

# Abstract interpretation of memory models

Pietro Ferrara

ETH Zurich  
Switzerland

IFIP WG 2.3 Programming Methodology, Lachen, Switzerland


# Multithreading

*“(...) in order for an application to take advantage of the dual-core capabilities, the application should be optimized for multithreading.”*

G. Koch. Discovering multi-core: extending the benefits of Moore's law. In Technology Intel Magazine. Intel, July 2005.

- Parallelism supported through **multithreading**
  - > Java and C#
- Implicit communications via **shared memory**
- Synchronization on **monitors**
- **Subtle** and problematic

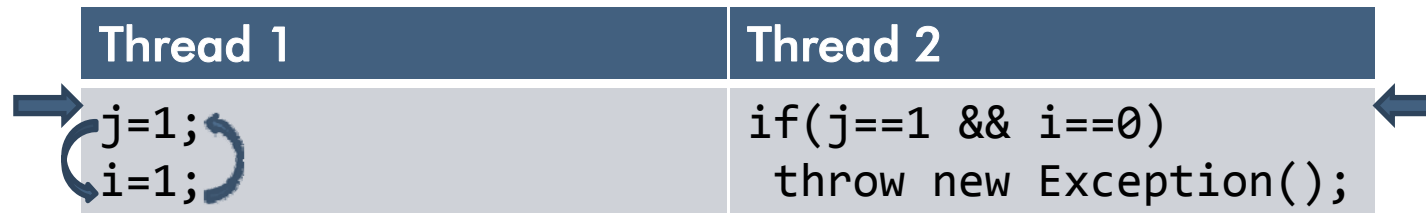
# An example

Thread 1	Thread 2
 <code>i=1;</code> <code>j=1;</code>	<code>if(j==1 &amp;&amp; i==0)</code> <code>throw new Exception();</code>

- At the beginning,  $i=0$ ,  $j=0$
- May the **exception** be thrown?
  - Sequential consistency: **no!**

Variable	Value
<code>i</code>	1
<code>j</code>	1

# An example



- At the beginning
- May in  
• Sequential
- Java: **yes!**
- E.g. swap of independent statements

The exception is thrown

# Memory models



- Describe how programs **interact** with memory
  - > Also in presence of data races
- Interesting only for multithreading
  - > Single thread: *“as if sequential”* semantics
- Define which behaviors are allowed
  - > Which parallel writes a read action can see
- Two **opposite** needs
  - > Make sure programs run quickly (as much as possible)
    - **Allow** optimizations
  - > Programmers understand programming languages
    - **Restrict** optimizations

# Memory models

- **Sequential consistency**
  - > Actions appear as executed exactly in the order of the original program
    - Too much restrictive (presence of data races)
- **Happens-before memory model (HBMM)**
  - > Happens-before order: transitive closure of
    - Program order
    - Synchronize-with order
- **Java memory model**
  - > Restriction of the HBMM
  - > Complex formalization

# Happens-before memory model

- A read  $r[v]$  can see a write  $w[v]$  if and only if
  - >  $r[v]$  does not happens before  $w[v]$
  - > There is no  $w'[v]$  such that
    - $w[v]$  happens before  $w'[v]$
    - $w'[v]$  happens before  $r[v]$

Thread 1	Thread 2
 <code>i=0;</code>	<code>int v=i;</code>
 <code>i=1;</code>	
<code>new Thread2().start();</code>	

The diagram shows two threads. Thread 1 has three lines of code: `i=0;`, `i=1;`, and `new Thread2().start();`. Thread 2 has one line of code: `int v=i;`. A red 'X' is next to `i=0;` and a green checkmark is next to `i=1;`. Two blue arrows point from the `int v=i;` line in Thread 2 to the `i=0;` and `i=1;` lines in Thread 1, indicating that Thread 2's read of `i` happens before both of Thread 1's writes to `i`.

- `i=0` happens before `i=1`
- `i=1` happens before `int v=i`

# Sound static analysis

- Infer and prove properties
  - > at **compile time** (static)
  - > respected by **all** possible executions (sound)
- Abstract interpretation
  - > Mathematical theory developed by P. & R. Cousot
  - > Define concrete semantics of programs
  - > Soundly **approximate** it
- Common approach:
  - > Semantics defined in a **fixpoint** form

# Contribution

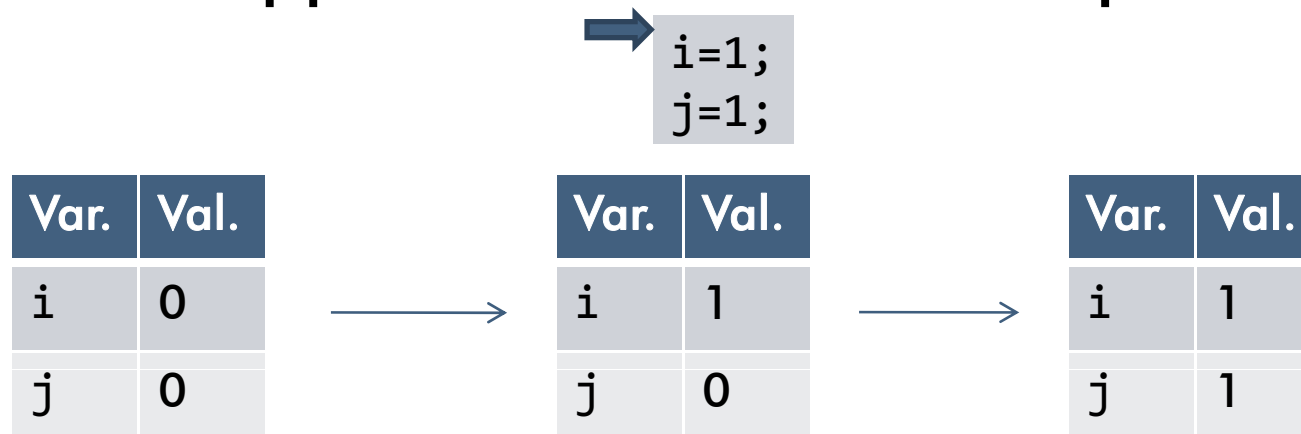
- Static analysis **sound** with respect to **MM**
  - > Based on abstract interpretation
  - > Instantiated on the **HBMM**
- Concrete semantics: formalization of the **MM**
  - > Trace semantics
  - > In fixpoint form
- Abstract semantics: **generic** with respect to
  - > Property of interest
  - > Numerical domain
- **Implemented**

# Outline

1. Introduction
2. Concrete semantics
3. Abstract semantics
4. Experimental results
5. Conclusion and future work

# Formalizing executions

- Trace semantics:
  - > Execution represented as an **ordered trace of states**
- Common approach in abstract interpretation



- It relies on the *“as if sequential”* semantics
- Set of states:  $\Sigma$
- Finite traces of states:  $\Sigma^{\vec{+}}$

# Multithread and multicore

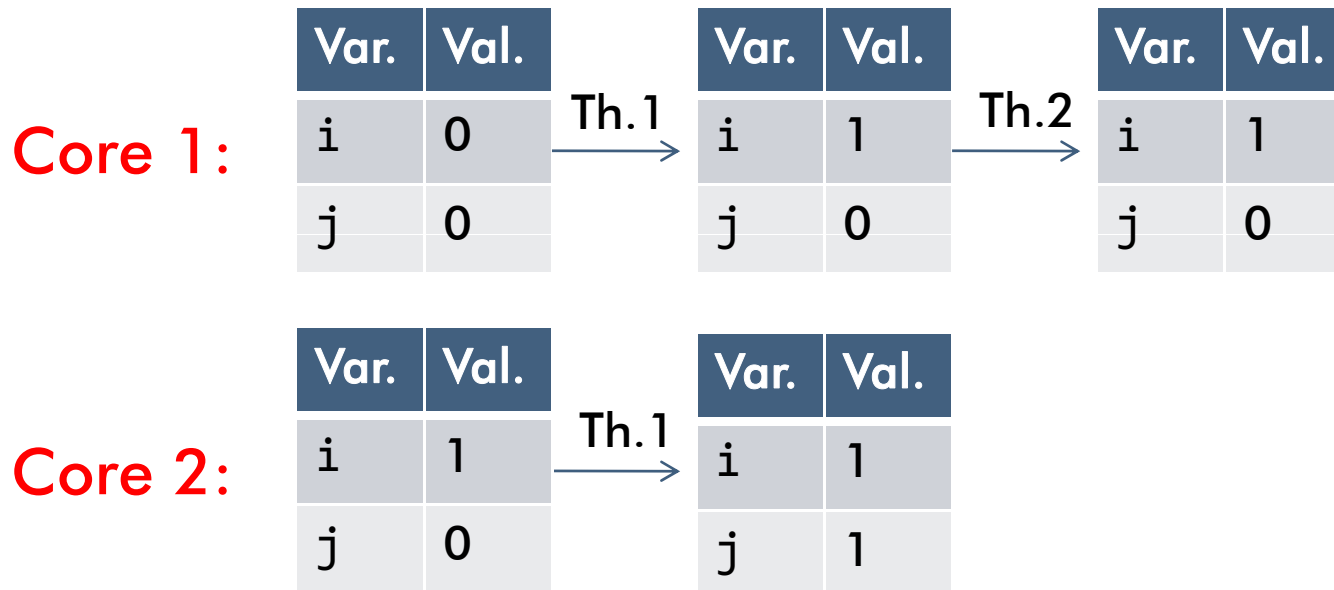
- Problem with multithreading:
  - > Several traces of execution
- Several architectures: single, dual, ... cores CPU
- Example: single core

Thread 1	Thread 2
<code>i=1;</code> <code>j=1;</code>	<code>if(j==1 &amp;&amp; i==0)</code> <code>throw new Exception();</code>

Var.	Val.		Var.	Val.		Var.	Val.		Var.	Val.
i	0	→ Th.1	i	1	→ Th.1	i	1	→ Th.2	i	1
j	0		j	0		j	1		j	1

# Dual core

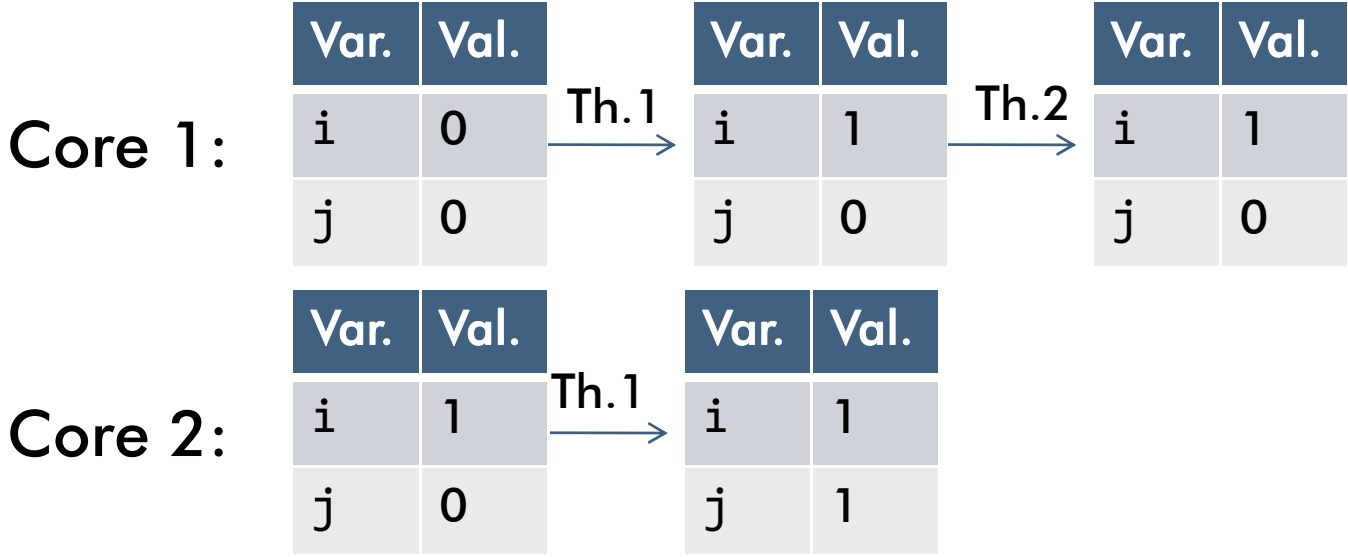
	Thread 1	Thread 2
C.1		
C.2		
	<code>i=1;</code> <code>j=1;</code>	<code>if(j==1 &amp;&amp; i==0)</code> <code>throw new Exception();</code>



# Thread partitioning

- **Main issue:**
  - > Different ways of executing a multithreaded program
  - > Execution relies on the **hardware** architecture
- **Multithreaded model**
  - > Several **threads** are executed in parallel
  - > *“As if sequential”* semantics for each thread locally
  - > Shared memory
- **Thread partitioning domain**
  - > It collects the trace of execution **for each thread**
  - > Abstract away the interleaving of **different threads**

# Thread partitioning domain



Thread 1:

Thread 2:

# Thread partitioning domain

- We collect for each thread its trace of execution

$$\Psi : [\text{Tid} \rightarrow \Sigma^{\vec{+}}]$$

- What is the state of a single thread?

- > **Private** memory

- E.g. local variables

- > **Shared** memory

- E.g. the heap

- **Memory model**

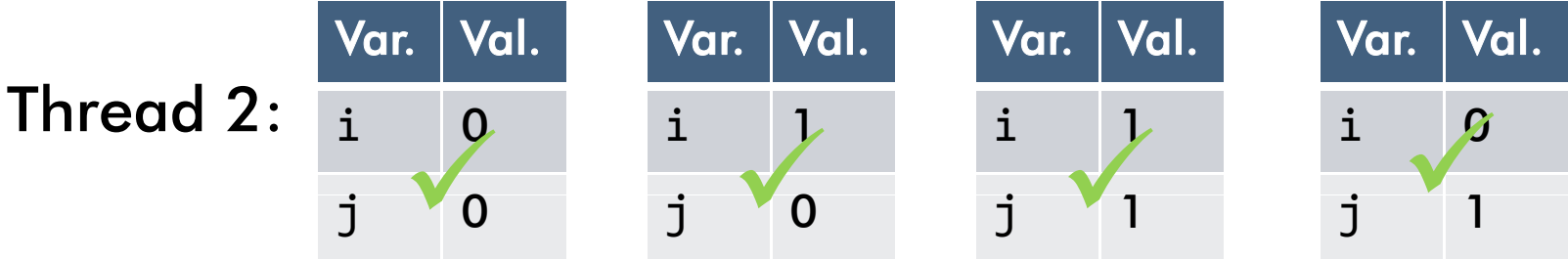
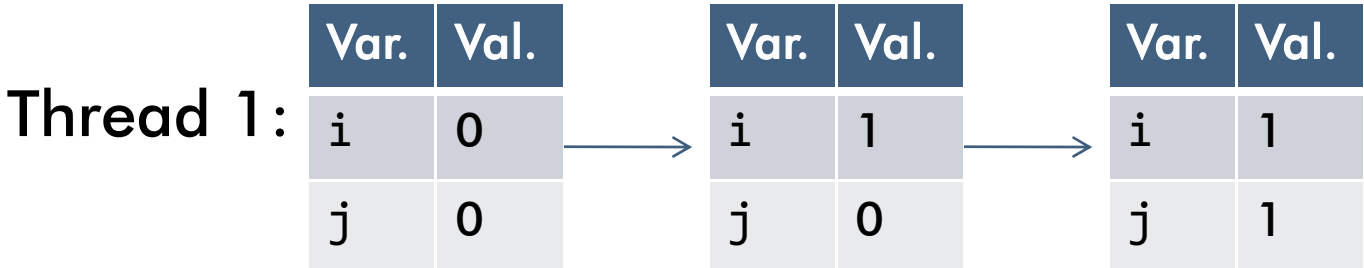
- > Define the possible **views** of the shared memory

- **Concrete semantics**

- > **Choose** one of the possible views of the memory

# Thread partitioning domain

Thread 1	Thread 2
<pre>i=1; j=1;</pre>	<pre>if(j==1 &amp;&amp; i==0)   throw new Exception();</pre>



HBMM defines the allowed views of shared memory

Pietro Ferrara: "Abstract interpretation of memory models"  
 IFIP WG 2.3 "Programming Methodology", Lachen, Switzerland

# Partial finite trace semantics

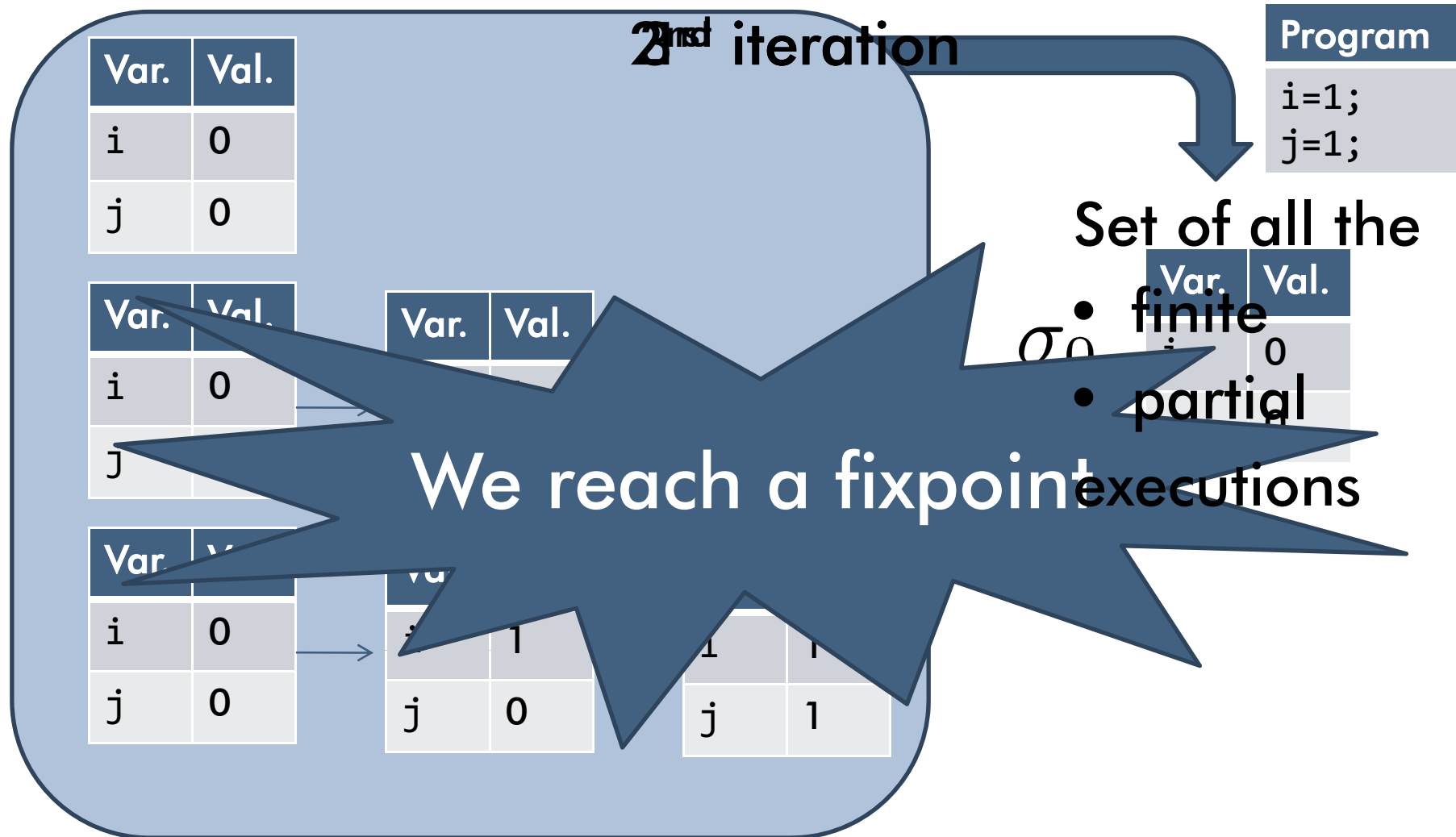
- Build up all the **partial finite executions** of a program
- **Fixpoint semantics** of thread  $t$

$$\mathcal{S}^\circ[t] = \text{lfp}_{\emptyset}^{\subseteq} \lambda T. \{ \sigma_0^t \} \cup \{ \sigma_0 \rightarrow \dots \rightarrow \sigma_n \rightarrow \sigma_{n+1} : \sigma_0 \rightarrow \dots \rightarrow \sigma_n \in T \wedge \sigma_{n+1} \in \mathcal{S}[\sigma_n] \}$$

where

- >  $\sigma_0^t$  is the initial state of thread  $t$
- >  $\mathcal{S}$  is the small step semantics of the programming language
- >  $\sigma_0 \rightarrow \dots \rightarrow \sigma_n$  is a trace composed by  $n$  states

# Partial finite trace semantics



# Something is missing

- Definition of single thread semantics

$$S^\circ[t] = \text{lfp}_{\emptyset}^{\subseteq} \lambda T. \{ \sigma_0^t \} \cup \{ \sigma_0 \rightarrow \dots \rightarrow \sigma_n \rightarrow \sigma_{n+1} : \\ \sigma_0 \rightarrow \dots \rightarrow \sigma_n \in T \wedge \sigma_{n+1} \in S[\sigma_n] \}$$

where

$\sigma_0$  is the initial state

$S$  is

$S[\sigma]$

It should take into  
account the MM

memory states

# Small step semantics

- We suppose to have the semantics of statements
  - > *"As if sequential"* semantics
  - > Atomic at multithreaded level
- If we have a thread local action
  - > Follow the semantics of statements
- If we read something from the shared memory
  - > MM provides all possible views of the shared memory
- Non deterministic choice
  - > Different threads' interleaving may cause to read different values from the shared memory

# Single-thread semantics

$$S^\circ \llbracket f, t \rrbracket = \text{lfp}_{\emptyset}^{\subseteq} \lambda T. \{ \sigma_0^t \} \cup \{ \sigma_0 \rightarrow \dots \rightarrow \sigma_n \rightarrow \sigma_{n+1} : \sigma_0 \rightarrow \dots \rightarrow \sigma_n \in T \wedge \sigma_{n+1} \in S^\circ \llbracket f, \sigma_n \rrbracket \}$$

The semantics extracts  
the values written in  
parallel by other threads

# Multithread semantics

- $S^{\circ}$  exposes values on the shared memory
  - > These values may be read by other threads
  - > They may cause new behaviors
    - Exposing eventually new values on the shared memory
- Multithread semantics
  - > Iterate the single thread semantics on each thread
  - > Each iteration may produce new values on the heap
  - > Stop when a fixpoint is reached
- Not only values on the shared memory
  - > New executions may create and launch new threads
    - Dynamic allocation of threads

# Multithread semantics

$$\begin{aligned} S^{\parallel} [f_0] &= \text{lfp}_{\emptyset}^{\subseteq} \lambda \Phi. \{f_0\} \cup \\ &\quad \{f_i : \exists f_{i-1} \in \Phi : \forall t \in \text{dom}(f_{i-1}) : \\ &\quad f_i(t) \in S^{\circ} [f_{i-1}, t] \cap \Sigma_{S^{\circ}}^{\vec{+}}\} \end{aligned}$$

where

- $\Sigma_{S^{\circ}}^{\vec{+}}$  contains the traces ending with a final state
- $f_0$  is the function relating
  - > each active thread at the beginning of the execution
  - > to a trace containing its initial state

# Multithread semantics in action

Thread 1	Thread 2
<pre>i=1; j=1;</pre>	<pre>if(j==1 &amp;&amp; i==0)   throw new Exception();</pre>

Initial state

Th.1:	<table><tr><td>i</td><td>0</td></tr><tr><td>j</td><td>0</td></tr></table>	i	0	j	0
i	0				
j	0				
Th.2:	<table><tr><td>i</td><td>0</td></tr><tr><td>j</td><td>0</td></tr></table>	i	0	j	0
i	0				
j	0				

1<sup>st</sup> iteration

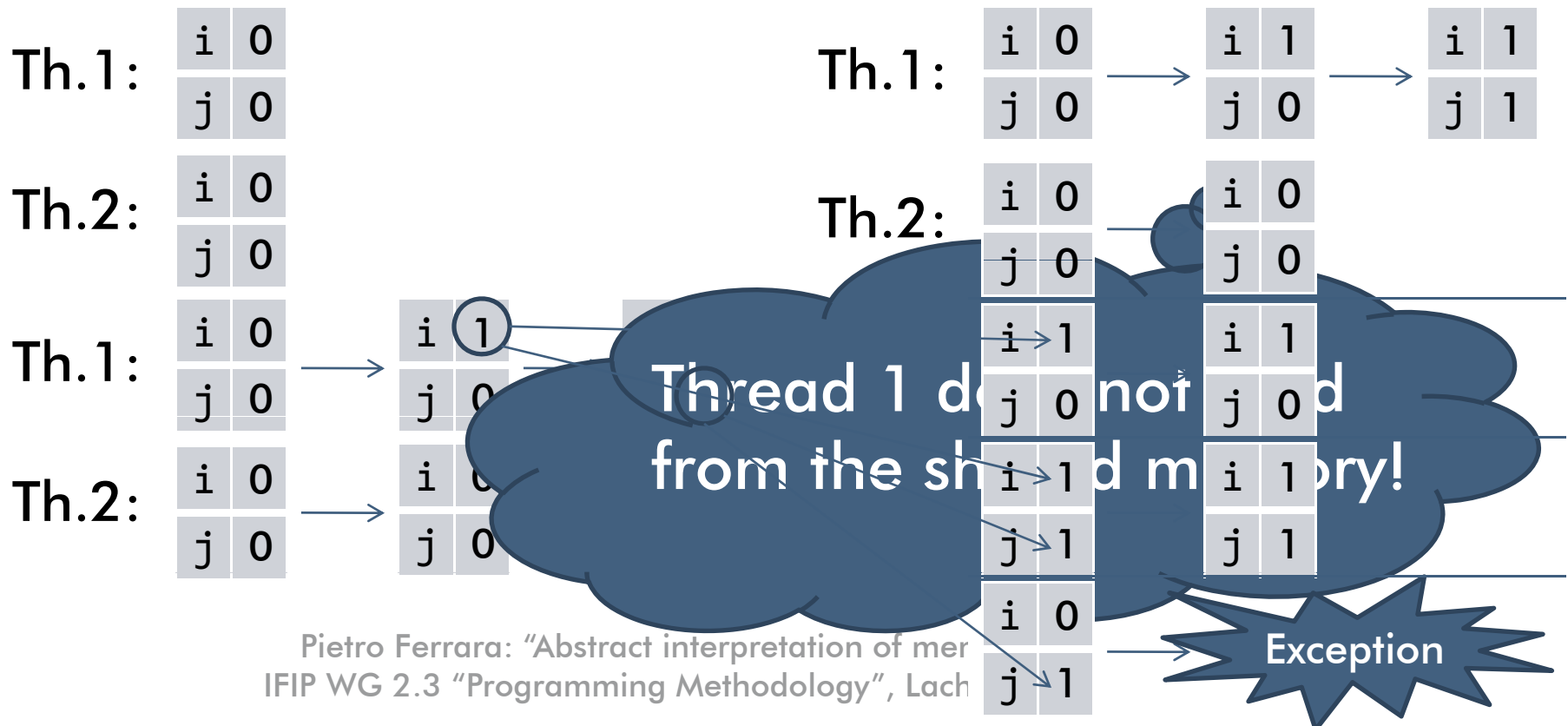
Th.1:	<table><tr><td>i</td><td>0</td></tr><tr><td>j</td><td>0</td></tr></table>	i	0	j	0	→	<table><tr><td>i</td><td>1</td></tr><tr><td>j</td><td>0</td></tr></table>	i	1	j	0	→	<table><tr><td>i</td><td>1</td></tr><tr><td>j</td><td>1</td></tr></table>	i	1	j	1
i	0																
j	0																
i	1																
j	0																
i	1																
j	1																
Th.2:	<table><tr><td>i</td><td>0</td></tr><tr><td>j</td><td>0</td></tr></table>	i	0	j	0	→	<table><tr><td>i</td><td>0</td></tr><tr><td>j</td><td>0</td></tr></table>	i	0	j	0						
i	0																
j	0																
i	0																
j	0																

# Multithread semantics in action

Thread 1	Thread 2
<code>i=1;</code> <code>j=1;</code>	<code>if(j==1 &amp;&amp; i==0)</code> <code>throw new Exception();</code>

Initial state  $\cup$  1<sup>st</sup> iteration

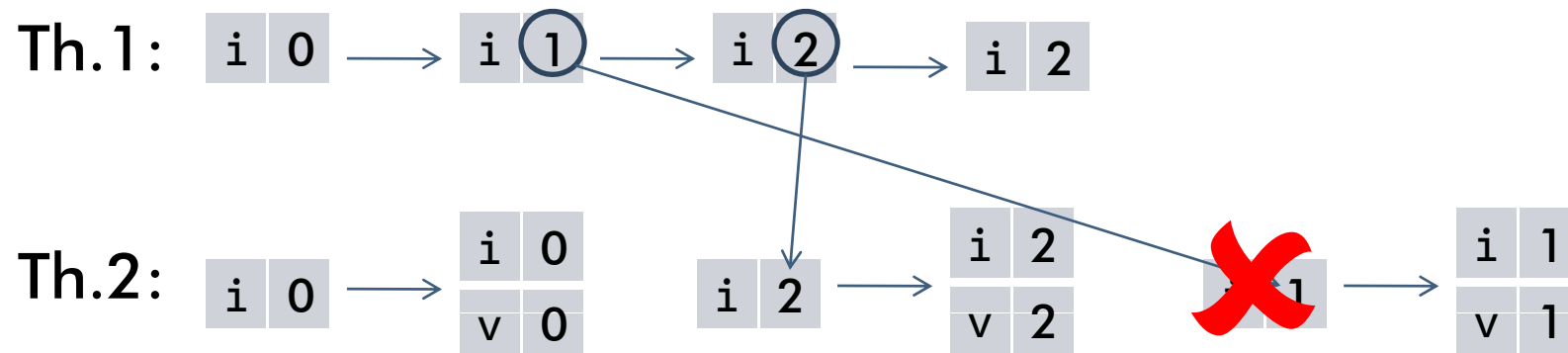
2<sup>nd</sup> iteration



# Synchronization patterns

- Our formalization supports
  - > Monitors

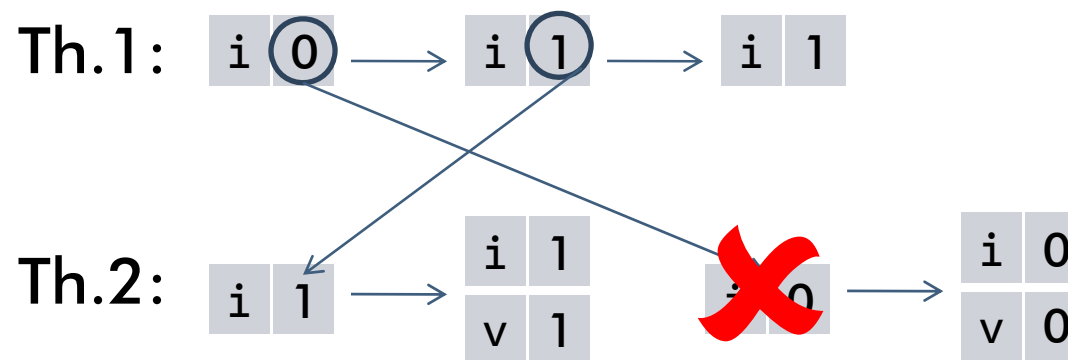
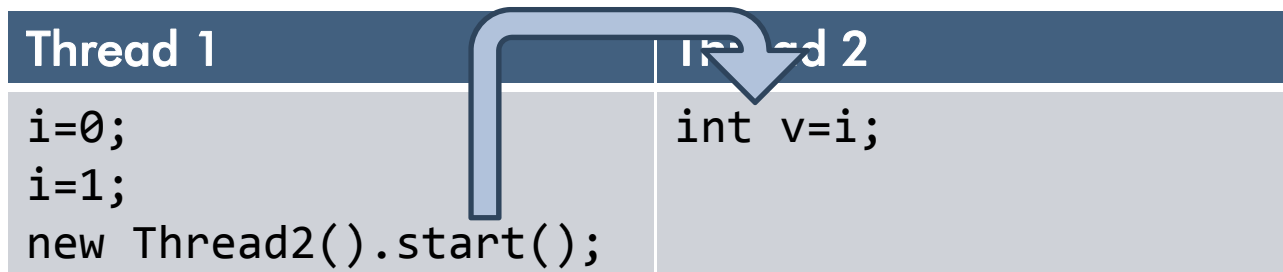
Thread 1	Thread 2
<code>lock(o); i=1; i=2; unlock(o);</code>	<code>lock(o); int v=i; unlock(o);</code>



Pietro Ferrara: "Abstract interpretation of memory models"  
IFIP WG 2.3 "Programming Methodology", Lachen, Switzerland

# Synchronization patterns

- Our formalization supports
  - > Dynamic launch of threads



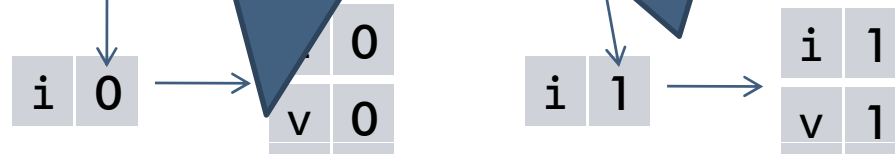
# Synchronization patterns

- Our formalization does not track information on
  - > Join of threads
  - > Lazy-vous synchronization (wait and notify)

The semantics produces also this execution

Th.1:

Th.2:



# Soundness

- **Final goal:**
  - > Formalize a **static analysis**
  - > Sound with respect to all the behaviors of a **MM**
- **We do not consider some synchronization actions**
  - > Then we allow **more behaviors**
- **This does not break the soundness**
  - > ... but it will produce **less precise** results
- **Future work**
  - > Track information on other synchronization patterns

# Outline

1. Introduction
2. Concrete semantics
3. Abstract semantics
4. Experimental results
5. Conclusion and future work

# Abstract semantics

- **Concrete semantics**
  - > Not computable
- **Abstract semantics**
  - > Approximation of the concrete semantics
  - > Computable
  - > Soundness proved formally
- **One abstract trace approximates all concrete traces**
  - > Common approach in abstract interpretation
  - > One abstract state for each statement in the program

# Abstract domain

$$\bar{\Psi} : [\bar{\text{Tid}} \rightarrow \bar{\Sigma}^{\vec{+}}]$$

- What is an abstract state?
  - > We do not define it in the details
  - > We interact with it using some **functions**
- **Generic analysis**
  - > E.g., it can be instantiated with different numerical domains
    - Intervals, signs, ...
    - Up to know, only non-relational numerical domains

# Single thread abstract semantics

$$\overline{S}^\circ[[\overline{f}, \overline{t}]] = \text{lfp}_{\emptyset}^{\sqsubseteq^{\vec{\tau}}} \lambda \overline{\tau}. \{ \overline{\sigma}_0^{\overline{t}} \} \sqcup^{\vec{\tau}} \{ \overline{\sigma}_0 \rightarrow \dots \rightarrow \overline{\sigma}_n \rightarrow \overline{\sigma}_{n+1} : \\ \overline{\tau} = \overline{\sigma}_0 \rightarrow \dots \rightarrow \overline{\sigma}_n \wedge \overline{\sigma}_{n+1} = \overline{S}[[\overline{f}, \overline{\sigma}_n]] \}$$

where

- $\sqcup^{\vec{\tau}}$  and  $\sqsubseteq^{\vec{\tau}}$  are respectively the applications of the upper bound and the partial order operator of states to all the elements of the traces

$$S^\circ[[f, t]] = \text{lfp}_{\emptyset}^{\subseteq} \lambda T. \{ \sigma_0^t \} \cup \{ \sigma_0 \rightarrow \dots \rightarrow \sigma_n \rightarrow \sigma_{n+1} : \\ \sigma_0 \rightarrow \dots \rightarrow \sigma_n \in T \wedge \sigma_{n+1} \in S[[f, \sigma_n]] \}$$

# Multithread semantics

$$\overline{S}^{\parallel} [\overline{f}_0] = \text{lfp}_{\emptyset}^{\sqsubseteq^{\Phi}} \lambda \overline{f}. \{ \overline{f}_0 \} \sqcup^{\Phi} \{ \overline{f}_i : \forall \overline{t} \in \text{dom}(\overline{f}) : \\ \overline{f}_i(\overline{t}) = \overline{S}^{\circ} [\overline{f}, \overline{t}] \sqcap^{\Phi} \overline{\Sigma}_{\overline{S}^{\circ}}^{\vec{\dagger}} \}$$

where  $\sqcup^{\Phi}, \sqcap^{\Phi}$  and  $\sqsubseteq^{\Phi}$  are respectively the applications of  $\sqcup^{\vec{\dagger}}, \sqcap^{\vec{\dagger}}$  and  $\sqsubseteq^{\vec{\dagger}}$  to the traces of all the threads

$$\overline{S}^{\parallel} [\overline{f}_0] = \text{lfp}_{\emptyset}^{\subseteq} \lambda \Phi. \{ \overline{f}_0 \} \cup \\ \{ \overline{f}_i : \exists \overline{f}_{i-1} \in \Phi : \forall \overline{t} \in \text{dom}(\overline{f}_{i-1}) : \\ \overline{f}_i(\overline{t}) \in \overline{S}^{\circ} [\overline{f}_{i-1}, \overline{t}] \cap \overline{\Sigma}_{\overline{S}^{\circ}}^{\vec{\dagger}} \}$$

# Abstract semantics in action

Thread 1	Thread 2
<code>i=1;</code> <code>j=1;</code>	<code>if(j==1 &amp;&amp; i==0)</code> <code>throw new Exception();</code>

Initial state

Th.1: 

i	[0..0]
j	[0..0]

Th.2: 

i	[0..0]
j	[0..0]

1<sup>st</sup> iteration

Th.1: 

i	[0..0]	→	i	[1..1]	→	i	[1..1]
j	[0..0]		j	[0..0]		j	[1..1]

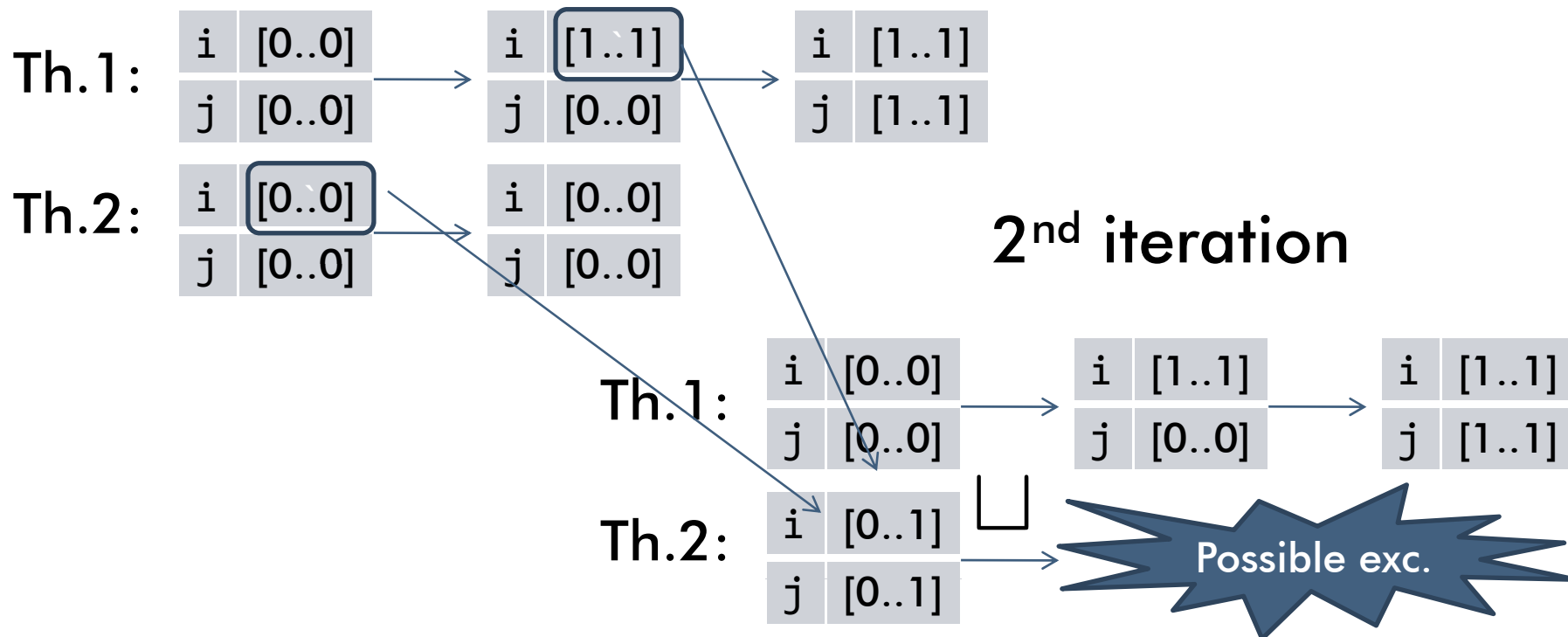
Th.2: 

i	[0..0]	→	i	[0..0]
j	[0..0]		j	[0..0]

# Multithread semantics in action

Thread 1	Thread 2
<code>i=1;</code> <code>j=1;</code>	<code>if(j==1 &amp;&amp; i==0)</code> <code>throw new Exception();</code>

Initial state  $\sqcup$  1<sup>st</sup> iteration



# Abstract threads

- We support **dynamic allocation** of threads
  - > **Unbounded** number of concrete threads
- **Finite approximation**
  - > **One** abstract thread may represent **many** concrete threads
  - > For instance, Java
    - Threads are objects
    - Threads are identified by reference
    - One abstract reference may represent many concrete references
- **Self interleaving of abstract threads**

# An example

Main Thread	Thread
<pre>i=0; for(int j=0;j&lt;input;j++)     new Thread().start();</pre>	<pre>i++;</pre>

MainThread: 

j	0
---	---

 → 

j	1
---	---

 → 

j	2
---	---

 → 

j	3
---	---

 →

Thread: 

i	0
---	---

 → 

i	1
---	---

Thread: 

i	1
---	---

 → 

i	2
---	---

Thread: 

i	2
---	---

 → 

i	3
---	---

# Self interleaving

Main Thread	Thread
<pre>i=0; for(int j=0;j&lt;input;j++)   new Thread().start();</pre>	<pre>i++;</pre>

MainThread:



Thread:



Thread:



Thread:



After K iterations we apply widening

5<sup>th</sup> iteration

.....

Thread:



n<sup>th</sup> iteration

# Java

- This framework has been **applied to Java**
- Everything is identified or accessed by reference
  - > Threads are objects
  - > The shared memory is the heap
  - > Monitors are defined on objects
- We need to develop a specific alias analysis
  - > That is, how we approximate reference
  - > **May** aliasing
    - Identify threads
    - Discover conflicts on the heap
  - > **Must** aliasing to check mutual exclusion on monitors

# Outline

1. Introduction
2. Concrete semantics
3. Abstract semantics
4. Experimental results
5. Conclusion and future work

# Implementation

- Implemented in **Checkmate**

<http://www.pietro.ferrara.name/checkmate>

- **Generic analyzer of multithreaded programs**
- It analyzes **Java bytecode**
- **Pluggable with different**
  - > **Numerical non-relational domains**
  - > **Properties**
  - > **Memory models**
    - **HBMM**
    - **Two relaxed versions**

# Demo

# Eclipse

Pietro Ferrara: "Abstract interpretation of memory models"  
IFIP WG 2.3 "Programming Methodology", Lachen, Switzerland

# Case studies

Th1	Th2
r1=x; y=1; r2=x;	x=1; r3=y;

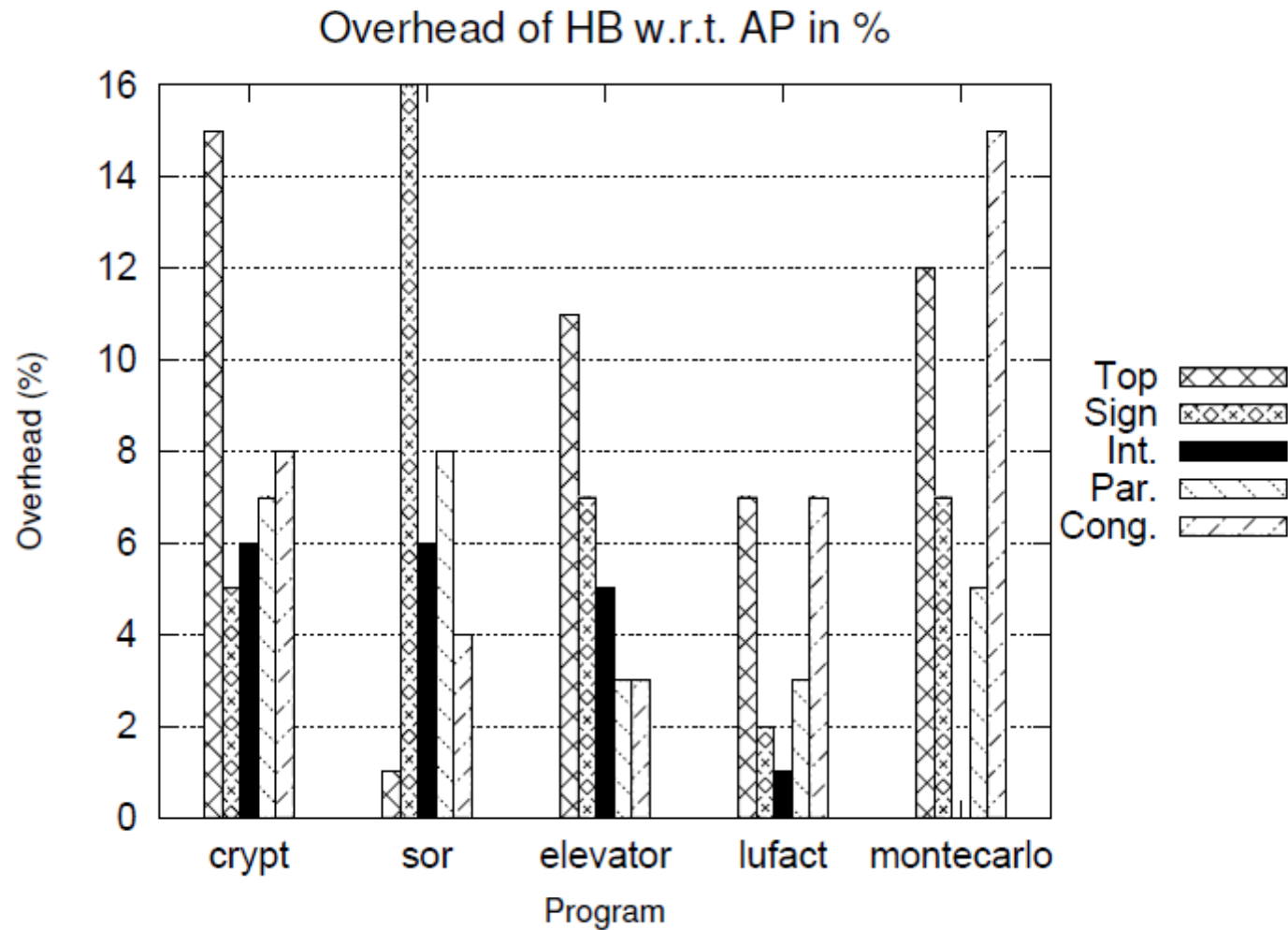
Th1	Th2
r1=x; if(r1!=0) y=42;	r2=y; if(r2!=0) x=42;

- Applied to eight case studies taken from JMM
- **Usually** we catch the behaviors precisely
- Too much approximated in some cases
  - > In particular when dealing with **volatile** variables
  - > Synchronization not taken into account
    - We may **extend** our formalization

# Benchmarks

- An **incremental** application
  - > HBMM compared to two more approximated MMs
  - > Study the overhead of more precise MMs
  - > Overhead of HBMM increases **linearly**
- Some well-known benchmarks
  - > We are in position to analyze programs composed by
    - An **unbounded** number of threads
    - Some **thousands** of bytecode statements
  - > Usually in a couple of minutes

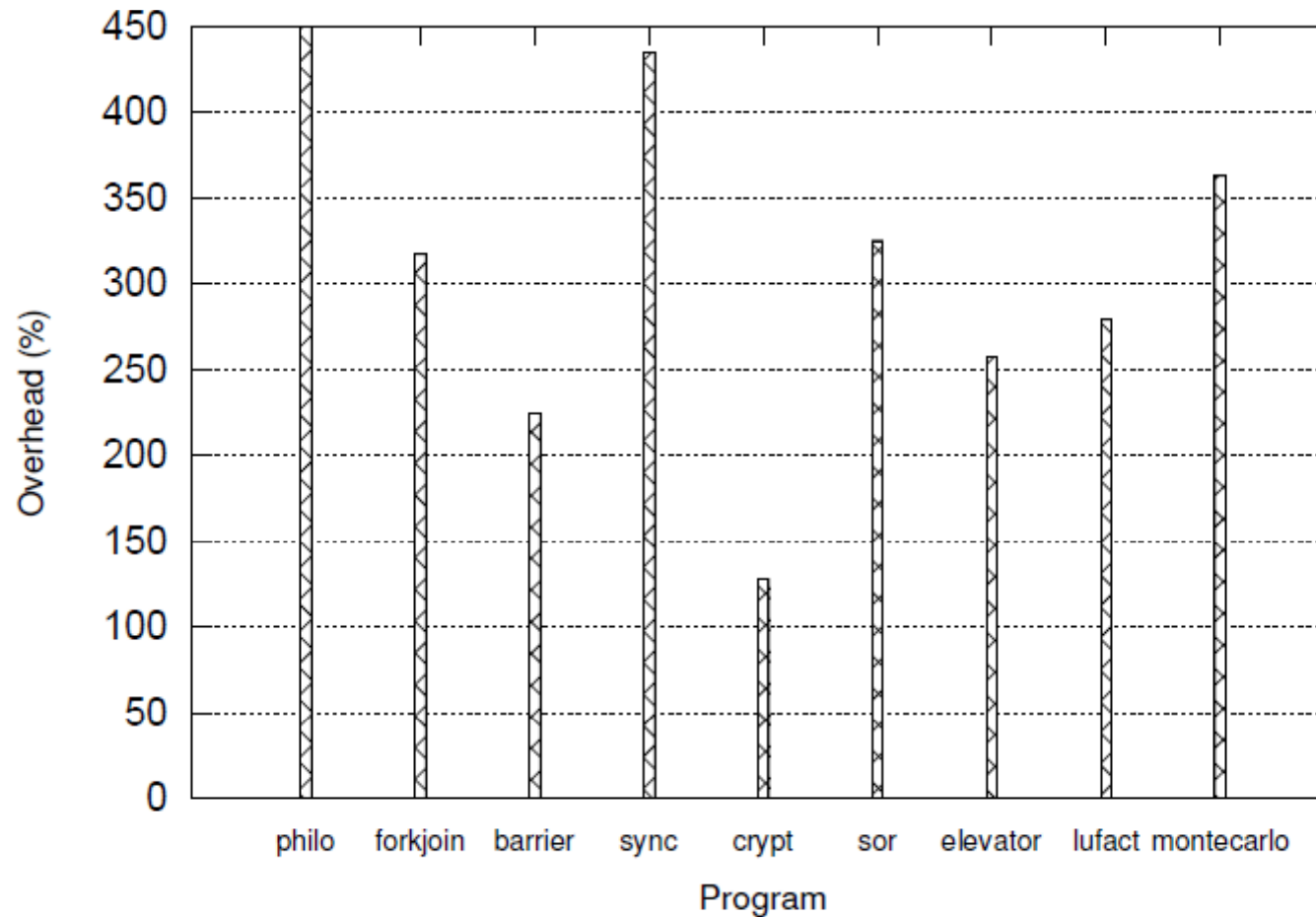
# Overhead of HBMM



Pietro Ferrara: "Abstract interpretation of memory models"  
IFIP WG 2.3 "Programming Methodology", Lachen, Switzerland

# Overhead multithread semantics

Overhead in % of multithread fixpoint computation using Intervals and HB memory model



Pietro Ferrara: "Abstract interpretation of memory models"  
IFIP WG 2.3 "Programming Methodology", Lachen, Switzerland

# Outline

1. Introduction
2. Concrete semantics
3. Abstract semantics
4. Experimental results
5. Conclusion and future work

# Conclusion

- We proposed a
  - > Generic static analysis
  - > Of multithreaded programs
  - > Sound with respect to HBMM
- The framework is more generic
  - > It could be applied to other MM
  - > It can be extended to other synchronization patterns
  - > It can be applied to different programming languages
- Applied to Java bytecode multithreaded programs
  - > Fast and precise

# Limits

- The overall analysis **does not scale up**
- Support only **few** synchronization patterns
  - > Monitors
  - > Launch of threads
- We do not support
  - > Join of threads
  - > volatile variables
  - > Rendez-vous style synchronizations
    - Difficult to track them precisely and efficiently

# Future work - Precision

- Support more synchronization patterns
  - > In particular, rendez-vous style patterns
- Support **relational** numerical domains
  - > Up to now, only non-relational domains
    - Sign, Intervals, ...
  - > Relational information required by industrial programs
- Apply the framework to **other memory models**
  - > Java MM
  - > Sequential consistency

# Future work - Scalability

- Usually, few static analyzers scale up
  - > Necessary modular reasoning
  - > For instance, relying on contracts
    - Clousot @ Microsoft
  - > Checkmate: 1<sup>st</sup> generic analyzer of multithreading
- Contracts to modularly reason on multithreading
  - > For instance, Chalice
- Idea: modular analysis based on these contracts



# Thank you!

Pietro Ferrara: "Abstract interpretation of memory models"  
IFIP WG 2.3 "Programming Methodology", Lachen, Switzerland