(4, 8)(8, 16) \in \mathscr{L}(Min_1(M))$ while it does not belong to $\mathscr{L}(M) = \mathscr{L}(Min(M))$.*

Of course when the value of $k$ increases it increases also the precision of the simplification wrt $\stackrel{.}{\equiv}_k$ by collapsing states that are equivalent, namely at the limit with $k$ increasing the $k$-minimization becomes minimization.

**THEOREM 5.6.** *Given two states $p$ and $q$ we have that $p \stackrel{.}{\equiv} q$ iff $\forall k \in \mathbb{N}. p \stackrel{.}{\equiv}_k q$.*

**EXAMPLE 5.7.** *Observe that if we compute the simplification of the SFA $M$ in the example in Fig. 4 wrt $\stackrel{.}{\equiv}_k$ and $k = 2$ we obtain the minimized SFA, namely $Min(M) = Min_2(M)$. Indeed, if we consider the language of words of length 2 that reach a given state we can no longer merge $q_6$ with $q_7$ and $q_8$.*

***k-Invariant.*** Minterms provide a systematic simplification of SFA based on the extraction of invariant properties that hold for the language of strings that reach (or start) from a given state. Consider an SFA $M = \langle \mathcal{A}, Q, q_0, F, \Delta \rangle$, and for every $q \stackrel{\psi}{\longrightarrow} p \in \Delta$ let $\mu(\psi)$ be the predicate $\psi$ written as a disjunction of minterms. Consider a state $q \in Q$. For every string $\mu(\psi_1) \dots \mu(\psi_k) \in \dot{\mathscr{L}}_k(q)$ of length $k$ that reaches the state $q$ we have that $IsSat(\bigwedge_{i \in [1, k]} \mu(\psi_i))$ is true iff all the disjunctions $\mu(\psi_i)$ of minterms share at least one minterm. This because, thanks to the properties of minterms, only

**Figure 4.** Minimization and $k$-Minimization

one minterm at the time can be true. Let

$$Inv\{\mu(\psi_i)\}_{i\in[1,k]} \overset{def}{=} \\ \{\ \varphi \in \text{MinTerms}\ |\ \forall i \in [1,k].\ IsSat(\varphi \wedge \mu(\psi_i))\ \}$$

It is the set of *all* the minterms shared by *all* the $\mu(\psi_i)$, which provides the invariant property of the corresponding string. Indeed, thanks to minterms this satisfiability can be checked syntactically. We can therefore define the following equivalence relation $\overset{inv}{\equiv}_k \subseteq Q \times Q$ such that $q \overset{inv}{\equiv}_k p$ iff $\Im_k(q) = \Im_k(p)$ where

$$\Im_k(q) = \{\ Inv\{\mu(\psi_i)\}_{i\in[1,k]}\ |\ \mu(\psi_1)\dots\mu(\psi_k) \in \dot{\mathscr{L}}_k(q)\ \}$$

Thus, $\overset{inv}{\equiv}_k$ collapses states reached by paths that have the same $k$-invariant property. We can observe that, if two states share the same $k$-language then they surely share the same $k$-invariant, while the opposite may not be true since the language fixes an order in the constraints that the commutativity of the conjunction relaxes.

THEOREM 5.8. *Given two states $p$ and $q$, and $k \in \mathbb{N}$, we have that $p\dot{\equiv}_k q$ implies $\forall k' \le k.\ p \overset{inv}{\equiv}_{k'} q$.*

EXAMPLE 5.9. *Consider again the automaton $M$ in Fig. 4. The minterms generated by its predicates are given in the table in Fig. 5: each $i$ denotes the minterm $M_i$ obtained as the conjunction between the constraint on $x$ and on $y$, for instance 3 stays for the minterm $M_3 = (x\ \text{even}\ \wedge x \ge 0 \wedge y\ \text{even}\ \wedge y < 0)$. In Fig. 5 we rewrite $M$ where on each edge the predicates are denoted as the set of the minterms specifying it. For instance $x \ge 0 \equiv \bigvee_{i\in[1,8]} M_i$. Note that on this automaton the $k$-invariant generates the same transformation as the $k$-minimization as showed in Fig. 4. Consider instead the automaton $M_1$ on the right. In this case, the languages recognized by $q_7$ and $q_8$ are different, for instance the trace $(-1,3)(3,-4)(3,5) \in \dot{\mathscr{L}}_3(q_7)$ is not in $\dot{\mathscr{L}}_3(q_8)$ since $(3,-4)$ does not satisfy the predicate between $q_2$ and $q_5$ in $M_1$, i.e., $y \ge 0$. If we consider 3-invariant then we observe that the invariant on the path $q_0q_2q_4q_7$ is $M_{13} \wedge M_{14}$ and the same is for the path $q_0q_2q_5q_8$, hence we can collapse the states $q_7$ and $q_8$.*

### 5.2 Topological abstraction of abstract SFA

It is worth noting that abstraction in SFA may influence the automata simplification. In this section we prove that the efficacy of simplification, and in particular of minimization and $k$-minimization, in SFA is strictly related with the degree of abstraction of their semantics or syntax.

EXAMPLE 5.10. *Consider the SFA $M$ in Fig. 4 and assume that we want to abstract from the parity of $y$. Hence we define abstraction $\eta_1$ on the predicates of $M$ as $\eta_1(y\ \text{odd}) = \eta_1(y\ \text{even}) = \eta_1((y+1)\ \text{odd}) = \text{true}$ and as the identity on the other predicates. In this example we do not abstract the semantics and we consider $\rho = \text{id}$. Let $M_{\eta_1}$ be the SFA wrt the considered abstraction (where the predicates of $M$ are substituted with their abstraction*

*according to $\eta_1$). By applying minimization to this SFA we obtain the SFA $Min(M_{\eta_1})$ depicted at the top left of Fig. 6. We can observe that, due to the predicate abstraction $\eta_1$, the minimization of $M_{\eta_1}$ collapses more states than the minimization of $M$ and therefore: $\mathscr{L}(Min(M)) \subseteq \mathscr{L}(Min(M_{\eta_1}))$. For example the string of pairs $(1,2)(2,4)(4,8)(8,16) \in \mathscr{L}(Min(M_{\eta_1}))$ while it does not belong to $\mathscr{L}(Min(M))$.*

*We have an analogous situation in the case of $k$-minimization. Consider the predicate abstraction $\eta_2$ such that $\eta_1(x\ \text{odd}) = \eta_1(x\ \text{even}) = \eta_1((x+1)\ \text{even}) = \text{true}$ and as the identity on the other predicates, and let $\rho = \text{id}$. By applying the simplification algorithm wrt. $\dot{\equiv}_k$ with $k = 1$ to the SFA $M_{\eta_2}$ we obtain the SFA at the top right of Fig. 6. Also in this case, due to the abstraction $\eta_2$ the simplification algorithm collapses more states and therefore: $\mathscr{L}(Min_1(M)) \subseteq \mathscr{L}(Min_1(M_{\eta_2}))$. For example the string of pairs $(1,2)(3,6)(5,10)(7,14) \in \mathscr{L}(Min_1(M_{\eta_2}))$ while it does not belong to $\mathscr{L}(Min_1(M))$.*

Let $\mathbb{S}$ denote the set of SFA and let us define the following ordering relation $\dot{\le}$ on $\mathbb{S}$ modeling precisely the relative precision of SFA with respect to language containment and size of the automaton, where given $M_1 = \langle \mathcal{A}, Q_1, q_0^1, F_1, \Delta_1 \rangle$ and $M_2 = \langle \mathcal{A}, Q_2, q_0^2, F_2, \Delta_2 \rangle \in \mathbb{S}$ we have that:

$$M_1 \dot{\le} M_2 \Leftrightarrow \mathscr{L}(M_1) \subseteq \mathscr{L}(M_2) \vee \\ \mathscr{L}(M_1) = \mathscr{L}(M_2) \wedge |Q_2| \le |Q_1|$$

It is immediate to observe that $\langle \mathbb{S}, \dot{\le} \rangle$ is a possibly non-complete lattice. Given the SFA simplification $Sim_R : \mathbb{S} \to \mathbb{S}$, a SFA $M = \langle \mathcal{A}, Q, q_0, F, \Delta \rangle$ and a $\langle\!\langle \rho \rangle\!\rangle\langle\!\langle \eta \rangle\!\rangle$-abstraction of the effective Boolean algebra $\mathcal{A}$ we wonder when the diagram in Fig. 7 commutes. In general we have that when we simplify the SFA after the abstraction of the underlying algebra we obtain an SFA that is more abstract than the one obtained by applying simplification before the abstraction. The intuition beyond this is that the abstraction of the underlying Boolean algebra could make equivalent edges of the original SFA that are not equivalent and this may cause the merge of states that would not be merged when simplifying original SFA.

PROPOSITION 5.11. *Given $M = \langle \mathcal{A}, Q, q_0, F, \Delta \rangle \in \mathbb{S}$, the closures $\eta \in uco(\wp(\Psi_{\mathcal{A}}))$ and $\rho \in uco(\wp(\mathfrak{D}_{\mathcal{A}}))$ and a relation $R$, we have that: $Sim_R(M)_\eta^\rho \dot{\le} Sim_R(M_\eta^\rho)$.*

EXAMPLE 5.12. *At the bottom left of Fig. 6 we show the result of abstracting the Boolean algebra after the SFA minimization. We observe that even if $Min(M_{\eta_1})$ and $Min(M)_{\eta_1}$ recognize the same language the automata obtained by minimizing after the abstraction of the underlying Boolean algebra has less states than the one computed by abstracting the Boolean algebra after the minimization. We have a similar result for $k$-minimization as we*

**Figure 5.** $k$-Invariant transformation



$$\eta_1(y\,\mathtt{even}) = \eta_1(y,\mathtt{odd}) = \eta_1((y+1)\,\mathtt{odd}) = \mathtt{true}$$
$$\eta_2(x\,\mathtt{even}) = \eta_2(x\,\mathtt{odd}) = \eta_2((x+1)\,\mathtt{even}) = \mathtt{true}$$

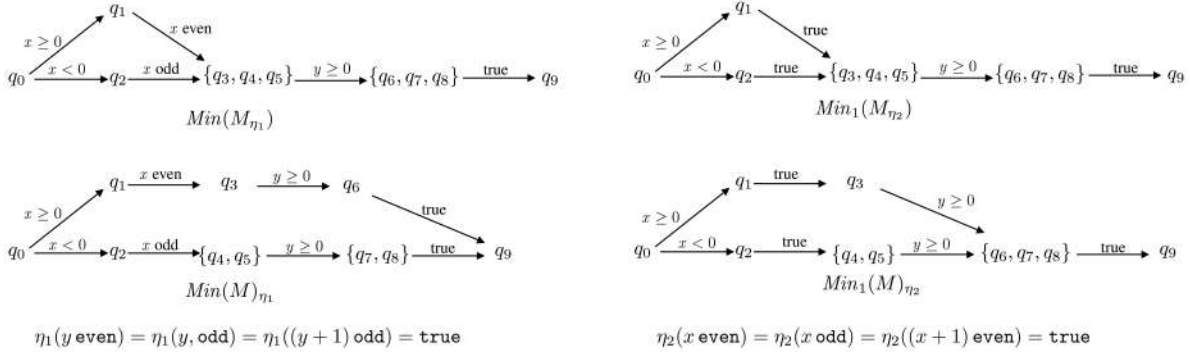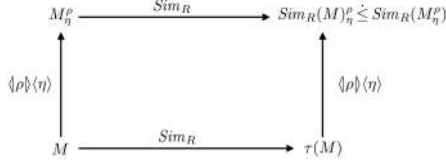**Figure 6.** Minimization and $k$-Minimization in presence of Abstraction



**Figure 7.** Completeness of SFA simplification

*can see by comparing the SFA at the bottom right and top right of Fig. 6.*

## 6. Programs as SFA

In this section we specify the approximate semantics of a program as the language recognized by a SFA. We consider programs in imperative computational model and assume to have access their correct control flow graph (CFG). The CFG of a program is a graph where nodes are given by sequences of non branching instructions. More formally, let $\mathbf{I}$ be the instruction set containing both branching and non-branching instructions. We denote with $\mathbb{I} \subseteq \mathbf{I}$ the set of non-branching instructions and with $\mathbb{C}$ the set of boolean expressions over program states that are guards of the branching instructions. Let $\mathtt{c}$ range over $\mathbb{C}$ and $\mathtt{b}$ range over $\mathbb{I}^*$. The CFG of a program $P \in \mathbf{I}^*$ is a graph $G_P = (N_P, E_P)$ where the set $N_P \subseteq \mathbb{I}^*$ of nodes specifies the basic blocks of $P$, namely the maximal sequences of sequential instructions of $P$, while the set of edges $E_P \subseteq N_P \times \mathbb{C} \times N_P$ denotes the guarded transitions of $P$. In particular, a labeled edge $(\mathtt{b}, \mathtt{c}, \mathtt{b}') \in E_P$ means that the execution of $P$ flows from $\mathtt{b}$ to $\mathtt{b}'$ when the execution of $\mathtt{b}$ leads to a program state that satisfies condition $\mathtt{c}$. When a basic block $\mathtt{b}$ has no outgoing edges in $E_P$ we say that it is final, denoted $\mathtt{b} \in Final[G_P]$. We denote with $in[\mathtt{b}]$ and $out[\mathtt{b}]$ respectively the entry and exit point of the basic block $\mathtt{b}$, and with $\mathbb{PP}[G_P]$ the block delimiters of $G_P$, namely the set of all the entry and exit points of the basic blocks of $G_P$, namely:

$$\mathbb{PP}[G_P] \stackrel{\text{def}}{=} \{\, in[\mathtt{b}] \,\big|\, \mathtt{b} \in N_P \,\} \cup \{\, out[\mathtt{b}] \,\big|\, \mathtt{b} \in N_P \,\}$$

Let $\Sigma$, ranged over by $s$, be the set of possible program states. Let $\mathtt{exec} : \mathbb{I}^* \longrightarrow \wp^{\mathtt{re}}(\Sigma \times \Sigma)$ be the function that defines the semantics of basic blocks, namely the pairs of input/output states that model the execution of sequences of instructions. When $(s, s') \in \mathtt{exec}(\mathtt{b})$ it means that the execution of the sequence of instructions $\mathtt{b}$ transforms state $s$ into state $s'$. Let us denote with $s \models \mathtt{c}$ the fact that the boolean condition $\mathtt{c}$ is satisfied by state $s \in \Sigma$.

We define the set of executions of the CFG of a program $P$ the sequences of basic blocks and guards that can be encountered along a path of $G_P = (N_P, E_P)$. Formally:

$$Exe[G_P] \stackrel{\text{def}}{=} \left\{\, \mathtt{b}_0\mathtt{c}_1\mathtt{b}_1\mathtt{c}_2 \ldots \mathtt{c}_k\mathtt{b}_k \,\middle|\, \begin{array}{l} \forall 0 \leq i < k : \\ (\mathtt{b}_i, \mathtt{c}_{i+1}, \mathtt{b}_{i+1}) \in E_P \end{array} \right\} \quad (3)$$

We consider a safety semantics, namely the semantics of all prefixes of execution traces of a given program $P$ [19]. The execution trace semantics of a program $P$, denoted $[\![P]\!]$, is therefore the set of all finite executions starting from the entry point of the starting basic block $\mathtt{b}_0$ in the CFG $G_P$ of $P$. Let $Init_P \subseteq \Sigma$ be the set of possible initial states of program $P$. Formally, for each $s_0 \in Init_P$:

$$[\![P]\!](s_0) \stackrel{\text{def}}{=} \{(s_0, s_1)(s_1, s_1)(s_1, s_2) \ldots (s_k, s_k)(s_k, s_{k+1}) \mid$$
$$\mathtt{b}_0\mathtt{c}_1\mathtt{b}_1 \ldots \mathtt{c}_k\mathtt{b}_k \in Exe[G_P],$$
$$\forall 0 < i \leq k : s_i \models \mathtt{c}_i, (s_{i-1}, s_i) \in \mathtt{exec}(\mathtt{b}_{i-1})\}$$

$$[\![P]\!] \stackrel{\text{def}}{=} \bigcup \{\, [\![P]\!](s_0) \,\big|\, s_0 \in Init_P \,\}$$

In order to define the SFA that corresponds to the CFG semantics of a given program we need to define an effective Boolean algebra that it is suitable for the representation of program execution. For this reason we define the following effective Boolean algebra where predicates are either basic blocks of instructions or guards of branching instructions, representing the syntactic structure of the program, and the denotations are pairs of input/output states:

$$\mathfrak{P} \stackrel{\text{def}}{=} \langle \Sigma \times \Sigma, \mathbb{I}^* \cup \mathbb{C}, \{\!|\cdot|\!\}, \bot, \top, \wedge, \vee, \neg \rangle$$

where the semantic function $\{\!|\cdot|\!\} : \mathbb{I}^* \cup \mathbb{C} \longrightarrow \wp^{\text{re}}(\Sigma \times \Sigma)$ is defined as follows for $\varphi \in \mathbb{I}^* \cup \mathbb{C}$:

$$\{\!|\varphi|\!\} \stackrel{\text{def}}{=} \left\{ \begin{array}{ll} \{\, (s,s') \mid (s,s') \in \texttt{exec}(\mathtt{b}) \,\} & \text{if } \varphi = \mathtt{b} \in \mathbb{I}^* \\ \{\, (s,s) \mid s \models \mathtt{c} \,\} & \text{if } \varphi = \mathtt{c} \in \mathbb{C} \end{array} \right.$$

we denote with $\{\!|\cdot|\!\}$ also its point-wise extension to $\wp^{\text{re}}(\mathbb{I}^* \cup \mathbb{C})$.

DEFINITION 6.1. *Let $P$ be a program with CFG $G_P$. The SFA associated with $P$ is*

$$M(P) \stackrel{\text{def}}{=} \langle \mathfrak{P}, \mathbb{PP}[G_P], in[\mathtt{b}_0], \{\, out[\mathtt{b}] \mid \mathtt{b} \in Final[G_P]\}, \Delta_P \rangle$$

*where $\mathtt{b}_0$ is the starting basic block of $G_P$ and $\Delta_P$ is defined as:*

$$\Delta_P \stackrel{\text{def}}{=} \{\, (in[\mathtt{b}], \mathtt{b}, out[\mathtt{b}]) \mid \mathtt{b} \in N_P \,\} \cup$$
$$\{\, (out[\mathtt{b}], \mathtt{c}, in[\mathtt{b}']) \mid (\mathtt{b}, \mathtt{c}, \mathtt{b}') \in E_P \,\}$$

PROPOSITION 6.2. *If $P$ is a program then $M(P)$ is a deterministic SFA. $M(P)$ is clean if no dead-block is included in $G_P$.*

The language $\mathscr{L}(M(P)) \in \wp^{\text{re}}((\Sigma \times \Sigma)^*)$ recognized by the SFA $M(P)$ approximates the concrete program semantics $[\![P]\!]$ in a language of sequences of infinitely many possible input/output relations associated with each basic block. This is formally stated by the following theorem.

THEOREM 6.3. *If $P$ is a program then for any $s_0 \in Init_P$: $[\![P]\!](s_0) \in \mathscr{L}(M(P))$.*

Given the SFA $M(P)$ that represents the CFG of a program $P$ then it is possible to approximate the semantics of $P$ by abstracting either the predicates, namely the syntax, or the semantics of the effective Boolean algebra underlying $M(P)$.

Let us consider the minimization simplifications. Given compatible abstractions $\rho \in uco(\wp(\Sigma \times \Sigma))$ and $\eta \in uco(\wp^{\text{re}}(\mathbb{I}^* \cup \mathbb{C}))$ and $k \in \mathbb{N}$ we have that

$$M(P) \stackrel{.}{\leq} Min(M_\eta^\rho(P)) \stackrel{.}{\leq} Min_k(M_\eta^\rho(P))$$

This provides a reduction of the original SFA, and therefore CFG, providing at the same time a unique approximate representation of the abstract semantics of $P$. This is possible thanks to the combined syntactic and semantic approximation, acting both on the code and on its interpretation. Two programs $P$ and $Q$ can then be considered *similar* if they have the same reduced abstract SFA up to $k \in \mathbb{N}$:

$$P \cong_k Q \text{ iff } k = \max \{\, n \mid Min_n(M_\eta^\rho(P)) = Min_n(M_\eta^\rho(Q)) \,\}$$

This weaker notion of similarity can be improved by considering minimal SFA as canonical representation of the approximate syntax and semantics of programs:

$$P \cong Q \text{ iff } Min(M_\eta^\rho(P)) = Min(M_\eta^\rho(Q))$$

The following theorem is therefore immediate by construction.

THEOREM 6.4. *Let $P$ and $Q$ be programs, then $P \cong Q$ iff $\forall k \in \mathbb{N}: P \cong_k Q$.*

It is clear that, for decidable $\langle\rho\rangle\langle\eta\rangle$-abstractions, there exists $k \in \mathbb{N}$ such that $P \cong_k Q \implies P \cong Q$.

# 7. Formal similarity analysis of executables

The idea of *BinJuice* [18] is that the *juice* of a binary forms a template that is expected to be identical regardless of code variations due to register renaming, memory address allocation, and constant replacement. Similar ideas have been employed in *BinDiff* [12] where executables are treated as graphs of graphs: a control flow graph where each block is itself represented as a graph, which is the sequence of its instructions. While the subset of *BinDiff* considered here is sound and semantic compatible, it is computationally expensive. For large size executables, this problem has been tackled in *BinJuice* which adds a further level of abstraction to make the resulting abstract SFA more compact. In contrast to other similar tools for similarity analysis such as *DarunGrim2*, *Rdiff*, *Patchdiff*, and *Radar2*, all designed to find differences in variants of the same program for the purpose of creating patches, *BinJuice* and *BinDiff* are motivated by a different problem: Find similar code in binaries that are not known to be related. This necessitates more advanced abstractions acting on both code and semantics, therefore better showing the potential of abstract SFA.

## 7.1 *BinJuice*

*BinJuice* performs symbolic transformations on the source disassembled binary in order to transform each basic block of assembly code into a corresponding symbolic representation. The idea of symbolic execution is that the operations encoded by the assembly instructions are immediately performed when the arguments are integers, in a sort of partial evaluation local to each basic block, otherwise the same operation keeps its symbolic structure. Consider for example the following fragment of binary code and the result of its disassembly:

| Binary | Assembly |
|---|---|
| 401290: b8 05 00 00 00 | mov eax,0x5 |
| 401295: c3 04 00 00 00 | add ebx,0x4 |
| 40129b: 6b c3 | |

*BinJuice* performs algebraic manipulation of instructions in order to reach a canonical form. Thus, the result of symbolic execution with algebraic simplification of the previous example is:

| Normalized State Updates | Constraints |
|---|---|
| eax=5 | |
| ebx=def(ebx)×5 + 20 | 20 = 4 × 5 |

where $\texttt{def(ebx)}$ denotes the value of $\texttt{ebx}$ before the execution of the basic block, namely at the entry of the basic block. The syntactic information lost during symbolic execution is actually added back by the constraints on numerical values. In other words, the symbolic execution of basic blocks augmented with numerical constraints is actually an isomorphism. The key abstraction in *BinJuice* is generalization, whose idea is to use typed logical variables in order to be independent from register names. The generalization is performed by consistently replacing register names with logical variables. The replacement is consistent in that two occurrences of the same register name are always replaced by the same variable. Observe that this replacement is a purely syntactic operation. In addition to abstracting the registers used, also constants are abstracted. *BinJuice* associates a type with each logical variable to keep track of type of the original register. In the example considered before the generalization phase of *BinJuice* produces the following juice:

| Juice |
|---|
| $A = V_1$ |
| $B = \texttt{def}(B) \times N_1 + N_2$ |
| constraints: $N_2 = N_1 \times N_3$ |
| types: $\texttt{type}(A) = \texttt{type}(B) = \texttt{reg32}$ |

Let us consider the function $\mathcal{G}$ that generalizes a single basic block.

$$\mathcal{G} : \mathbb{I}^* \longrightarrow \wp^{\mathrm{re}}(SUpd) \times \wp^{\mathrm{re}}(\widehat{\mathbb{C}}) \times \wp^{\mathrm{re}}(\mathcal{T})$$

where $\wp^{\mathrm{re}}(SUpd)$ is the domain of normalized symbolic updates while $\wp^{\mathrm{re}}(\widehat{\mathbb{C}})$ is the set of constraints where register names and numerical values have been replaces by symbolic variables, and $\wp^{\mathrm{re}}(\mathcal{T})$ denotes the domain of type declarations.

We say that $\mathcal{G}(\mathsf{b})$ is the *juice* of the basic block $\mathsf{b}$. Observe that $\mathcal{G}$ acts as an abstraction since there may be more than one basic block sharing the same juice. In particular, $\mathcal{G}$ can be associated with an upper closure $\mathscr{G} \in uco(\wp^{\mathrm{re}}(\mathbb{I}^*))$ as follows:

$$\mathscr{G}(B) \overset{\text{def}}{=} \left\{ \mathsf{b}' \mid \exists \mathsf{b} \in B.\, \mathcal{G}(\mathsf{b}) = \mathcal{G}(\mathsf{b}') \right\}$$

approximating in one single symbolic representation all basic blocks that have the same juice. We can therefore model the generalization process that *BinJuice* operates on the CFG of the disassembled binaries as an $\langle\mathscr{G}\rangle$-abstraction of the predicates of the effective Boolean algebra $\mathfrak{P}$ introduced in Section 6 for representing the CFG of programs as SFA. Here, we consider the extension of $\mathscr{G}$ to branching conditions on which it behaves like identity, $\mathscr{G} \in uco(\wp^{\mathrm{re}}(\mathbb{I}^* \cup \mathbb{C}))$. The resulting *BinJuice* symbolic automaton on the Boolean algebra $\mathfrak{P}_{\mathscr{G}}$ associated with a disassembled program $P$ is:

$$M_{\mathscr{G}}(P) = \langle \mathfrak{P}_{\mathscr{G}}, \mathbb{PP}[G_P], in[b_0], \{\, out[b] \mid b \in Final[G_P]\}, \Delta_{\mathscr{G}} \rangle$$

where $\mathfrak{P}_{\mathscr{G}} = \langle \Sigma \times \Sigma, \mathscr{G}(\wp^{\mathrm{re}}(\mathbb{I}^* \cup \mathbb{C})), \{\mid \cdot \mid\}, \bot, \top, \wedge, \vee, \neg \rangle$,

$$\Delta_{\mathscr{G}} = \left\{ \, (in[\mathsf{b}], \mathscr{G}(\mathsf{b}), out[\mathsf{b}]) \mid (in[\mathsf{b}], \mathsf{b}, out[\mathsf{b}]) \in \Delta_P \right\}$$
$$\cup \left\{ \, (out[\mathsf{b}], \mathsf{c}, in[\mathsf{b}']) \mid (out[\mathsf{b}], \mathsf{c}, in[\mathsf{b}']) \in \Delta_P \right\}$$

and the semantic function $\{\mid \cdot \mid\}$ is the same as defined in Section 6 but now with a reduced abstracted domain:

$$\{\mid \cdot \mid\} : \mathscr{G}(\wp^{\mathrm{re}}(\mathbb{I}^* \cup \mathbb{C})) \longrightarrow \wp^{\mathrm{re}}(\Sigma \times \Sigma)$$

We observe that, $\mathscr{G}$ is neither syntactic nor semantic compatible (Def. 3.6, Def. 3.9) since:

(1) it collapses simplified updates with different semantics by abstracting values and variables, for example $\mathscr{G}(eax = 5) = \mathscr{G}(eax = 7) = (X = N)$;

(2) it still distinguishes between different simplified updates sharing the same semantics, as for example $\mathscr{G}(eax = ebx * 2) \neq \mathscr{G}(eax = ebx + ebx)$. But also

$$\mathscr{G}(eax = 2 * ebx + 10, \text{constraint: } 10 = 5 * 2) =$$
$$X = N_1 * Y + N_2, \text{constraint: } N_2 = N_3 * N_4$$

and

$$\mathscr{G}(eax = 2 * ebx + 10, \text{constraint: } 10 = 5 + 5) =$$
$$X = N_1 * Y + N_2, \text{constraint: } N_2 = N_3 + N_4$$

Indeed, $\mathscr{G}$ is not comparable with $\Omega(\mathtt{id})$.

PROPOSITION 7.1. *$\mathscr{G}$ is is neither syntactic nor semantic compatible.*

This observation is also related to the incorrectness of *BinJuice* in detecting similar basic blocks indeed *BinJuice* can lead to both false positives (blocks miss-classified as equivalent) and false negatives (blocks that are erroneously classified as different).

As observed before there are two causes of semantic incompatibility: (1) merging updates with different semantics and (2) distinguishing updates with the same semantics. We are interested in over-approximating $\Omega(\mathtt{id})$, namely obtaining a closure $\eta$ such that $\Omega(\mathtt{id}) \sqsubseteq \eta$, therefore avoiding (2) yet keeping (1).

A possible way for making $\mathscr{G}$ semantic compatible is to erase from the domain $\wp^{\mathrm{re}}(\mathbb{I}^* \cup \mathbb{C})$ all the elements that have the same generalized symbolic updates but different constraints. Namely by erasing all syntactic constraints. In the example above, it means for instance to restrict to the blocks that have generalized update $X = N_1 * Y + N_2$ while abstracting from the constraints on $N_2$. It is possible to prove that *BinJuice* is semantic compatible when considering this restricted domain of blocks. This highlights the fact that *BinJuice* is sensible to the structure of the constraints. Indeed, the constraints keep track of how the numerical values present in the update have been computed and is therefore tight to the particular way in which the basic block has computed them. This means that *BinJuice* can be foiled by an attacker that changes the structure of the constraints.

Define $\pi_1 : \wp^{\mathrm{re}}(SUpd) \times \wp^{\mathrm{re}}(\widehat{\mathbb{C}}) \times \wp^{\mathrm{re}}(\mathcal{T}) \to \wp^{\mathrm{re}}(SUpd)$ as the projection on the first element of the tuple of the juice. Based on this, given $\mathsf{b} \in \mathbb{I}^*$ we define the predicate abstraction $\mathcal{U}[b] \in uco(\wp(\mathbb{I}^*))$ that keeps only the blocks that have the same generalized update of $\mathsf{b}$ and abstract in $\top$ every other block:

$$\mathcal{U}[\mathsf{b}](\mathsf{b}') \begin{cases} \mathsf{b}' & \text{if } \pi_1(\mathscr{G}(\mathsf{b})) = \pi_1(\mathscr{G}(\mathsf{b}')) \\ \top & \text{otherwise} \end{cases}$$

As expected, for every basic block $\mathsf{b}$ we have that the predicate abstraction $\mathscr{G} \circ \mathcal{U}[\mathsf{b}]$, that extracts the juice of blocks that have the same generalized updates of $\mathsf{b}$, is such that $\mathcal{U}(\mathcal{U}[\mathsf{b}])$ is syntactic $\langle \mathscr{G} \circ \mathcal{U}[\mathsf{b}]\rangle$-compatible, as stated by the following result.

THEOREM 7.2. $\forall b \in \mathbb{I}^*$ *we have that* $\Omega(\mho(\mathcal{U}[b])) \sqsubseteq \mathscr{G} \circ \mathcal{U}[b]$

This result is a direct consequence of the definitions of $\mathcal{U}[\mathsf{b}]$ and of $\mathscr{G}$, and by Prop. 3.15-(2). Once again, this formally proves that *BinJuice* over-approximates the set of blocks with the same semantics when we restrict to blocks that have the same symbolic update.

## 7.2 *BinDiff*

We consider a subset of *BinDiff*, employing *instruction permutation* and *same string reference* (i.e., instructions and nodes can be matched by common string references, e.g., indicating functions that all contain code referring to the same string). All these equivalences correspond straightforwardly to abstractions of the SFA acting at syntactic and topological level. Consider the SFA $M(P)$ and the following abstractions:

**Permutation.** Let $\tau : \mathbb{I} \longrightarrow \mathcal{T}$ be a function associating the mnemonic op-code in $\mathcal{T}$ at each instruction in $\mathbb{I}$. Consider the lift of $\tau$ to multi-sets. Define an equivalence relation on basic blocks, viz., predicates in $M(P)$, such that for any $\mathsf{b}, \mathsf{b}' \in \mathbb{I}^* \cup \mathbb{C}$: $\mathsf{b} \equiv \mathsf{b}'$ if $\tau(\mathsf{b}) = \tau(\mathsf{b}')$. This clearly induces a partition which is a (partitioning) closure operator, denoted $\eta_\tau$ on predicates in $\mathbb{I}^* \cup \mathbb{C}$. In other words, $\tau$ forgets the order and the arguments of instructions. It is therefore clear that $\eta_\tau(\mathsf{b}) = \eta_\tau(\mathsf{b}') \;\not\Rightarrow\; \{\mid\mathsf{b}\mid\} = \{\mid\mathsf{b}'\mid\}$, namely $\eta_\tau$ may collapse blocks with different semantics meaning that it is not semantic compatible, i.e., $\eta_\tau \not\sqsubseteq \Omega(\mathtt{id})$. On the other hand, since $\eta_\tau$ observes precisely the multi-set of instructions, we could have blocks with the same semantics but written with different sets of instructions, i.e., $\{\mid\mathsf{b}\mid\} = \{\mid\mathsf{b}'\mid\} \not\Rightarrow \eta_\tau(\mathsf{b}) = \eta_\tau(\mathsf{b}')$ meaning that $\eta_\tau$ fails also the syntactic compatibility.

**Same reference.** Let $\mathcal{N}$ be a set of strings and $\xi : \mathbf{I} \longrightarrow \wp(\mathcal{N})$ the function associating with each basic block $\mathsf{b}$ the set of strings of $\mathcal{N}$ appearing in $\mathsf{b}$. This is clearly the left-adjoint of a GC, therefore inducing a closure $\eta_\xi$ on predicates which is also a partition. This abstraction forgets any instruction considering only a set of string manipulated in the block. Again, it is quite straightforward to observe that this abstraction can both collapse blocks

with different semantics and distinguish blocks with the same semantics, for instance a string may be computed without writing it explicitly. Hence also $\eta_\xi$ fails both the compatibilities.

In order to make *permutations* syntactic compatible, we can indeed restrict the domain of the permutation abstraction similarly to what we have done on *BinJuice* and forcing syntactic compatibility. Let $Instr(B) \stackrel{\text{def}}{=} \{\, \mathtt{b}' \mid \exists \mathtt{b} \in B.\ \eta_\tau(\mathtt{b}') = \eta_\tau(\mathtt{b}) \,\}$ and

$$\mathcal{S}[\mathtt{b}](\mathtt{b}') = \begin{cases} \mathtt{b}' & \text{if } Instr(\mathtt{b}) = Instr(\mathtt{b}') \\ \top & \text{otherwise} \end{cases}$$

As expected, for every basic block $\mathtt{b}$ we have that the predicate abstraction $Instr \circ \mathcal{S}[\mathtt{b}]$, that collects blocks that have the same set of instructions of $\mathtt{b}$, is such that $\mho(\mathcal{S}[\mathtt{b}])$ is syntactic $\langle Instr \circ \mathcal{S}[\mathtt{b}]\rangle$-compatible, as stated by the following result which is a consequence of the definitions of $\mathcal{S}[\mathtt{b}]$ and of $Instr$, and by Prop. 3.15-(2).

THEOREM 7.3. $\forall b \in \mathbb{I}^*$ *we have that* $\Omega(\mho(\mathcal{S}[b])) \sqsubseteq Instr \circ \mathcal{S}[b]$.

## 8. Related Works

To the best of our knowledge, this is the first application of abstract interpretation to symbolic finite automata and of abstract symbolic automata to similarity analysis of binary executables. The most related work is [14], where the authors introduced the notion of *lattice automata*. Lattice automata, like SFA, allow languages over an infinite alphabet. In contrast to abstract SFA, lattice automata do not distinguish between symbolic/syntactic abstractions and semantic ones. Indeed transitions in lattice automata are constrained by elements in an atomic lattice $L$, which provide precisely the allowed alphabet-set along that transition. SFA are in this context strictly more general as they separate the symbolic constraints and their semantics, allowing in principle separate approximations for them.

The idea of approximating the program's data in a so called *predicate abstraction* is nowadays common practice in static program analysis. The roots of this idea are in automatic software verification (see [1, 13]). Observe that, given a program $P$, predicate abstraction abstracts the semantics (states) of $P$ into a set of predicates $E$ and then it derives an abstract program $P$-bool that models how the execution of $P$ affects $E$. Thus, predicate abstraction corresponds to a semantic abstraction $\rho$ that groups states w.r.t. to $E$, and $P$-bool is a possible way of representing the syntactic compatible abstraction $\Omega(\rho)$ of $P$. As observed in [2] predicate abstraction considers only finite abstractions, while the semantic abstraction of denotations in abstract SFA can be an infinite domain. Moreover, predicate abstraction does not allow to change the CFG of the program.

The relation between the approximation of symbolic/syntactic structures and their semantics is well known in the literature (see [6] and [11] for a recent account). In particular in [22] the authors study this relation for the systematic synthesis of optimal symbolic predicate transformers, as introduced in [20]. None of these consider the case of abstract interpretation of SFA. In [7] the authors model disassembled binaries as finite state automata (FSA). A widening of FSA is introduced for extracting syntactic code invariants in self-modifying metamorphic programs. This construction lacks of abstractions concerning the semantics of sequences of instructions.

## 9. Conclusion

We have studied how to weaken symbolic finite automata by abstract interpretation. The results is a general theory of approximated SFA which is parametric on the chosen abstraction. The purpose is to provide a compact and effective representation of code approximations acting both at syntactic and semantic level. Interestingly, for a Turing complete programming language, there is no syntactic

abstraction which induces a compatible semantic abstraction. This follows from a simple padding argument, and it is indeed a common underlying problem in most known methods for program similarity analysis, such as *BinDiff* and *BinJuice*. Observe that the existing tools either abstract the syntax independently from the semantics, like in *BinDiff*, or represent into the syntax the abstraction of the semantics, like in predicate abstraction. In the first case we risk to fall far away form the meaning of the program to analyze, in the second case the analysis may be too much bound to the semantics without having the possibility of exploiting better syntax properties necessary in similarity analysis (e.g. in *BinDiff* and *BinJuice*). Compatibility bridges these two aspects. By semantic compatibility the abstraction of the syntax distinguishes programs with the same abstract semantics, namely the abstract program provides an under-approximation of the program behavior. By syntactic compatibility the abstraction of the syntax collapses programs with different semantics, hence capturing behaviors that are not related with the program to analyze, therefore providing an over-approximation of the program behavior. Interestingly, in our model we can restrict the form of predicates in order to have compatibility. This is what we proved in *BinJuice* and *BinDiff*, thus showing the limits of existing tools for code similarity and the possibility of systematically deriving conditions for making them syntactic compatible.

Another direction of future research is in the use of topological abstractions of SFA for extracting signatures of self-modifying code as recently studied in [7]. This requires the extension of widening operations, such as those introduced in [5, 10, 14], to abstract SFA. In this case approximate SFA provide advanced signatures in metamorphic malware analysis, incorporating both properties of way code changes during program execution (the invariant of the metamorphic engine) and additional semantic information, such as the values passed in system-calls. This may reduce the false positives occurring in [7] in signature-based detection of metamorphic malware.

## Acknowledgments

## References

[1] T. Ball, R. Majumdar, T. D. Millstein, and S. K. Rajamani. Automatic predicate abstraction of C programs. In M. Burke and M. L. Soffa, editors, *PLDI*, pages 203–213. ACM, 2001. ISBN 1-58113-414-2.

[2] P. Cousot. Verification by abstract interpretation. In *Verification: Theory and Practice, Essays Dedicated to Zohar Manna on the Occasion of His 64th Birthday*, volume 2772, pages 243–268. Springer, 2003.

[3] P. Cousot and R. Cousot. Abstract interpretation: A unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Conference Record of the 4th ACM Symposium on Principles of Programming Languages* (*POPL '77*), pages 238–252. ACM Press, 1977.

[4] P. Cousot and R. Cousot. Systematic design of program analysis frameworks. In *Conference Record of the 6th ACM Symposium on Principles of Programming Languages* (*POPL '79*), pages 269–282. ACM Press, 1979.

[5] P. Cousot and R. Cousot. Formal language, grammar and set-constraint-based program analysis by abstract interpretation. In *Proceedings of the Seventh ACM Conference on Functional Programming Languages and Computer Architecture*, pages 170–181. ACM Press, New York, NY, 25–28 June 1995.

[6] P. Cousot, R. Cousot, and L. Mauborgne. Theories, solvers and static analysis by abstract interpretation. *J. ACM*, 59(6):31, 2012.

[7] M. Dalla Preda, R. Giacobazzi, S. K. Debray, K. Coogan, and G. M. Townsend. Modelling metamorphism by abstract interpretation. In *Proc. of the 19th Int. Static Analysis Symp. (SAS '10)*, volume 6337 of *Lecture Notes in Computer Science*, pages 218–235. Springer-Verlag, Berlin, 2010.

[8] L. D'Antoni and M. Veanes. Equivalence of extended symbolic finite transducers. In N. Sharygina and H. Veith, editors, *CAV*, volume 8044 of *Lecture Notes in Computer Science*, pages 624–639. Springer, 2013. ISBN 978-3-642-39798-1.

[9] L. D'Antoni and M. Veanes. Minimization of symbolic automata. In S. Jagannathan and P. Sewell, editors, *POPL*, pages 541–554. ACM, 2014. ISBN 978-1-4503-2544-8.

[10] V. D'Silva. Widening for automata. Diploma Thesis, Institut Fur Informatick, Universitat Zurich, 2006.

[11] V. D'Silva, L. Haller, and D. Kroening. Abstract satisfaction. In S. Jagannathan and P. Sewell, editors, *POPL*, pages 139–150. ACM, 2014. ISBN 978-1-4503-2544-8.

[12] H. Flake. Structural comparison of executable objects. In U. Flegel and M. Meier, editors, *DIMVA*, volume 46 of *LNI*, pages 161–173. GI, 2004. ISBN 3-88579-375-X.

[13] C. Flanagan and S. Qadeer. Predicate abstraction for software verification. In *Proc. of Conf. Record of the 29th ACM Symp. on Principles of Programming Languages (POPL '02)*, pages 191–202. ACM Press, 2002.

[14] T. L. Gall and B. Jeannet. Lattice automata: A representation for languages on infinite alphabets, and some applications to verification. In H. R. Nielson and G. Filé, editors, *SAS*, volume 4634 of *Lecture Notes in Computer Science*, pages 52–68. Springer, 2007. ISBN 978-3-540-74060-5.

[15] D. Gao, M. Reiter, and D. Song. BinHunt: Automatically finding semantic differences in binary programs. In *Proceedings of the 10th International Conference on Information and Communications Security*, ICICS '08, pages 238–255. Springer-Verlag, 2008.

[16] R. Giacobazzi, F. Ranzato, and F. Scozzari. Making abstract interpretation complete. *Journal of the ACM*, 47(2):361–416, March 2000.

[17] S. Hunt and I. Mastroeni. The PER model of abstract non-interference. In C. Hankin and I. Siveroni, editors, *Proc. of The 12th Internat. Static Analysis Symp. (SAS '05)*, volume 3672 of *Lecture Notes in Computer Science*, pages 171–185. Springer-Verlag, 2005.

[18] A. Lakhotia, M. Dalla Preda, and R. Giacobazzi. Fast location of similar code fragments using semantic 'juice'. In *2nd Workshop on Program Protection and Reverse Engineering PPREW 2013*. ACM, 2013.

[19] I. Mastroeni and R. Giacobazzi. An abstract interpretation-based model for safety semantics. *Int. J. Comput. Math.*, 88(4):665–694, 2011.

[20] T. W. Reps, S. Sagiv, and G. Yorsh. Symbolic implementation of the best transformer. In B. Steffen and G. Levi, editors, *VMCAI*, volume 2937 of *Lecture Notes in Computer Science*, pages 252–266. Springer, 2004. ISBN 3-540-20803-8.

[21] H. Rogers. *Theory of recursive functions and effective computability*. The MIT press, 1992.

[22] A. V. Thakur, M. Elder, and T. W. Reps. Bilateral algorithms for symbolic abstraction. In A. Miné and D. Schmidt, editors, *SAS*, volume 7460 of *Lecture Notes in Computer Science*, pages 111–128. Springer, 2012. ISBN 978-3-642-33124-4.

[23] M. Veanes, P. Hooimeijer, B. Livshits, D. Molnar, and N. Bjørner. Symbolic finite state transducers: algorithms and applications. In J. Field and M. Hicks, editors, *POPL*, pages 137–150. ACM, 2012. ISBN 978-1-4503-1083-3.

[24] M. Ward. The Closure Operators of a Lattice. *Annals of Mathematics*, 43(2):191–196, 1942.

[25] Zynamics. BinDiff3.2manual., 2004. URL `http://www.zynamics.com/bindiff/manual/`.