

Static analysis of the determinism of multithreaded programs

Pietro Ferrara *

École Polytechnique

F-91128 Palaiseau (France)

Dipartimento di Informatica

Università Ca' Foscari di Venezia

I-30170 Venezia (Italy)

Pietro.Ferrara@polytechnique.edu

Abstract

Threads communicate implicitly through shared memory. Because of the random interleaving during their parallel execution, nondeterministic behaviors possibly arise that is why multithreaded programming is strictly more difficult than programming in sequential languages. Moreover the random interleaving may lead to subtle bugs, that are really hard to be detected and fixed. We propose a novel deterministic property focused on multithreading. We define it as difference among concrete traces, and then we abstract it in two separate steps in order to statically analyze it. At the intermediate level of abstraction, we propose the new idea of weak determinism. We sketch how the proposed property may be used in order to semi-automatically parallelize sequential programs. Finally, we present some experimental results when applying the analysis to a set of well-known benchmarks. We believe that our approach, dealing directly with the source of the problem (i.e. the nondeterministic interactions via shared memory) is in position to bypass the actual limits of the static analysis of multithreaded programs, mostly focused on properties like data race condition and absence of deadlocks.

1. Introduction

While the improvement of single-core processors performances is slowing down, multi-core solutions appear to be the most promising way to extend the Moore's Law into the future [10]. The current trend of the CPU market underlines this idea: for instance the new family of Intel® processors (Intel® Core™) is dual or quad core, and the same happens for the AMD™ ones. Even if any application that runs on

*Work partially supported by MIUR PRIN'07 Project SOFT - Tecniche formali orientate alla sicurezza

a single-core processor runs also on a multi-core one, in order to fully take advantage of these capabilities the applications must be optimized for multithreading. It is necessary to develop programs with explicit parallelism, otherwise a single application would not take relevant advantage from the multi-core hardware.

For developers, reasoning about concurrency and multithreading is strictly more difficult than working on sequential programs [18]. A first level of complexity is due to the interleaving of the executions of different threads, and to their implicit communications through shared memory. A more subtle level of complexity is induced by the specification of the memory model of the programming language [13], that is often not completely clear and understandable for the developer.

Both these issues are related to a common problem: the non-determinism induced by the parallel execution of many threads. In fact, some unexpected and counter-intuitive behaviors arise, like the ones due to the local optimization of compilers, invisible at single thread level.

A first way to restrict the non-determinism of multithreaded programs has been the static or dynamic checking on the concurrent actions [16]. Another way has been to define specific models that are in position to guarantee more restrictive multithreaded executions. Even if the work of the researchers in this field has been huge, there are still important limits in the proposed static analyses [18].

In this context, the contribution of this paper is mainly addressed into three distinct directions: (i) defining and abstracting in two steps the deterministic property in order to statically analyze it; (ii) defining on the first level of approximation the novel idea of weak determinism; (iii) sketching how determinism and its relaxations may be used in order to semi-automatically parallelize sequential programs.

In this way, we first define the determinism on a concrete domain and semantics, and then we abstract it, proving for-

mally the soundness of our approach. In particular, at the first level of abstraction for each variable we trace a value for each thread that may have written it, while at the second level we trace a single value that approximates all the possible values and the set of threads that may have written in parallel on that variable.

Then we define the weak determinism, a novel idea that can be seen as a way of relaxing the deterministic property. We sketch how the deterministic property may be used in order to semi-automatically parallelize sequential programs. The analysis of these properties has been implemented, and we report some experimental results when applying it to a set of well-known benchmarks. As far as we know, our work is the first attempt to define and statically analyze the determinism of a multithreaded program generic on the programming language.

The rest of this section introduces a running example, some basic concepts of abstract interpretation, and the notation. Section 2 presents the concrete semantics, while Section 3 and 4 present respectively the first and second level of abstraction. The deterministic property is defined by Section 5, while the weak determinism is introduced in Section 6. Section 7 sketches how to use the deterministic property in order to semi-automatically parallelize sequential programs. Section 8 presents the experimental results. The related works are presented by section 9, and Section 10 concludes.

1.1. The running example

```
public class Account {
    public int amount=0;
    public void withdraw(int money) {
        synchronized(this) {
            account.amount-=money;
        }
    }
    public void deposit(int money) {
        synchronized(this) {
            account.amount+=money;
        }
    }
}
```

(a) The class Account

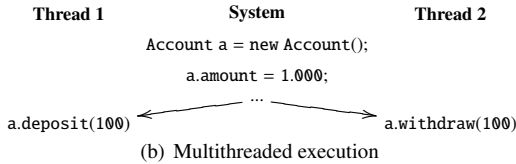


Figure 1. A running example

In order to illustrate the concepts presented throughout all the paper, we will always refer to the following running example written in Java style code.

Figure 1.(a) depicts the class Account that stores in a pro-

ected field amount an amount of money. Two methods are defined (`withdraw` and `deposit`), and they are used to withdraw and deposit money; they are both synchronized on the monitor defined on the current object. Figure 1.(b) depicts the multithreaded execution: the system is going to execute in parallel a deposit of 100€ and a withdraw of the same amount of money. At the beginning the amount of money contained in the account is 1.000€.

1.2. Abstract interpretation

Abstract interpretation is a theory to define and soundly approximate the semantics of a program [3, 4]. Roughly, a concrete semantics, aimed at specifying the runtime properties of interests, is defined; then it is approximated through one or more steps in order to finally obtain an abstract semantics that is computable, but still precise enough to capture the property of interest. The abstract semantics must be composed of an abstract domain, and an abstract transition function. Usually, each state of the concrete domain is composed by a set of elements (e.g. all the possible computational states), that is approximated by a unique element in the abstract domain.

Formally, the concrete domain D forms a complete lattice $\langle D, \sqsubseteq_C, \perp_C, \top_C, \sqcup_C, \sqcap_C \rangle$. On this domain, a semantics \mathbb{S} is defined. The concrete elements are related to the abstract domain $\langle D^\#, \sqsubseteq_A, \perp_A, \top_A, \sqcup_A, \sqcap_A \rangle$ by a concretization function γ and an abstraction function α . In order to obtain a sound analysis, we require that they form a Galois connection; we denote it by $\langle D, \sqsubseteq_C, \perp_C, \top_C, \sqcup_C, \sqcap_C \rangle \xleftrightarrow[\alpha]{\gamma} \langle D^\#, \sqsubseteq_A, \perp_A, \top_A, \sqcup_A, \sqcap_A \rangle$.

Since the most part of abstract domains does not satisfy the ascending chain condition, a widening operator ∇_A is required. This is an upper bound operator such that for all increasing chains $a_0^\# \sqsubseteq_A \dots a_n^\# \sqsubseteq_A \dots$ the increasing chain defined as $w_0^\# = a_0^\#, \dots w_{i+1}^\# = w_i^\# \nabla_A a_{i+1}^\#$ is not strictly increasing.

Finally, the abstract semantics $\mathbb{S}^\#$ has to soundly approximate the concrete one, i.e. $\forall d^\# \in D^\# : \alpha \circ \mathbb{S}[\gamma(d^\#)] \sqsubseteq_A \mathbb{S}^\#[d^\#]$.

1.3. Notation

Let S be a generic set. We denote by S^\ddagger the set of all the finite traces composed by elements in S . A trace is an ordered sequence of elements. We denote by $len : [S^\ddagger \rightarrow \mathbb{N}]$ the function that given a trace returns its length, i.e. the number of elements contained in it. We denote by $istate : [(S^\ddagger \times \mathbb{N}) \rightarrow \mathbb{N}]$ the function that returns the i -th element of the given trace, where i is the integer value passed to the function.

As already stated in the introduction, we will introduce two levels of abstraction, and consequently we will denote by

S and s concrete sets and elements, by $S^\#$ and $s^\#$ sets and elements of the first level of abstraction, and by S° and s° sets and elements of the second level. In the same way, we will denote by $func$ a function on concrete sets, by $func^\#$ the corresponding function at the first level of abstraction, and $func^\circ$ at the second level.

2. Syntax and concrete semantics

In this section we present a syntax focused on the write and read operations on shared memory, and the augmented concrete domain and semantics, in order to trace for each value the thread that wrote it in the shared memory.

2.1. Syntax

For the sake of readability, we consider a very restricted syntax, focused on the interactions with a shared memory, i.e. the read and the write actions. In this way we consider only statements of the form `sh_var = value` and `read(sh_var)`.

The statement `sh_var = value` writes on a shared variable `sh_var` a given value, while `read(sh_var)` reads a shared variable `sh_var` and returns its value. Note that this syntax can be easily extended with other statements (`if`, `while`, etc..).

2.2. Concrete domain

The concrete domain we consider is strictly focused on the shared memory. For each value, we trace the thread that wrote it. Since we are in the concrete context, at each point we know exactly which thread wrote a value, and we may relate each shared variable to a pair composed by its value and an identifier of that thread.

Formally, let T be the set of the identifiers of threads, Var be the set of shared variables, and V be the set of concrete values. The shared memory is a function that relates each variable to a pair composed by its value and the identifier of the thread that wrote it ($S : [Var \rightarrow (V \times T)]$). Note that this representation of shared memory is quite different with respect to the approach usually adopted by current programming languages (e.g. `Java`), where the shared memory is the heap, and the values are accessed by reference. At static level, this approach would lead to multiple issues orthogonal w.r.t. our goal, in particular alias analysis, that deserve to be considered separately [6].

The idea behind a multithreaded program is to execute each thread in parallel. The thread-partitioned trace domain is aimed at formalizing this concept. The same approach has been adopted in previous works [7, 6].

The thread-partitioned domain represents an execution as

a function that, given a thread, returns the trace containing only the shared memory locations obtained by its computation ($\Psi : [T \rightarrow S^\#]$). Each trace collects the execution of the thread to which it is related. All the possible executions of a program are represented by a set of these functions. As usual in abstract interpretation, the concrete domain is the lattice obtained using powerset operators ($\langle \wp(\Psi), \subseteq, \emptyset, \wp(\Psi), \cup, \cap \rangle$).

2.3. Transfer function

In order to define the concrete semantics, we first introduce the $\mathbb{W} : [(S \times T) \rightarrow S]$ of `sh_var = value` defined as $\mathbb{W}[\text{sh_var} = \text{value}](s, t) = s[\text{sh_var} \mapsto (\text{value}, t)]$.

The value of the shared variable when evaluating a shared memory s is given by the function $\mathbb{R} : [S \rightarrow V]$ of `eval(sh_var)` defined as $\mathbb{R}[\text{eval}(\text{sh_var})](s) = \pi_1(s(\text{sh_var}))$ where π_1 is the projection on the first component.

These two functions may be used in order to define a transfer function for a more complex language; usually in the abstract interpretation framework such a function is the basic step in order to define the concrete semantics as the result of fixpoint computation.

2.4. The running example

Thread 1 : [a.amount \mapsto (1.000, System)] \rightarrow [a.amount \mapsto (1.100, Thread 1)]
Thread 2 : [a.amount \mapsto (1.100, Thread 1)] \rightarrow [a.amount \mapsto (1.000, Thread 2)]
Thread 1 : [a.amount \mapsto (900, Thread 2)] \rightarrow [a.amount \mapsto (1.000, Thread 1)]
Thread 2 : [a.amount \mapsto (1.000, System)] \rightarrow [a.amount \mapsto (900, Thread 2)]

Figure 2. The concrete semantics

Figure 2 presents the two multithreaded concrete executions that can be obtained using the domain and the semantics just introduced. The first multithreaded trace represents the execution in which `Thread 1` is executed before `Thread 2`, while the second state depicts the opposite situation. Note that other executions are not possible, as both the methods are synchronized on the same monitor.

3. A value for each thread (abstraction 1)

In this section we present the first level of abstraction, where each shared variable is related to an abstract value for each thread that may write on that. Our analysis is parameterized on abstract domain that approximates numerical val-

ues (e.g. sign domain [3] that traces if they are positive or negative).

3.1. Abstract domain (first level)

We consider the following Galois connection between the concrete domain of values and its abstract counterpart:

$$\langle \wp(\mathbb{V}), \subseteq, \emptyset, \wp(\mathbb{V}), \cup, \cap \rangle \xleftrightarrow[\alpha_V]{\gamma_V} \langle \mathbb{V}^\#, \sqsubseteq_{V^\#}, \perp_{V^\#}, \top_{V^\#}, \sqcup_{V^\#}, \sqcap_{V^\#} \rangle$$

This means that the non-relational abstract domain for numerical values has to soundly approximate the concrete values. We also suppose that the identifiers of threads are the same in the concrete and in the abstract. This hypothesis may appear restrictive: current programming languages, like Java, identify threads by reference; the syntactic identification of threads is quite naive. On the other hand, we are presenting a new property on multithreaded programs; this property can be straightly extended also to a context on which threads are identified, for instance, by reference. In this case, it suffices to parameterize the analysis on the abstraction of references (e.g. a point-to analysis). For sake of simplicity, we suppose that threads are statically identified; on the other hand, this supposition still preserves the generality of our analysis. Moreover a previous work [6] has already proposed an alias analysis that is in position to statically identify Java threads.

In the concrete context, different executions may contain values written by different threads because of their random interleavings. This first level of abstraction traces a value for each thread; in this way it gathers all the values written by the same thread in the same abstract element.

Formally, the shared memory $\mathbf{s}^\#$ relates each shared variable to a function that maps a thread to the abstraction of the values that it may have written, i.e. $\mathbf{S}^\# : [\text{Var} \rightarrow [\text{T} \rightarrow \mathbb{V}^\#]]$.

As in the concrete semantics, we represent an execution as a function that, given a thread, returns the trace containing only the shared memories obtained by its computation ($\Psi^\# : [\text{T} \rightarrow \mathbf{S}^{\#^\tau}]$), where $\mathbf{S}^{\#^\tau}$ denotes the set of all the finite traces with elements in $\mathbf{S}^\#$.

3.2. Upper bound operator

The upper bound operator on shared memories keeps all the values written by different threads; if there is a value related to the same thread in the two shared memories we make the join between these two values. With an abuse of notation, we suppose that if $\text{var} \notin \text{dom}(\mathbf{s})$, then

$$\text{dom}(\mathbf{s}(\text{var})) = \emptyset.$$

$$\begin{aligned} \mathbf{s}_1^\# \sqcup_{\mathbf{S}^\#} \mathbf{s}_2^\# &= \mathbf{s}^\# : \forall \text{var} \in \text{dom}(\mathbf{s}_1^\#) \cup \text{dom}(\mathbf{s}_2^\#), \\ &\quad \forall t \in \text{dom}(\mathbf{s}_1^\#(\text{var})) \cup \text{dom}(\mathbf{s}_2^\#(\text{var})), \\ \mathbf{s}^\#(\text{var})(t) &= \begin{cases} \mathbf{s}_1^\#(\text{var})(t) & \text{if } \text{var} \notin \text{dom}(\mathbf{s}_2^\#) \vee \\ & t \notin \text{dom}(\mathbf{s}_2^\#(\text{var})) \\ \mathbf{s}_2^\#(\text{var})(t) & \text{if } \text{var} \notin \text{dom}(\mathbf{s}_1^\#) \vee \\ & t \notin \text{dom}(\mathbf{s}_1^\#(\text{var})) \\ \mathbf{s}_1^\#(\text{var})(t) \sqcup_{V^\#} \mathbf{s}_2^\#(\text{var})(t) & \text{otherwise} \end{cases} \end{aligned}$$

The upper bound operator on traces simply applies the upper bound operator of shared memories on all the elements of all the traces of the two states.

$$\begin{aligned} \tau_1^\# \sqcup_{\mathbf{S}^\#} \tau_2^\# &= \sigma_0^\# \rightarrow \dots \rightarrow \sigma_i^\# : i = \max(\text{len}(\tau_1^\#), \text{len}(\tau_2^\#)), \\ \forall j \in [0..i] : \sigma_j^\# &= \begin{cases} \text{istate}(\tau_1^\#, j) & \text{if } j < \text{len}(\tau_1^\#) \wedge \\ & \sqcup_{\mathbf{S}^\#} \text{istate}(\tau_2^\#, j) \quad j < \text{len}(\tau_2^\#) \\ \text{istate}(\tau_2^\#, j) & \text{if } j \geq \text{len}(\tau_1^\#) \\ \text{istate}(\tau_1^\#, j) & \text{if } j \geq \text{len}(\tau_2^\#) \end{cases} \end{aligned}$$

The upper bound operator on the multithreaded state is the pointwise extension of the upper bound of traces on all the elements of the codomain.

$$\begin{aligned} f_1^\# \sqcup_{\Psi^\#} f_2^\# &= f^\# : \forall t \in \text{dom}(f_1^\#) \cup \text{dom}(f_2^\#) : \\ f^\#(t) &= \begin{cases} f_1^\#(t) \sqcup_{\mathbf{S}^\#} f_2^\#(t) & \text{if } t \in \text{dom}(f_1^\#) \cap \text{dom}(f_2^\#) \\ f_1^\#(t) & \text{if } t \in \text{dom}(f_1^\#) \wedge t \notin \text{dom}(f_2^\#) \\ f_2^\#(t) & \text{if } t \in \text{dom}(f_2^\#) \wedge t \notin \text{dom}(f_1^\#) \end{cases} \end{aligned}$$

3.3. Abstraction function

The abstraction function maps a set of concrete shared memories into an abstract memory that relates each variable to a function that associates each thread with the abstraction of the values it may have written.

$$\begin{aligned} \alpha_S : [\wp(\mathbf{S}) \rightarrow \mathbf{S}^\#] \\ \alpha_S(\mathbf{S}) &= f^\# : \forall \text{var} \in \bigcup_{\mathbf{s} \in \mathbf{S}} \text{dom}(\mathbf{s}), \forall t \in \bigcup_{\mathbf{s} \in \mathbf{S}} \text{dom}(\mathbf{s}(\text{var})) : \\ f^\#(\text{var})(t) &= \alpha_V(\{v : \exists \mathbf{s} \in \mathbf{S} : \mathbf{s}(\text{var}) = (v, t)\}) \end{aligned}$$

The abstraction of a set of traces produces an abstract trace such that its i -th element is the abstraction of the i -th elements of all the given concrete traces.

$$\begin{aligned} \alpha_{S^\#} : [\wp(\mathbf{S}^\#) \rightarrow \mathbf{S}^{\#^\tau}] \\ \alpha_{S^\#}(\mathbf{T}) &= \sigma_0^\# \rightarrow \dots \rightarrow \sigma_i^\# : i = \max(\bigcup_{\tau \in \mathbf{T}} \text{len}(\tau)), \\ \forall j \in [0..i] : \sigma_j^\# &= \alpha_S(\bigcup_{\tau \in \mathbf{T} : \text{len}(\tau) > j} \text{istate}(\tau, j)) \end{aligned}$$

When considering traces, the abstraction function returns an unique function that abstracts together all the traces produced by the same thread in the different executions.

$$\begin{aligned} \alpha_\Psi : [\wp(\Psi) \rightarrow \Psi^\#] \\ \alpha_\Psi(\Phi) &= f^\# : \forall t \in \bigcup_{f \in \Phi} \text{dom}(f) : \\ f^\#(t) &= \alpha_{S^\#}(\bigcup_{f \in \Phi : t \in \text{dom}(f)} f(t)) \end{aligned}$$

Theorem 1 (Soundness of $\langle \Psi^\#, \sqsubseteq_{\Psi^\#} \rangle$)

$$\langle \wp(\Psi), \subseteq \rangle \xleftrightarrow[\alpha_\Psi]{\gamma_\Psi} \langle \Psi^\#, \sqsubseteq_{\Psi^\#} \rangle$$

3.4. Transfer function

In correspondence to the semantic functions \mathbb{W} and \mathbb{R} , we introduce their abstract counterpart $\mathbb{W}^\#$ and $\mathbb{R}^\#$, as follows.

$\mathbb{W}^\# : [(\mathbb{S}^\# \times \mathbb{T}) \rightarrow \mathbb{S}^\#]$ of $\text{sh_var} = \text{value}$ is defined as

$$\begin{aligned} \mathbb{W}^\#[[\text{sh_var} = \text{value}]](\mathbf{s}^\#, t) &= \\ &= \mathbf{s}^\#[\text{sh_var} \mapsto [t \mapsto \alpha_V(\text{value})]] \end{aligned}$$

$\mathbb{R}^\# : [\mathbb{S}^\# \rightarrow \mathbb{V}^\#]$ of $\text{eval}(\text{sh_var})$ is defined as

$$\mathbb{R}^\#[[\text{eval}(\text{sh_var})]](\mathbf{s}^\#) = \bigsqcup_{t \in \text{dom}(\mathbf{s}^\#(\text{sh_var}))} \mathbf{s}^\#(\text{sh_var})(t)$$

As in the concrete semantics, these two functions may be used in order to define a transfer function for a more complex language.

Lemma 1 (Soundness) $\mathbb{W}^\#$ and $\mathbb{R}^\#$ are respectively the abstraction of \mathbb{W} and \mathbb{R} .

3.5. The running example

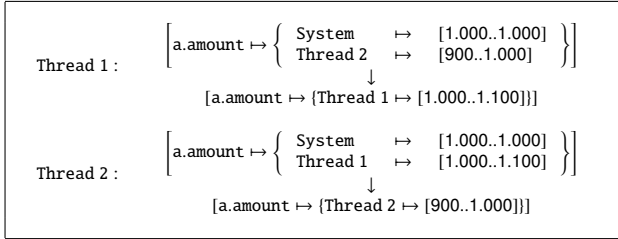


Figure 3. The abstract semantics

Figure 3 depicts the abstract semantics obtained by using the domain and the semantics just introduced. In order to capture numerical information we use the interval domain [3], where $[a..b]$ represents the set of concrete values S such that $\forall n \in S : a \leq n \leq b$.

4. Just one value (abstraction 2)

In this section we present the second level of abstraction, in which the values written by different threads collapse in the same abstract element.

4.1. Abstract domain (second level)

A shared memory $\mathbf{s}^\circ \in \mathbb{S}^\circ$ relates each variable var to a pair composed by a value and the set of threads that may have written it in var , formally $\mathbb{S}^\circ : [\text{Var} \rightarrow (\mathbb{V}^\# \times \wp(\mathbb{T}))]$. We represent an execution as a function $f^\circ \in \mathbb{S}^\circ$ that, given a thread, returns the trace containing only the shared memories obtained by its computation ($\Psi^\circ : [\mathbb{T} \rightarrow \mathbb{S}^{\circ\ddagger}]$).

4.2. Upper bound operator

The upper bound operator on shared memories makes the upper bound between values and the set union between the sets of threads that may have written the given value.

$$\begin{aligned} \mathbf{s}_1^\circ \sqcup_{\mathbb{S}^\circ} \mathbf{s}_2^\circ = \mathbf{s}^\circ : \forall \text{var} \in \text{dom}(\mathbf{s}_1^\circ) \cup \text{dom}(\mathbf{s}_2^\circ) : \mathbf{s}^\circ(\text{var}) &= \\ &= \begin{cases} (\pi_1(\mathbf{s}_1^\circ(\text{var})) \sqcup_{\mathbb{V}^\#} \pi_1(\mathbf{s}_2^\circ(\text{var})), & \text{if } \text{var} \in \text{dom}(\mathbf{s}_1^\circ) \cap \text{dom}(\mathbf{s}_2^\circ) \\ \pi_2(\mathbf{s}_1^\circ(\text{var})) \cup \pi_2(\mathbf{s}_2^\circ(\text{var})) & \text{dom}(\mathbf{s}_2^\circ) \\ \mathbf{s}_1^\circ(\text{var}) & \text{if } \text{var} \in \text{dom}(\mathbf{s}_1^\circ) \wedge \text{var} \notin \text{dom}(\mathbf{s}_2^\circ) \\ \mathbf{s}_2^\circ(\text{var}) & \text{if } \text{var} \in \text{dom}(\mathbf{s}_2^\circ) \wedge \text{var} \notin \text{dom}(\mathbf{s}_1^\circ) \end{cases} \end{aligned}$$

The upper bound operator on traces simply applies the upper bound operator of shared memories on all the elements of all the traces of the two states.

$$\begin{aligned} \tau_1^\circ \sqcup_{\mathbb{S}^\circ\ddagger} \tau_2^\circ = \sigma_0^\circ \rightarrow \dots \rightarrow \sigma_i^\circ : i = \max(\text{len}(\tau_1^\circ), \text{len}(\tau_2^\circ)), \\ \forall j \in [0..i] : \sigma_j^\circ = \begin{cases} \text{istate}(\tau_1^\circ, j) \sqcup_{\mathbb{S}^\circ} \text{istate}(\tau_2^\circ, j) & \text{if } j < \text{len}(\tau_1^\circ) \wedge j < \text{len}(\tau_2^\circ) \\ \text{istate}(\tau_1^\circ, j) & \text{if } j < \text{len}(\tau_1^\circ) \wedge j \geq \text{len}(\tau_2^\circ) \\ \text{istate}(\tau_2^\circ, j) & \text{if } j < \text{len}(\tau_2^\circ) \wedge j \geq \text{len}(\tau_1^\circ) \end{cases} \end{aligned}$$

The upper bound operator on the multithreaded state is the pointwise extension of the upper bound of traces on all the elements of the codomain.

$$\begin{aligned} f_1^\circ \sqcup_{\Psi^\circ} f_2^\circ = f^\circ : \forall t \in \text{dom}(f_1^\circ) \cup \text{dom}(f_2^\circ) : \\ f^\circ(t) = \begin{cases} f_1^\circ(t) \sqcup_{\mathbb{S}^\circ\ddagger} f_2^\circ(t) & \text{if } t \in \text{dom}(f_1^\circ) \cap \text{dom}(f_2^\circ) \\ f_1^\circ(t) & \text{if } t \in \text{dom}(f_1^\circ) \wedge t \notin \text{dom}(f_2^\circ) \\ f_2^\circ(t) & \text{if } t \in \text{dom}(f_2^\circ) \wedge t \notin \text{dom}(f_1^\circ) \end{cases} \end{aligned}$$

4.3. Abstraction function

The abstraction function abstracts an element of the first level of abstraction.

$$\begin{aligned} \alpha_{\mathbb{S}^\#} : [\mathbb{S}^\# \rightarrow \mathbb{S}^\circ] \\ \alpha_{\mathbb{S}^\#}(\mathbf{s}^\#) = \mathbf{s}^\circ : \forall \text{var} \in \text{dom}(\mathbf{s}^\#) : \\ \mathbf{s}^\circ(\text{var}) = (\bigsqcup_{t \in \text{dom}(\mathbf{s}^\#(\text{var}))} \mathbf{s}^\#(t), \text{dom}(\mathbf{s}^\#(\text{var}))) \end{aligned}$$

The abstraction of traces is the pointwise application of the abstraction of the shared memory to all the elements of all the traces.

$$\begin{aligned} \alpha_{\mathbb{S}^\#\ddagger} : [\mathbb{S}^\#\ddagger \rightarrow \mathbb{S}^{\circ\ddagger}] \\ \alpha_{\mathbb{S}^\#\ddagger}(\sigma_0^\# \rightarrow \dots \rightarrow \sigma_i^\#) = \alpha_{\mathbb{S}^\#}(\sigma_0^\#) \rightarrow \dots \rightarrow \alpha_{\mathbb{S}^\#}(\sigma_i^\#) \end{aligned}$$

The abstraction on the multithreaded state is the pointwise application of the abstraction of traces.

$$\begin{aligned} \alpha_{\Psi^\#} : [\Psi^\# \rightarrow \Psi^\circ] \\ \alpha_{\Psi^\#}(f^\#) = f^\circ : \forall t \in \text{dom}(f^\#) : f^\circ(t) = \alpha_{\mathbb{S}^\#\ddagger}(f^\#(t)) \end{aligned}$$

Theorem 2 (Soundness of $\langle \Psi^\circ, \sqsubseteq_{\Psi^\circ} \rangle$)

$$\langle \Psi^\#, \sqsubseteq_{\Psi^\#} \rangle \xleftarrow[\alpha_{\Psi^\#}]{\gamma_{\Psi^\#}} \langle \Psi^\circ, \sqsubseteq_{\Psi^\circ} \rangle$$

4.4. Transfer function

The transfer function $\mathbb{W}^\circ : [(S^\circ \times T) \rightarrow S^\circ]$ of `sh_var = value` is defined as

$$\mathbb{W}^\circ \llbracket \text{sh_var} = \text{value} \rrbracket (s^\circ, t) = s^\circ[\text{sh_var} \mapsto (\text{value}^\#, \{t\})]$$

Indeed the function $\mathbb{R}^\circ : [S^\circ \rightarrow V^\#]$ of `eval(sh_var)` is defined as

$$\mathbb{R}^\circ \llbracket \text{eval}(\text{sh_var}) \rrbracket (s^\circ) = \pi_1(s^\circ(\text{sh_var}))$$

Lemma 2 (Soundness) \mathbb{W}° and \mathbb{R}° are respectively the abstraction of $\mathbb{W}^\#$ and $\mathbb{R}^\#$.

4.5. The running example

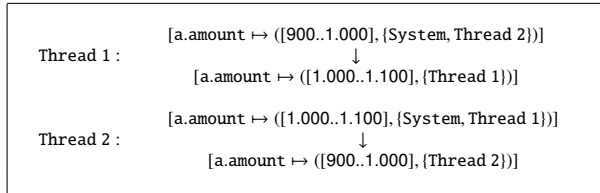


Figure 4. The abstract semantics on the running example

Figure 4 presents the abstract semantics obtained using the domain and the semantics just introduced. Also in this case we use the interval domain in order to capture numerical information.

5. The deterministic property

In the last three sections, we have presented the concrete domain and semantics and the two levels of abstraction. In this section we define on them the deterministic property.

5.1. Determinism

A program is generally said to be deterministic “if given an input it returns always the same output”[5]. This definition does not specify what is the input and the output of a program. Someone may think that the output is what it writes on the screen, some others that it is the state of the memory at the end of the execution, or during it. In a multithreaded context, a too restrictive application of the concept of determinism may lead to define that a multithreaded program is deterministic if all the statements, even of different threads, are executed always in the same total order.

We face this problem through the thread-partitioned domain, that abstracts away the inter-thread order in which the statements are executed. At this level a program is deterministic if, given an input, returns the same output for any order of execution.

Moreover our work is focused on the non-determinism induced by the random interleavings of the execution of different threads. As we do not make any supposition or restriction on the programming language and its semantics, many different forms of non-determinism may exist: a random number produced by the system, and the inputs received by an operator are two typical examples of non-determinism not induced by the multithreaded execution. We want to catch and analyze only the non-deterministic behaviors due to the parallel executions of threads, discarding all the other cases.

This goal can be reached by using the information collected by our analysis, as we trace for each value on the shared memory which thread may have written it: it suffices to check if at a given point (i.e. for a given state of execution of a given thread) a value could have been written by two different threads.

5.2. Formal definition of determinism on the concrete domain

We first define the determinism as difference between two states, and then we apply it to a set of elements of the concrete multithreaded semantics.

Definition 1 (Determinism on shared memory) *Given two states, they underline a non deterministic behavior due to the multithreaded execution if the two states relate the same variable to values written by different threads.*

$$\begin{aligned}
DS &: [S \times S \rightarrow \{\text{true}, \text{false}\}] \\
DS(s_1, s_2) &= \text{false} \\
&\Downarrow \\
\exists \text{var} \in \text{dom}(s_1) \cap \text{dom}(s_2) &: s_1(\text{var}) = (\text{val}_1, t_1), \\
s_2(\text{var}) &= (\text{val}_2, t_2), t_1 \neq t_2
\end{aligned}$$

Definition 2 (Determinism on multithreaded state)

$$\begin{aligned}
D &: [\wp(\Psi) \rightarrow \{\text{true}, \text{false}\}] \\
D(\theta) &= \text{false} \\
&\Downarrow \\
\exists f_1, f_2 \in \theta &: \exists t \in \text{dom}(f_1) \cap \text{dom}(f_2) : \\
\tau_1 = f_1(t), \tau_2 = f_2(t), \exists i \in [0.. \min(\text{len}(\tau_1), \text{len}(\tau_2))] &: \\
DS(\text{istate}(\tau_1, i), \text{istate}(\tau_2, i)) &= \text{false}
\end{aligned}$$

5.3. First level of abstraction

$$\begin{aligned}
DS^\# &: [S^\# \rightarrow \{\text{true}, \text{false}\}] \\
DS^\#(s^\#) &= \text{false} \\
&\Downarrow \\
\exists \text{var} \in \text{dom}(s^\#) &: |\text{dom}(s^\#(\text{var}))| > 1 \\
\\
D^\# &: [\Psi^\# \rightarrow \{\text{true}, \text{false}\}] \\
D^\#(f^\#) &= \text{false} \\
&\Downarrow \\
\exists t \in \text{dom}(f^\#) &: \tau^\# = f^\#(t), \exists i \in [0..len(\tau^\#)] : \\
DS^\#(\text{istate}(\tau^\#, i)) &= \text{false}
\end{aligned}$$

Lemma 3 (Soundness) $D^\#$ is sound w.r.t. D , i.e.:

$$\forall \theta \in \wp(\Psi) : D(\theta) = \text{false} \Rightarrow D^\#(\alpha_\Psi(\theta)) = \text{false}$$

5.4. Second level of abstraction

$$\begin{aligned}
DS^\circ &: [S^\circ \rightarrow \{\text{true}, \text{false}\}] \\
DS^\circ(s^\circ) &= \text{false} \Leftrightarrow \exists \text{var} \in \text{dom}(s^\circ) : |\pi_2(s^\circ(\text{var}))| > 1 \\
\\
D^\circ &: [\Psi^\circ \rightarrow \{\text{true}, \text{false}\}] \\
D^\circ(f^\circ) &= \text{false} \\
&\Downarrow \\
\exists t \in \text{dom}(f^\circ) &: \tau^\circ = f^\circ(t), \exists i \in [0..len(\tau^\circ)] : \\
DS^\circ(\text{istate}(\tau^\circ, i)) &= \text{false}
\end{aligned}$$

Lemma 4 (Soundness) D° is sound w.r.t. $D^\#$, i.e.:

$$\forall f^\# \in \Psi^\# : D^\#(f^\#) = \text{false} \Rightarrow D^\circ(\alpha_{\Psi^\#}(f^\#)) = \text{false}$$

5.5. The running example

The deterministic property on the concrete semantics discovers that the running example introduced by Section 1.1 does not respect it. In particular, the two executions differ on the values read by both the threads: in the first one, Thread 1 reads a value written by System and Thread 2 the one written by Thread 2, while in the second execution Thread 1 reads the value of Thread 2 and Thread 2 the one of System. In Figure 5 we underline the identifiers of threads in the cases in which a non-deterministic behavior arises.

Since the two level of abstractions are sound w.r.t. determinism, they do not validate the deterministic property.

6. Weak determinism

In this section we introduce a new concept of determinism which is weaker than the previous definitions and that is defined on the first level of abstraction.

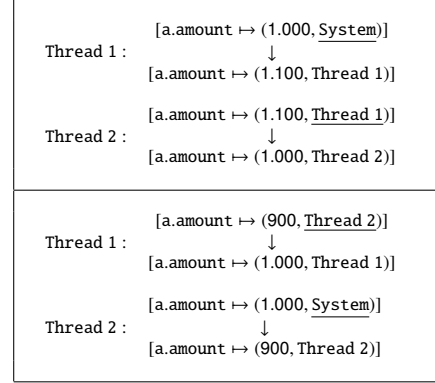


Figure 5. The non-deterministic behaviors in the concrete semantics

6.1. Approximating numerical values

The first level of abstraction is parameterized on a non-relational abstract domain that approximates the numerical values. Moreover, we collect for each thread an abstract value that approximates all the concrete values it may have written at that point. Through these two observations we are in position to define a new property: the weak determinism. The idea is that two different concrete values do not produce different observable behaviors if their abstraction is the same. This concept has to be tuned on an abstract domain: for instance, with the sign domain it would mean that at a given point all the values written in parallel on a given variable of the shared memory have the same sign.

6.2. Formal definition

Definition 3 (Weak determinism on shared memory)

$$\begin{aligned}
ADS^\# &: [S^\# \rightarrow \{\text{true}, \text{false}\}] \\
ADS^\#(s^\#) &= \text{false} \\
&\Downarrow \\
\exists \text{var} \in \text{dom}(s^\#) &: |\text{dom}(s^\#(\text{var}))| > 1 \wedge \\
\exists t_1, t_2 \in \text{dom}(s^\#(\text{var})) &: s^\#(\text{var})(t_1) \neq s^\#(\text{var})(t_2)
\end{aligned}$$

Definition 4 (Weak determinism on multithreaded state)

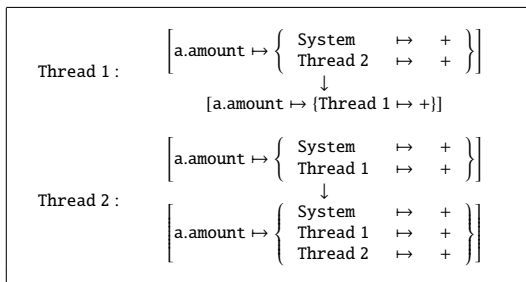
$$\begin{aligned}
AD^\# &: [\Psi^\# \rightarrow \{\text{true}, \text{false}\}] \\
AD^\#(f^\#) &= \text{false} \\
&\Downarrow \\
\exists t \in \text{dom}(f^\#) &: \tau^\# = f^\#(t), \exists i \in [0..len(\tau^\#)] : \\
ADS^\#(\text{istate}(\tau^\#, i)) &= \text{false}
\end{aligned}$$

Note that, even if a program is not deterministic in the concrete, the weak determinism may be validated it, as the values written by different threads may produce the same abstract element.

6.3. Example 2

```
public void withdraw2(int money) {
    synchronized(this) {
        int temp=account.amount-money;
        if(temp>0)
            account.amount=temp;
    }
}
```

(a) The source code



(b) The abstract semantics

Figure 6. Example 2

Suppose we now modify the `withdraw` method presented in section 1.1 as depicted in Figure 6.(a).

This method allows the withdrawing only if there is enough money in the bank account. Suppose we analyze the multi-threaded program presented by the running example when replacing the `withdraw` method with this one. We use the sign domain to capture numerical information; in this way we capture for each numerical value if it is positive (+), negative (−) or equal to zero (0). The result of the analysis at the first level of abstraction is depicted by Figure 6.(b); we simplify the elements depicting just the results of the `withdraw` method, ignoring the internal computation of the thread.

Applying the weak deterministic property on that abstract value we obtain that such a property is validated: in fact in all the possible executions the value stored by the field `amount` is always positive. On the other side, the full determinism is not guaranteed, as the amount of money may have been written in parallel by different threads.

7. From determinism to semi-automatic parallelization

In this section we sketch how the determinism of a program may be useful in order to semi-automatically parallelize it.

7.1. Motivation

As already explained by the introduction, multi-core architectures have recently appeared in a broad market, and

the only way to take advantage from this technology is to develop applications that perform as many tasks in parallel as possible. On the other hand, reasoning about parallel applications is strictly more difficult than on sequential programs. The logical consequence is that (semi-) automatic tools in order to find and detect possible parallelizable fragments of sequential code are particularly useful.

7.2. Determinism and parallelism

By definition, a sequential program performs a set of sequential operations. Usually, an operation reads some data written or modified by previous statements and it writes values that will be used during the rest of the computation. When each sequential operation deals with a disjointed set of variables, the program is trivially parallelizable. However, the latter condition does not apply to the most part of sequential programs.

Through our deterministic property, given two subsets of a sequential program, we can analyze them as if they were executed in parallel, and check if there are some non-deterministic behaviors due to the parallel execution. If it is not the case, the two parts of the sequential program can run in parallel without inducing any new behavior because of the parallel execution.

7.3. Relaxing the deterministic property

In order to find parallelizable blocks in a program, we need also to relax the deterministic property in order to focus only on the critical parts of the program. A first approach consists in restricting the check of the deterministic property only: (i) on the traces of some threads (e.g. only the ones that performs critical operations); (ii) on a subset of the shared memory’s variables (e.g. only the amount of money of a bank account, and not the name of the owner); (iii) on a subset of the states (e.g. only the ones that write on the shared memory). A priori, we cannot know which components may be ignored, and this should be an input by the developer, as he is the only one that knows which part are critical and which are not.

Another approach is to use the weak determinism presented by Section 6.

These two approaches may be combined together. For instance, we can express something like “check the weak determinism through the sign domain on the amount of money”.

8. Experimental results

The analyzer presented by [6] (that works on Java multi-threaded programs at bytecode level) has been extended in order to study the deterministic property; in addition, it has

Program	# St.	# Th.	An.	Det.	Weak Det.
philos	213	2	< 1"	< 1"	< 1"
tsp	1936	2	24"	< 1"	< 1"
elevator	1829	2	15"	< 1"	< 1"
barrier	363	3	< 1"	< 1"	< 1"
forkjoin	170	2	< 1"	< 1"	< 1"
sync	320	3	< 1"	< 1"	< 1"
sor	1121	2	4"	< 1"	< 1"
crypt	2636	3	4"	< 1"	< 1"
lufact	3732	2	26"	< 1"	< 1"
montecarlo	3864	2	46"	4"	4"
moldyn	9029	2	13'51"	50"	1'01"

Table 1. Experimental results

already been extended also to be sound w.r.t. the happens-before memory model, following the theoretical results of [7]. We tested the analysis of deterministic and weak deterministic properties on a set of well-known benchmarks: `philos`, `tsp`, and `elevator` are taken from [19], while the others (`barrier`, `forkjoin`, `sync`, `sor`, `crypt`, `lufact`, `montecarlo`, and `moldyn`) are taken from the Java Grande Forum Benchmark Suite [1]. All tests have been performed on an AMD Athlon 64 3000+ processor with 2 GB of RAM running a GNU/Linux operating system with the Java virtual machine version 1.6.0_06-b02.

Table 1 depicts the experimental results. The first column reports the name of the program, while the following two columns report respectively the number of bytecode statements and of abstract threads analyzed. The fourth column depicts the time required in order to build up an abstraction of the program following the approach introduced by [7], while the last two columns report respectively the time spent to analyze the deterministic and the weak-deterministic property on the results obtained by this abstraction.

All the programs contain a potentially unbounded number of concrete threads; these are abstracted into a finite set through an ad-hoc alias analysis (note that in Java threads are objects and so identified by reference) [6]. The time required in order to validate the properties is much less than the one required by the analysis in order to build up an abstraction of the program; in average, it requires the 6% of the overall computational time. In addition, in the worst case in order to analyze almost 10k of Java bytecode statements the overall analysis require about 15'. We think that these experimental results are particularly encouraging, considering also that the analysis has not been optimized at all, and we think we may easily obtain better computational times.

9. Related works

During the last decades, a topic largely debated by researchers has been how parallel threads can communicate in order to bound the nondeterministic behaviors. Many different approaches have been proposed, but none has been able to solve this problem.

A first way has been the static or dynamic checking on the actions performed by different threads, like general and data races [16]. The idea is that races cause nondeterminism, and then they are symptom of bugs. On one hand, the absence of general races is a too restrictive condition for multithreaded programs that communicate asynchronously through shared memory. On the other hand, avoiding data races requires sometimes unnecessary synchronization actions, and it does not guarantee at all the determinism of a program. In addition, a race does not take into account which statements and threads cause it, with which area of the shared memory it deals, if the informations written in parallel are different and how much. Our intuition is to directly study nondeterminism, and in this way we can relax the property on a critical subset of the shared memory, considering only some statements and threads, looking also at the (abstract) values written and read in parallel. In this way, our approach is strictly more flexible than general and data races.

Another big amount of work has been focused on the consistency conditions that multithreaded programs have to provide. The most known model is sequential consistency [12], that has been proved to be too restrictive for current programming language. For instance, the Java memory model [14] is more relaxed than it, and more restrictive than the happens-before memory model [11]. Our approach is orthogonal w.r.t. these consistency conditions, as they define which states of shared memory are visible during the execution of different threads. In addition a previous work [7] has already defined a static analysis sound w.r.t. the happens-before memory model, that is an over-approximation of the Java one.

A similar approach has been the definition of a specific model in order to reduce the non-determinism of multithreaded programs. The Software Transactional Memory (STM) [17, 9] allows to specify that the execution of a given part of program is atomic, i.e. it is viewed as a unique operation by other threads. A large work on the semantics of parallel languages has been developed by trace theory [15], as [8] that presents a deterministic semantics for a simple parallel programming language. In this case, our idea is to distinguish the model of execution and the determinism of a program, as they are two orthogonal issues. Our deterministic and weak deterministic prop-

erties may be applied to analyses sound w.r.t. sequential consistency, happens-before memory model, Java memory model, software transactional memory, etc etc..

10. Conclusion and future work

In this paper, we show how to define and statically analyze the determinism of multithreaded programs. The main novelty w.r.t. previous approaches is that our work deals directly with the determinism, and in particular with the behaviors due to the multithreaded executions, instead of, for instance, defining a restricted model. In this way we are able to define and statically analyze which programs are deterministic w.r.t. their multithreaded executions and which are not. In addition, since we are dealing at the source of the problem, our approach may lead to new and unexplored ideas, as the weak determinism and the semi-automatic parallelization of sequential programs. The soundness of our analysis has been formally proved, and we presented some experimental results when applying the analysis to a set of well-known benchmarks.

As future works, we want to investigate how to apply this property in order to propose to the developers possible parallelizations of sequential programs. In particular, we may rely on the information obtained through the cost analysis [2] in order to detect which sequential statements are computationally heavier, and to check if their parallelization induces non-deterministic behaviors.

References

- [1] Java Grande Forum Benchmark Suite. At <http://www.epcc.ed.ac.uk/research/activities/java-grande/>.
- [2] E. Albert, P. Arenas, S. Genaim, G. Puebla, and D. Zardini. Cost analysis of java bytecode. In Springer, editor, *Proceedings of ESOP '07*, LNCS, 2007.
- [3] P. Cousot and R. Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In ACM Press, editor, *Proceedings of POPL '77*, 1977.
- [4] P. Cousot and R. Cousot. Systematic design of program analysis frameworks. In ACM Press, editor, *Proceedings of POPL '79*, 1979.
- [5] Perry A. Emrath and David A. Padua. Automatic detection of nondeterminacy in parallel programs. In ACM Press, editor, *Proceedings of PADD '88*, 1988.
- [6] P. Ferrara. A fast and precise analysis for data race detection. In Elsevier, editor, *Proceedings of Bytecode '08*, volume ENTCS, 2008.
- [7] P. Ferrara. Static analysis via abstract interpretation of the happens-before memory model. In Springer, editor, *Proceedings of TAP '08*, LNCS, 2008.
- [8] Paul Gastin and Michael W. Mislove. A truly concurrent semantics for a simple parallel programming language. In Springer, editor, *Proceedings of CSL '99*, 1999.
- [9] Tim Harris, Simon Marlow, Simon Peyton-Jones, and Maurice Herlihy. Composable memory transactions. In ACM Press, editor, *Proceedings of PPOPP '05*, 2005.
- [10] G. Koch. Discovering multi-core: extending the benefits of Moore's law. In *Technology Intel Magazine*. Intel, July 2005.
- [11] L. Lamport. Time, clocks, and the ordering of events in a distributed system. In ACM Press, editor, *Commun. ACM*, 1978.
- [12] L. Lamport. How to make a multiprocessor computer that correctly executes multiprocess programs. In *IEEE Trans. Computers*, 1979.
- [13] Edward A. Lee. The problem with threads. In IEEE Computer Society Press, editor, *Computer*, 2006.
- [14] J. Manson, W. Pugh, and S. V. Adve. The Java memory model. In ACM Press, editor, *Proceedings of POPL '05*, 2005.
- [15] A Mazurkiewicz. Trace theory. In Springer, editor, *Advances in Petri nets 1986*, 1987.
- [16] R. H. B. Netzer and B. P. Miller. What are race conditions?: Some issues and formalizations. In ACM Press, editor, *ACM Lett. Program. Lang. Syst.*, 1992.
- [17] Nir Shavit and Dan Touitou. Software transactional memory. In ACM Press, editor, *Symposium on Principles of Distributed Computing*, 1995.
- [18] Herb Sutter and James Larus. Software and the concurrency revolution. In ACM Press, editor, *ACM Queue*, 2005.
- [19] C. Von Praun and T. R. Gross. Object race detection. In ACM Press, editor, *Proceedings of OOPSLA 01*, 2001.