

A Structured and High-Level Definition of Java
and of its Provably Correct and Secure Implementation
on the Java Virtual Machine

(The ASM Java/JVM Project)

Egon Börger

Dipartimento di Informatica, Università di Pisa
<http://www.di.unipi.it/~boerger>

Reference:

Java and the Java Virtual Machine - Definition, Verification, and Validation

R. Stärk, J. Schmid, E. Börger

Springer-Verlag , 2001

see <http://www.inf.ethz.ch/~jbook/>

Goal: Real-life Industrial Case Study Book

Illustrate through a relevant & complex example how to
enhance practical syst design & analysis using ASMs

for rigorous high-level modeling

linked seamlessly to executable code

in a verifiable and validatable way

- developing succinct ground models with precise, unambiguous, yet understandable meaning

to provide the possibility for implementation independent system analysis and validation

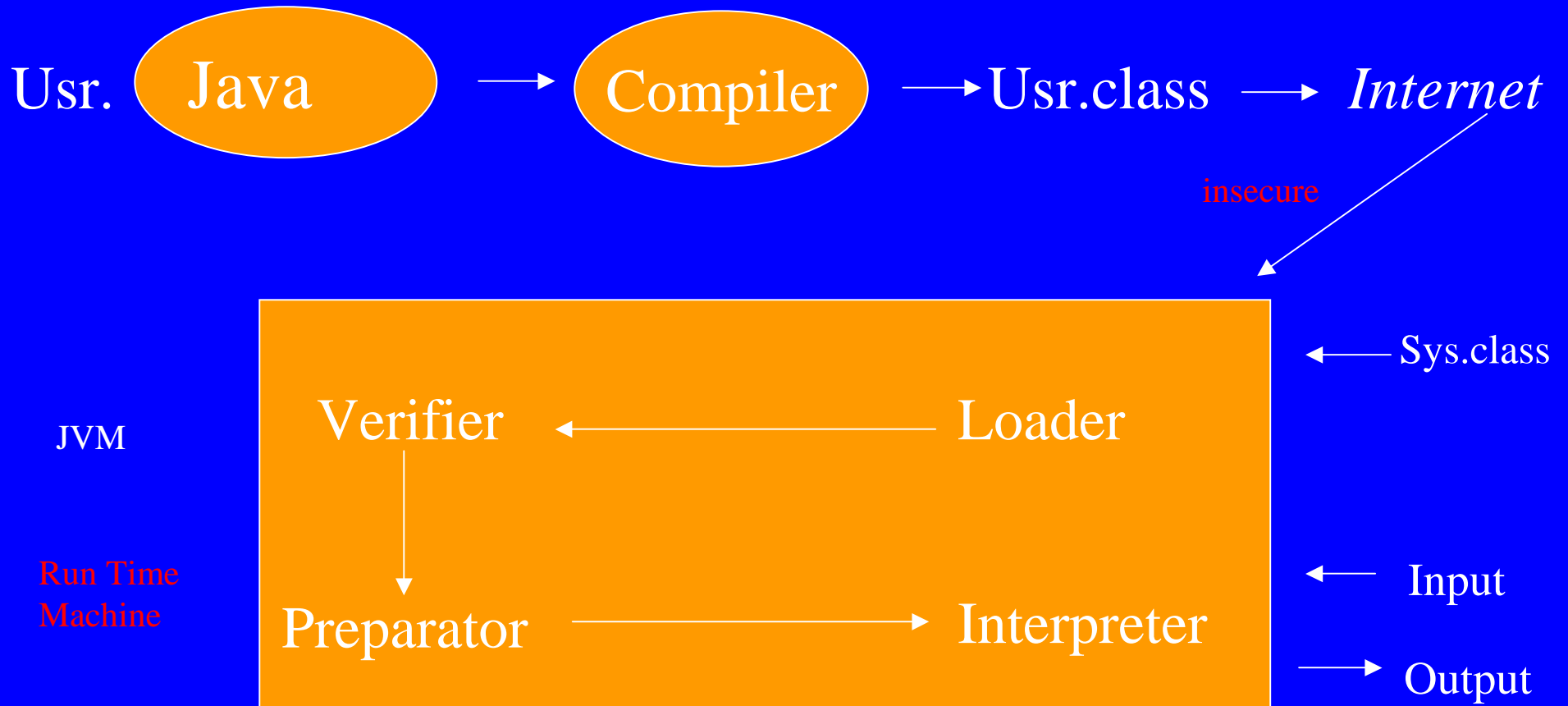
- refining & structuring models into a system (hierarchy) of (sub)models, modularizing orthogonal design decisions (“for change”), justifying them as correct
 - linking the ground model to the implementation
 - documenting the entire design for reuse and maintenance

Method: Separate & Combine Different Concerns using ASMs

- Separating **orthogonal design decisions**
 - to keep design space open (specify for change, avoiding premature design decisions)
 - to structure design space (rigorous interfaces for system (de)composition)
- Separating **design from analysis**
 - separating validation (by simulation) from verification (by proofs)
 - separating verification levels (degrees of proof detail)
 - reasoning for human inspection (design justification)
 - rule based reasoning systems
 - interactive systems
 - automatic tools: model checkers, automatic theorem provers
- **Crossing system levels by most general abstraction and refinement notions** offered by ASMs, tunable to the given problem

The Problem

Java/JVM claimed by SUN to be a safe and secure, platform independent programming env for Internet: **correctness problem** for **compiler**, **loader** (name space support), **verifier**, **access right checker** (security manager), **interpreter**.



Specific Goal of the ASM Java/JVM Project

Abstract (platform independent), **rigorous but transparent, modular definition** providing basis for mathematical and experimental **analysis**

- Reflecting SUN's design decisions (faithful ground model)
- Offering correct high-level understanding (to be practically useful for programmers)
- Providing rigorous, implementation independent basis for
 - Analysis and Documentation (for designers) through
 - Mathematical verification
 - Experimental validation
 - Comparison of different implementations
 - Implementation (compiln, loading, bytecode verification, security schemes)

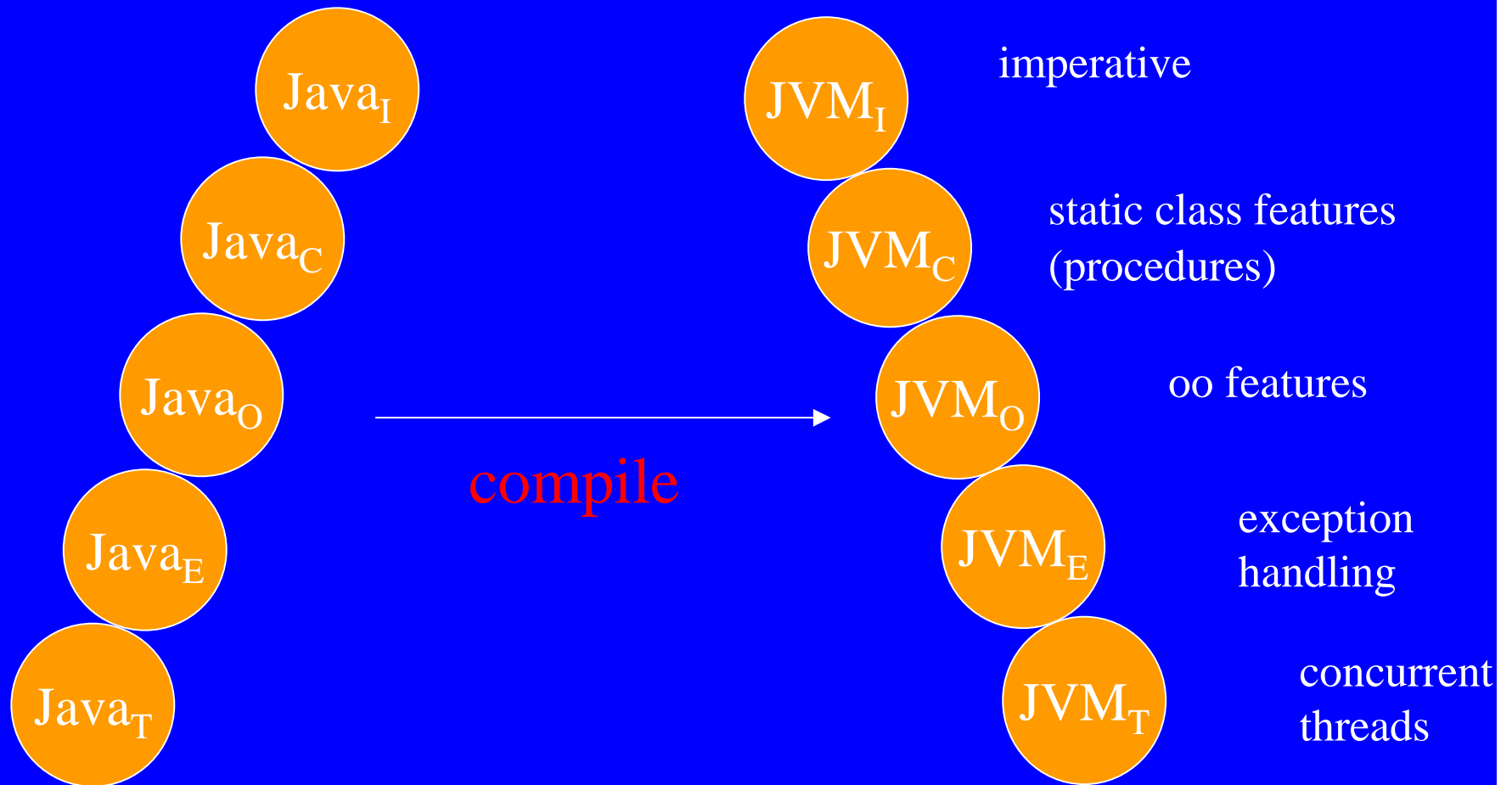
Main Result

A Structured and High-Level Definition of Java
and of its Provably Correct and Secure Implementation
on the Java Virtual Machine

Theorem. Under explicitly stated conditions, any
well-formed and well-typed Java program:

- upon correct compilation
- passes the verifier
- is executed on the JVM
- executes
 - without violating any run-time checks
 - correctly wrt Java source pgm semantics

Language driven decomposition of Java, JVM, compilation



Split into horizontal language **components** (conservative extensions)

The language driven decomposition of execJava and its submachines

execJava =

execJava_I	imperative control constructs
execJava_C	static class features (modules)
execJava_O	oo features
execJava_E	exception handling
execJava_T	concurrent threads

execJava_I =

execJavaExp_I	expression evaluation
execJavaStm_I	statement execution

NB. Grouping similar instructions into one parameterized abstract instr

Pgm exec as walk thru annotated abstract syntax tree

STATE defined by $\text{pos} : \text{Pos}$ $\text{restbody} : \text{Pos} \rightarrow \text{Phrase} \cup \text{Val} \cup \text{Abr}$

MACROS: $\text{context}(\text{pos}) =$ if $\text{restbody}/\text{pos} \in \text{Exp} \cup \text{Bstm}$ or $\text{pos} = \text{first}$
then $\text{restbody}/\text{pos}$
else $\text{restbody}/\text{up}(\text{pos})$

Replacing a phrase (in the current pos) by its result:

$\text{yield}(\text{result}) = \text{restbody} := \text{restbody}[\text{result}/\text{pos}]$

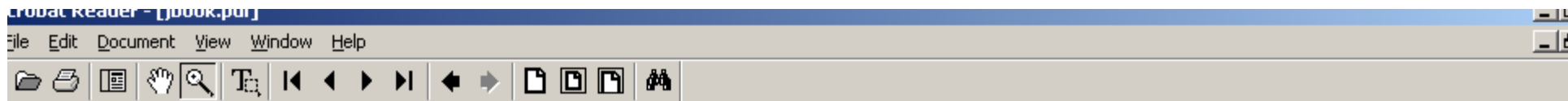
Passing the result of a phrase (in the current pos) to its parent phrase:

$\text{yieldUp}(\text{result}) = \text{restbody} := \text{restbody}[\text{result}/\text{up}(\text{pos})]$
 $\text{pos} := \text{up}(\text{pos})$

Being positioned on a direct subphrase of a structure $f(\dots t \dots)$:

$s = f(\dots \blacktriangleright t \dots)$ stands for $s = f(\dots t \dots) \ \& \ \text{pos} = \blacktriangleright \ \& \ \text{restbody}(\text{pos}) = t$

Phrase: exps & block stms **Val:** values **Abr:** reasons for abruption



$execJavaExp_I = \text{case } context(pos) \text{ of}$
 $lit \rightarrow yield(JLS(lit))$

$loc \rightarrow yield(locals(loc))$

$uop^\alpha exp \rightarrow pos := \alpha$

$uop^{\blacktriangleright} val \rightarrow yieldUp(JLS(uop, val))$

$^\alpha exp_1 bop^\beta exp_2 \rightarrow pos := \alpha$

$^{\blacktriangleright} val bop^\beta exp \rightarrow pos := \beta$

$^\alpha val_1 bop^{\blacktriangleright} val_2 \rightarrow \text{if } \neg(bop \in divMod \wedge isZero(val_2)) \text{ then}$
 $yieldUp(JLS(bop, val_1, val_2))$

$loc = ^\alpha exp \rightarrow pos := \alpha$

$loc = ^{\blacktriangleright} val \rightarrow locals := locals \oplus \{(loc, val)\}$
 $yieldUp(val)$

$^\alpha exp_0 ? ^\beta exp_1 : ^\gamma exp_2 \rightarrow pos := \alpha$

$^{\blacktriangleright} val ? ^\beta exp_1 : ^\gamma exp_2 \rightarrow \text{if } val \text{ then } pos := \beta \text{ else } pos := \gamma$

$^\alpha True ? ^{\blacktriangleright} val : ^\gamma exp \rightarrow yieldUp(val)$

$^\alpha False ? ^\beta exp : ^{\blacktriangleright} val \rightarrow yieldUp(val)$

```

execJavaStmJ = case context(pos) of
;      → yield(Norm)
α exp; → pos := α
► val; → yieldUp(Norm)

break lab;      → yield(Break(lab))
continue lab;   → yield(Continue(lab))
lab : α stm      → pos := α
lab : ► Norm      → yieldUp(Norm)
lab : ► Break(labb) → if lab = labb then yieldUp(Norm)
                                     else yieldUp(Break(labb))
lab : ► Continue(labc) → if lab = labc then yield(body/pos)
                                     else yieldUp(Continue(labc))
phrase(► abr) → if pos ≠ firstPos ∧ propagatesAbr(restbody/up(pos)) then
                                     yieldUp(abr)

{ }                → yield(Norm)
{α1 stm1 ... αn stmn} → pos := α1
{α1 Norm ... ► Norm}     → yieldUp(Norm)
{α1 Norm ... ► Normαi+1 stmi+1 ... αn stmn} → pos := αi+1

if (α exp)β stm1 else γ stm2 → pos := α
if (► val)β stm1 else γ stm2 → if val then pos := β else pos := γ
if (α True) ► Norm else γ stm → yieldUp(Norm)
if (α False)β stm else ► Norm → yieldUp(Norm)

while (α exp)β stm → pos := α
while (► val)β stm → if val then pos := β else yieldUp(Norm)
while (α True) ► Norm → yieldUp(body/up(pos))

Type x; → yield(Norm)

```

propagatesAbr
iff phrase is no:
labeled stm

Later refined : no
static initializer
try stm
finally stm
synchronized stm

The $\text{execJava}_{C/O}$ extensions

$\text{execJava}_C =$

execJavaExp_C

extending expression evaluation

execJavaStm_C

extending statement execution

Adding

- class fields (global variables)
- class method invocation/return (procedures)
- class initializers (module initializers)

$\text{execJava}_O =$

execJavaExp_O

adding instance fields/methods

Fields treated similarly to local vars (with local replaced by **global**), but: one has to **initialize** each class at its first active use, i.e. when for the first time accessing (or assigning to) some of its fields or calling some of its methods (after **left-to-right arg evaluation**) (or upon creation in Java₀)

```

execJavaExpC = case context(pos) of
  c.f          → if initialized(c) then yield(globals(c/f)) else initialize(c)
  c.f = α exp → pos := α
  c.f = ▶ val → if initialized(c) then
                globals(c/f) := val
                yieldUp(val)
                else initialize(c)

  c.mα(exps) → pos := α
  c.m▶(vals) → if initialized(c) then invoke(up(pos), c/m, vals)
                else initialize(c)

  () → yield([])
  (α1 exp1, ..., αn expn) → pos := α1
  (α1 val1, ..., ▶ valn) → yieldUp([val1, ..., valn])
  (α1 val1, ..., ▶ vali, αi+1 expi+1 ... αn expn) → pos := αi+1

```

Execution of initialization code for a class is started only when the **superclass** is already initialized, and also at the top of the class hierarchy. A class becomes initialized upon exiting from its initialization method.

```

execJavaStmC = case context(pos) of
  static αstm      → let c = classNm(meth)
                    if c = Object ∨ initialized(super(c)) then pos := α
                    else initialize(super(c))
  static ▶ Return → exitMethod(Norm)
                    classState(classNm(meth)) := Initialized

  return αexp;      → pos := α
  return ▶ val;     → yieldUp(Return(val))
  return;          → yield(Return)
  lab : ▶ Return   → yieldUp(Return)
  lab : ▶ Return(val) → yieldUp(Return(val))
  Return          → if pos = firstPos ∧ ¬null(frames) then
                    exitMethod(Norm)
  Return(val)     → if pos = firstPos ∧ ¬null(frames) then
                    exitMethod(val)

  ▶ Norm; → yieldUp(Norm)

```

instance field values of objects stored (using `setField`) in & retrieved (using `getField`) from **Heap**, under the ref of the object; default values assigned upon creation.

The class of new parametrized class instances is initialized before parameter evaln.

this stored as local var ; bound by inst meth call & by return from a constructor (to the newly created object, see the extension of `exitMethod`)

```
execJavaExp0 = case context(pos) of
  this → yield(locals( "this" ))
```

```
new c/mα(exps) → if initialized(c) then pos := α else initialize(c)
new c/m►(vals) → create ref
  heap(ref) := Object(c, {(f, defaultVal(type(f)))
    | f ∈ instanceFields(c)}))
  invoke(up(pos), c/m, [ref] · vals)
```

```
α exp.c/f → pos := α
► ref.c/f → if ref ≠ null then yieldUp(getField(ref, c/f))
```

```
α exp1.c/f = β exp2 → pos := α
► ref.c/f = β exp → pos := β
α ref.c/f = ► val → if ref ≠ null then
  setField(ref, c/f, val)
  yieldUp(val)
```

```
α exp instanceof c → pos := α
► ref instanceof c → yieldUp(ref ≠ null ∧ classOf(ref) ≤h c)
```

```
(c)α exp → pos := α
(c)► ref → if ref = null ∨ classOf(ref) ≤h c then yieldUp(ref)
```

```
α exp.c/mβ(exps) → pos := α
► ref.c/mβ(exps) → pos := β
α ref.c/m►(vals) → if ref ≠ null then
  let c' = case callKind(up(pos)) of
    Virtual → lookup(classOf(ref), c/m)
    Super → lookup(super(classNm(meth)), c/m)
    Special → c
  invoke(up(pos), c'/m, [ref] · vals)
```


The $\text{execJava}_{E/T}$ extensions

$\text{execJava}_E =$

execJavaExp_E for evaluation of run-time exceptions

execJavaStm_E for execution of exception statements

$\text{execJava}_T =$

execJavaStm_T for synchronization statements

(as part of execJavaThread)

Abrs in try stms:
caught excs lead to
 catch code exec,
othr abrs propagate

catch code yields
 up Norm or an abr

For finally stms:
 abrs suspend upon
 entering finally stm

exiting propagates up
 the suspended abr
 (resumed) or a new abr

Uncaught excs
 propagate up the method
 call stack; in static class
 initializers they make the
 class unusable

```

execJavaStmE = case context(pos) of
  throw α exp; → pos := α
  throw ▶ ref; →
    if ref = null then fail(NullPointerException) else yieldUp(Exc(ref))

  try α stm catch ... → pos := α
  try ▶ Norm catch ... → yieldUp(Norm)
  try ▶ Exc(ref) catch (c1 x1)β1 stm1 ... catch (cn xn)βn stmn →
    if ∃ 1 ≤ j ≤ n : classOf(ref) ≤h cj then
      let j = min{i | classOf(ref) ≤h ci}
      pos := βj
      locals := locals ⊕ {(xj, ref)}
    else yieldUp(Exc(ref))
  try ▶ abr catch (c1 x1)β1 stm1 ... catch (cn xn)βn stmn → yieldUp(abr)
  try α Exc(ref) ... catch (ci xi) ▶ Norm ... → yieldUp(Norm)
  try α Exc(ref) ... catch (ci xi) ▶ abr ... → yieldUp(abr)

α stm1 finally β stm2 → pos := α
▶ Norm finally β stm → pos := β
▶ abr finally β stm → pos := β
α s finally ▶ Norm → yieldUp(s)
α s finally ▶ abr → yieldUp(abr)

Exc(ref) → if pos = firstPos ∧ ¬null(frames) then
  exitMethod(Exc(ref))

lab : ▶ Exc(ref) → yieldUp(Exc(ref))
static ▶ Exc(ref) →
  classState(classNm(meth)) := Unusable
  if classOf(ref) ≤h Error then exitMethod(Exc(ref))
  else fail(ExceptionInInitializerErr)
  
```

Examples of run-time exceptions

```
execJavaExpE = case context(pos) of  
  αval1 bop ▶ val2 → if bop ∈ divMod ∧ isZero(val2) then  
    fail(ArithmeticException)  
  ▶ ref.c/f → if ref = null then fail(NullPointerException)  
  αref.c/f = ▶ val → if ref = null then fail(NullPointerException)  
  αref.c/m ▶ (vals) → if ref = null then fail(NullPointerException)  
  (c) ▶ ref → if ref ≠ null ∧ classOf(ref)  $\not\subseteq_h$  c then  
    fail(ClassCastException)
```

where **fail (exc)** = yield (throw new exc() ;)

When classes become unusable, their initialization is impossible, so that **initialize(c)** is extended by the following:

if classState(c) = Unusable then fail (NoClassDefFoundError)

Theorem: Java is type safe

- i.e. when a legal well-typed Java pgm is executed:
 - run-time vals of static/instance fields/array elems are compatible with their declared types
 - references to objects are in the heap (no dangling pointers)
 - run-time positions satisfy compile-time constraints (reachable, definitely assigned vars are well-defined local vars with vals of compile-time type,...)
 - positions of normally completing stms are compile-time normal
 - evaluated exprs/returned vals have compile-time compatible type
 - abruptions (jump,return,exc) have compile-time compatible type
 - stacks do not overflow nor underflow, ...
- **Proof:** induction on Java ASM runs, based upon a rigorous definition of the rules for definite assignment

Extending execJava, to become component of ExecJavaThread

```
execJavaStmT = case context(pos) of  
  synchronized (α exp) β stm → pos := α  
  synchronized (▶ ref) β stm →  
    if ref = null then fail(NullPointerException)  
    else  
      if ref ∈ sync(thread) then  
        sync(thread) := [ref] · sync(thread)  
        locks(ref) := locks(ref) + 1  
        pos := β  
      else  
        exec(thread) := Synchronizing  
        syncObj(thread) := ref  
        cont(thread) := (frames, (meth, restbody, β, locals))  
  synchronized (α ref) ▶ Norm → releaseLock(Norm)  
  synchronized (α ref) ▶ abr → releaseLock(abr)  
  
  static ▶ abr → notifyThreadsWaitingForInitialization  
  abr → if pos = firstPos ∧ null(frames) then killThread
```

Abstract scheduling of Multiple Threads: inserting execJava into ExecJavaThread

Thread scheduling separated from thread execution

ExecJavaThread \equiv

choose q **in** $dom(exec), runnable(q)$

if $q=thread$ and $exec(q)=Active$

then execJava

else

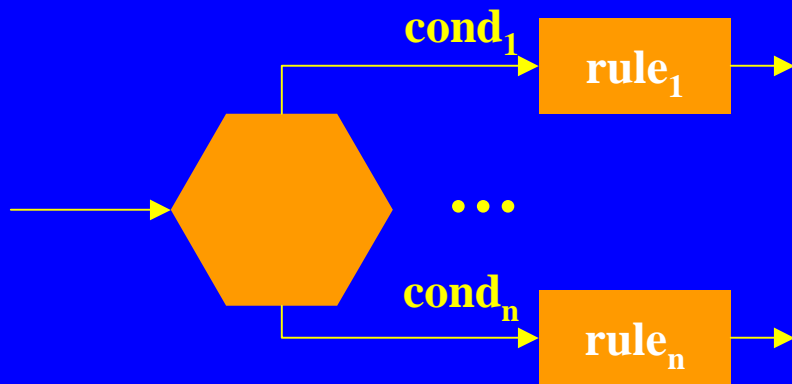
if $exec(q)=Active$ **then**

$cont(thread) := (frames, (methd, restbody, pos, locals))$

$thread := q$

run(q)

Diagram notation for Control State ASMs



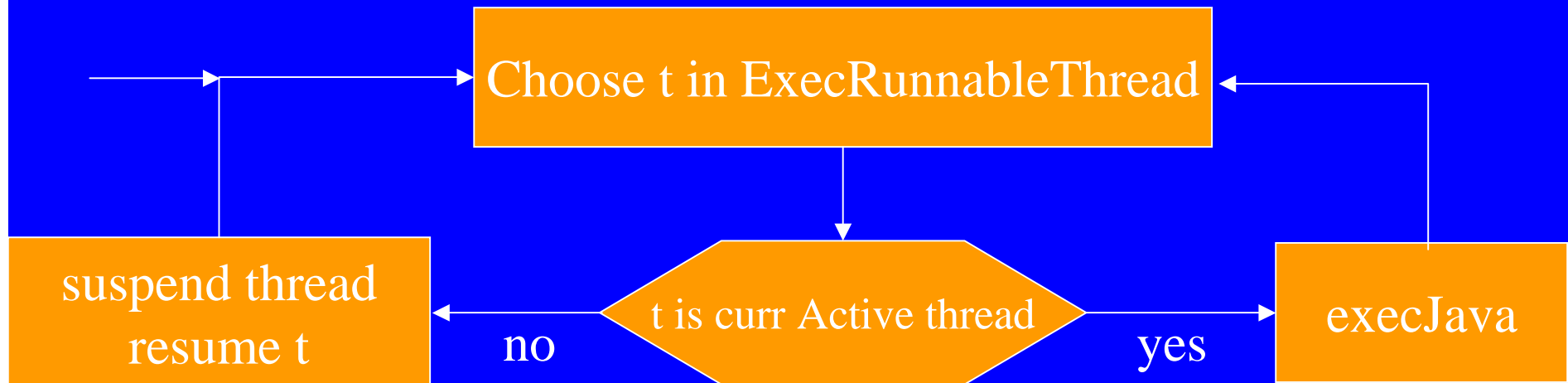
UML: combined
branching/action
nodes

meaning

labeling of the arrows
by “control” states
often suppressed

```
if  $ctl = i$  then
  if  $cond_1$  then  $rule_1$ 
                     $ctl := j_1$ 
                    ...
  if  $cond_n$  then  $rule_n$ 
                     $ctl := j_n$ 
```

Defining execJavaThread as control state ASM



Thread scheduling separated from thread execution

t in **ExecRunnableThread** = (t in $\text{dom}(\text{exec})$ & $\text{runnable}(t)$)

t is **curr Active thread** = ($t = \text{thread}$ & $\text{exec}(t) = \text{Active}$)

suspend thread = if $\text{exec}(\text{thread}) = \text{Active}$
then $\text{cont}(\text{thread}) := (\text{frames}, \text{currframe})$

resume(t) = $\text{thread} := t$
 $\text{run}(t)$

Theorem: Correctness of Thread Synchronization in Java

- Runtime threads are valid threads (of type `Thread`).
- If the execution state of a thread is `Not Started`, then the thread is not synchronized on any object and is not in the wait set of any object.
- If the state of a thread is `synchronizing`, then the thread is not already synchronized on the object it is competing for.
- If a thread is synchronized on an object, then the object is a valid reference in the heap.
- If a thread is waiting for an object, then it is synchronized on and is in the wait set of the object (without holding the lock of the object).
- If a thread has been notified on an object, then it is no longer in the wait set of the object. It is still synchronized on the object, but it does not hold the lock of the object.

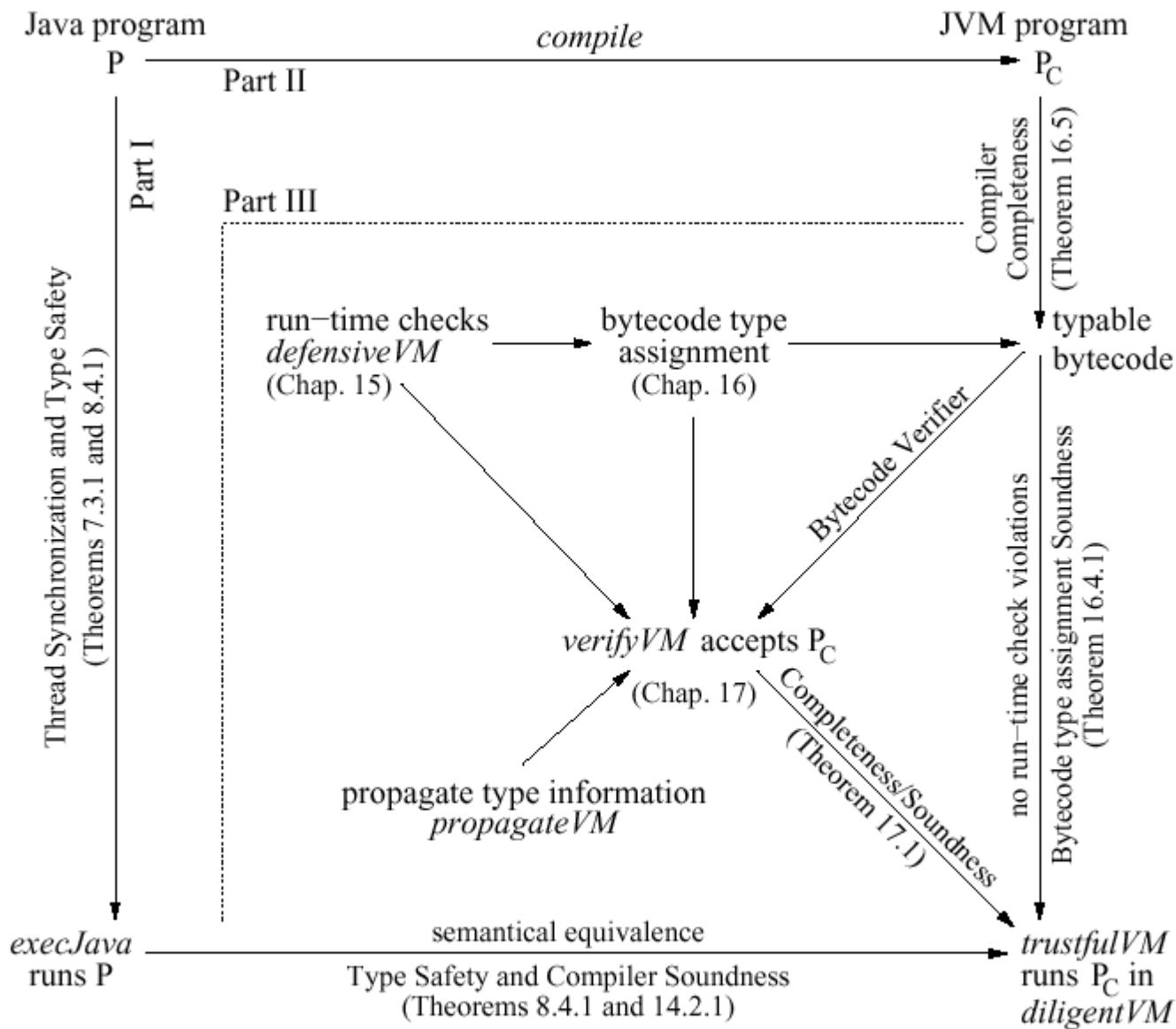
Theorem: Correctness of Thread Synchronization in Java (Cont'd)

- A thread cannot be in the wait set of two different objects.
- If a thread has terminated normally or abruptly, then it does not hold the lock of any object.
- If a thread holds the lock of an object, then the lock counter of the object is exactly the number of occurrences of the object in the list of synchronized objects of the thread.
- It is not possible that at the same time, two different threads hold the lock of the same object.
- If the lock counter of an object is greater than zero, then there exists a thread which holds the lock of the object.
- ...

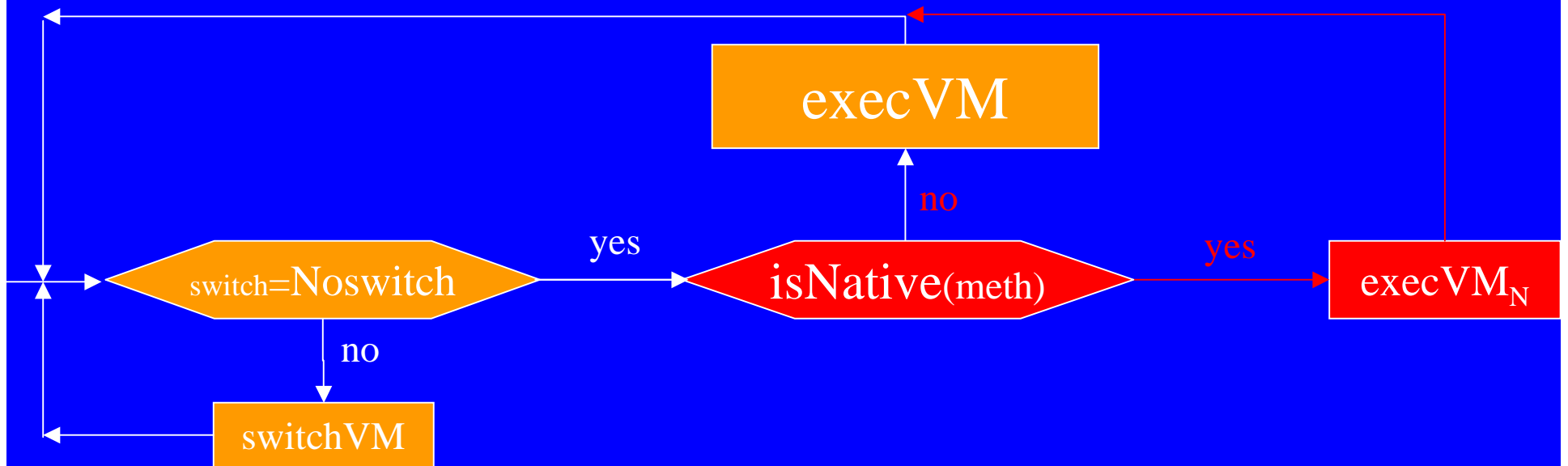
PROOF. Induction on Java ASM runs.

Security Driven JVM Decomposition

- **trustfulVM**: defines the execution functionality incrementally from language layered submachines **execVM**, **switchVM**
- **defensiveVM**: defines the constraints to be checked, in terms of **trustfulVM** execution, from the language layered submachine **check**; calls **trustfulVM** for execution
- **diligentVM**: checks the constraints at link-time, using a language layered submachine **verifyVM**; calls **trustfulVM** for execution
- **verifyVM** built up from language layered submachines **check**, **propagateVM**, **succ**
- **dynamicVM**: dynamic loading and linking of classes



Stepwise refinement of trustfulVM



execVM and switchVM incrementally extended (language driven)

$\text{trustfulVM}_I = \text{execVM}_I \subseteq \text{execVM}_C \subseteq \text{execVM}_O \subseteq \text{execVM}_E$

$\text{execVM}_N \subseteq \text{execVM}_D$ defining instructionwise **changes of current frame**

$\text{switchVM}_C \subseteq \text{switchVM}_E \subseteq \text{switchVM}_D$ defining **changes of frame stack**
 reflecting meth call/return, class initialization, capturing exceptions, class load/linking

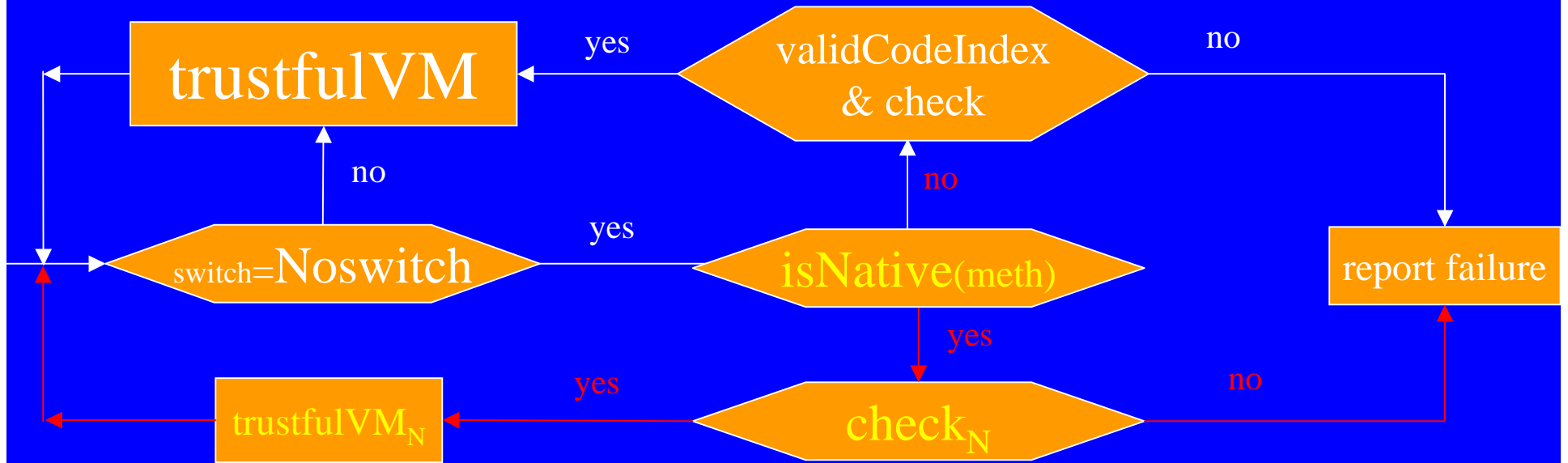
Stating rigorously and proving the Correctness of compiling from Java to JVM

- With respect to the ASM models for Java and JVM, and wrt the definition of **compile from Java to JVM** code, including the exception table, the execution of P in Java and the execution of **compile(P)** in Trustful VM are equivalent (in a sense made precise), for arbitrary pgms P.
- **PROOF.** By induction on the runs of the Java/JVM ASMs, using the type safety theorem.
- NB. This includes the **correctness of exception handling**
see Börger E., Schulte W., A Practical Method for Specification and Analysis of Exception Handling -- A Java/JVM Case Study. IEEE Transactions of Software Engineering, Vol.26, No.10, October 2000 (Special Issue on Exception Handling, eds. D.Perry, A.Romanovsky, A.Tripathi).

Deriving the Bytecode Verifier Conditions from Type Checking Runtime Constraints

- **Defensive VM:** Checks at run-time, before every execution step, the “structural constraints” which describe the verifier functionality (restrictions on run-time data: argument types, valid Ret addresses, resource bounds,...) guaranteeing “safe” execution
- **Static constraints (well-formedness) checked at link-time.**
- **Theorem:** If Defensive VM executes P successfully, then so does Trustful VM, with the same semantical effect.

Stepwise refinement of defensive VM



check incrementally extended, language driven as for trustfulVM

i.e. $check_I$ extended by $check_C$

extended by $check_O$

extended by $check_E$

extended by $check_N$

extended by $check_D$

Bytecode Type Assignments

- Link-time verifiable **type assignments** (conditions) extracted from **checking** function of the Defensive VM

Main problem: return addresses of Jsr(s), reached using Ret(x)

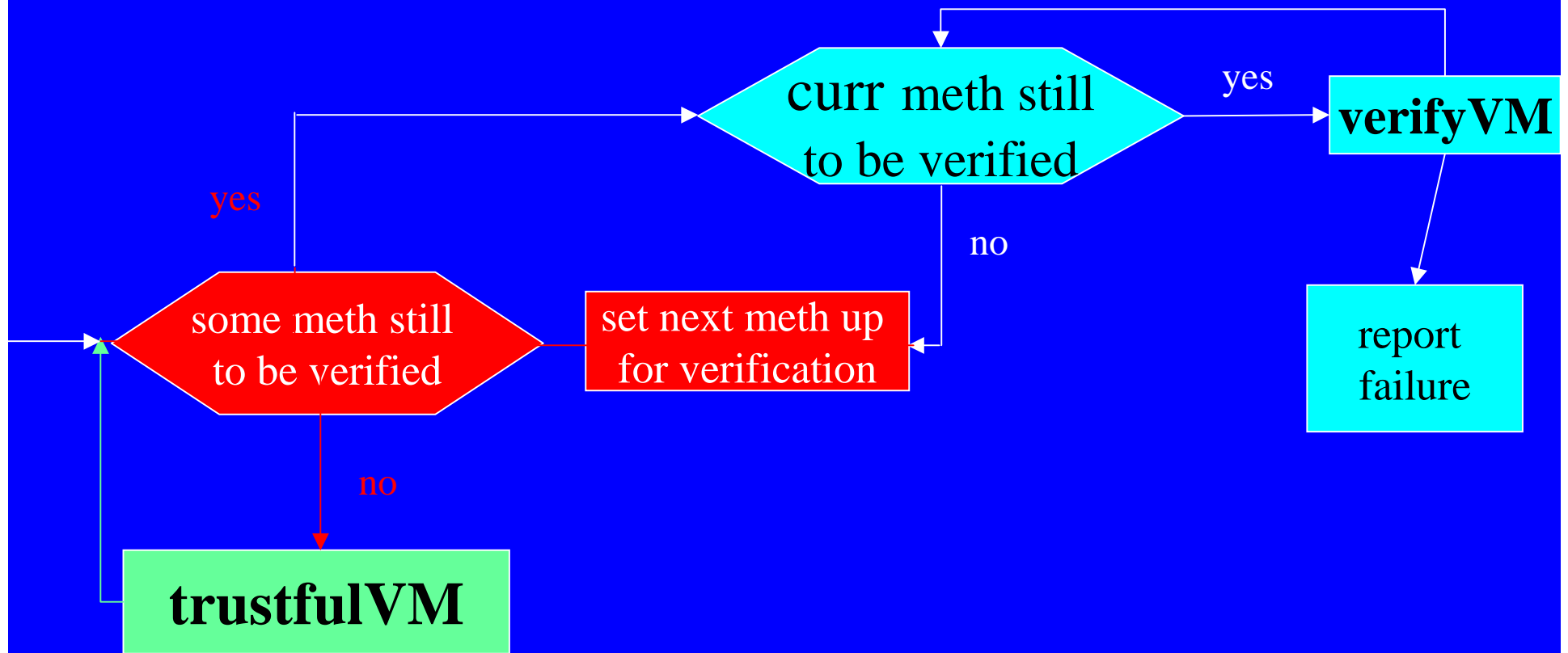
- **Soundness Theorem**: If P satisfies the type assignment conditions, then Defensive VM executes P without violating any run-time check.

Proof by induction on runs of the Defensive VM

- **Completeness Theorem**: Bytecode generated by **compile** from a legal Java program does have type assignments.

Inductive proof introduces **certifying compiler** assigning to each byte code instr also a type frame, which then can be shown to constitute a type assignment for the compiled code

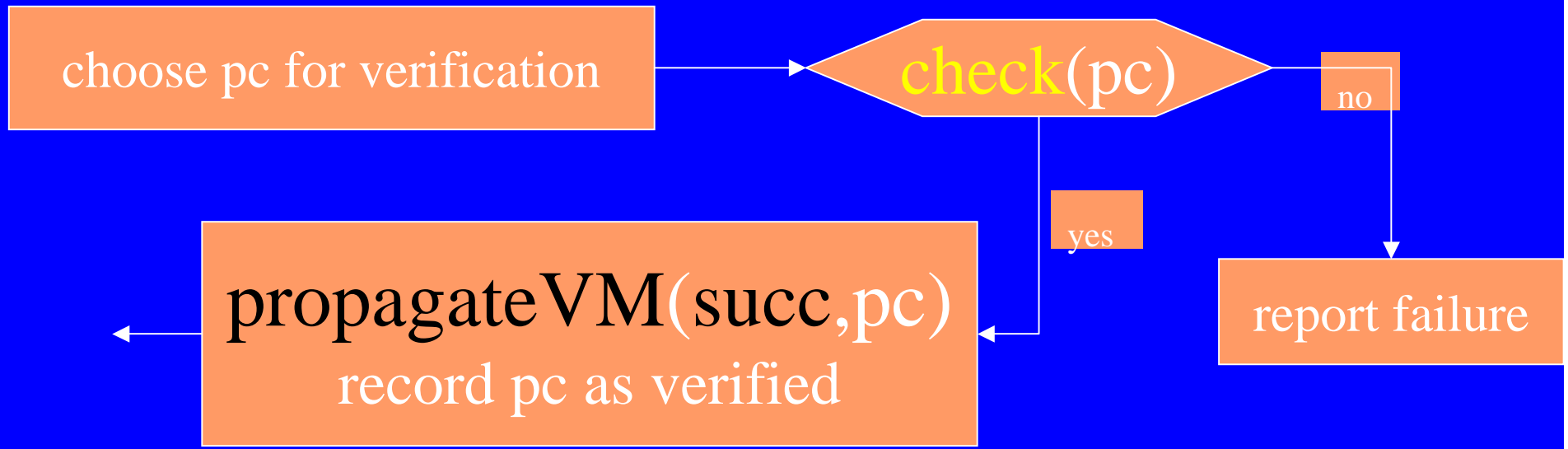
Stepwise refinement of diligent $VM_{I,C,O,E}$



switch VM_C in trustfulVM is refined to also link classes before their initialization, where the linking submachine triggers verifyVM

verifyVM decomped into lang layered check, succ, propagate

Stepwise refinement of verifyVM



propagateVM the checked type frame from pc to all possible successor frames, simulating execVM on types frames

propagateVM and succ incrementally extended

$$\text{succ}_I \subseteq \text{succ}_C \subseteq \text{succ}_O \subseteq \text{succ}_E$$

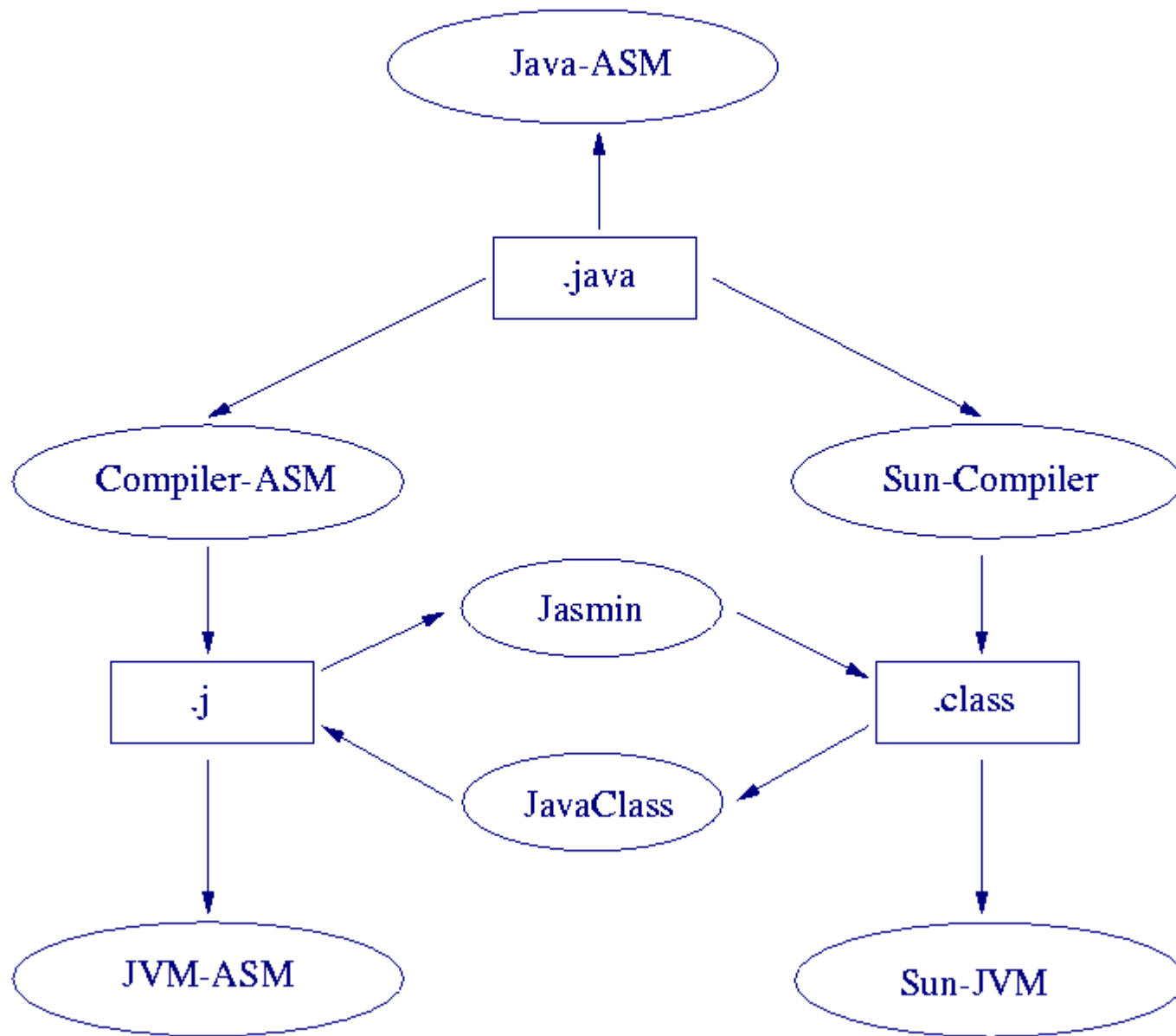
$$\text{propagate}_I \subseteq \text{propagate}_E$$

Proving Bytecode Verifier Complete and Correct

- **Bytecode Verifier Soundness Theorem:** For any program P , the bytecode verifier either rejects P or during the verification satisfies the type assignment conditions for P .
- **Bytecode Verifier Completeness Theorem:** If P has a type assignment, then the Bytecode Verifier does not reject P and computes a most specific type assignment.

Validating Java, JVM, compile

- AsmGofer: ASM programming system, extending TkGofer to execute ASMs (with Haskell definable external fcts)
- Provides **step-by-step execution**, with GUIs to support **debugging** of Java/JVM programs.
- Allows for the executable ASM models of Java/JVM:
 - to execute the Java source code P (**no counterpart in SUN env**)
 - to compile Java pgms P to bytecode compile(P) (in textual representation, using JASMIN to convert to binary class format)
 - to execute the bytecode programs compile(P)
E.g. our Bytecode Verifier rejects Saraswat's program
- Developed by Joachim Schmid, available at www.tydo.de/AsmGofer



Java and the Java Virtual Machine. Definition, Verification, and Validation

R. Stärk, J. Schmid, E. Börger

Springer-Verlag , 2001.

see <http://www.inf.ethz.ch/~jbook/>

For ASMGofer see www.tydo.de/AsmGofer/

My home page <http://www.di.unipi.it/~boerger>