

Applied Static Analysis 2016

Dr. Michael Eichberg (Organizer)

Johannes Lerch, Ben Hermann, Sebastian Proksch, Karim Ali Ph.D.



1

Java Bytecode

A HARDWARE- AND OPERATING SYSTEM-
INDEPENDENT BINARY FORMAT, KNOWN AS
THE CLASS FILE FORMAT.

The Java® Virtual Machine Specification
Java SE 8 Edition
Specification: JBR-000337, Java® SE 8 Release Contents Specification ("Specification") Version: 8
Status: Proposed Final Draft
Release: January 2014
Copyright © 1997, 2014, Oracle America, Inc. and/or its affiliates. All rights reserved. 500 Oracle Parkway, Redwood
City, California 94065, U.S.A.

Structure of the Java Virtual Machine (JVM)

- Data types:
 - Primitive Types:
`boolean`, `byte`, `short`, `int`, `char` (computational type `int`; cat. **1**)
`long`, (computational type `long`; cat. **2**)
`float`, (computational type `float`; cat. **1**)
`double`, (computational type `long`; cat. **2**)
`return address` (computational type `return address`; cat. **1**)
 - Reference Types: (computational type `reference value`; cat. **1**)
`class`,
`array`,
`interface types`

3

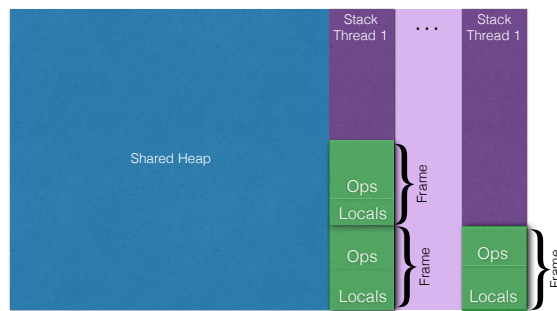
“boolean” only has special support w.r.t. the creation of arrays.

Structure of the Java Virtual Machine (JVM)

- Run-time Data Areas
 - the `pc` (program counter) register contains the address of the instruction that is currently executed by a thread; each thread has its own `pc`
 - each JVM thread has a **private stack** which holds local variables and partial results
 - the **heap** which is shared among all threads
 - **frames** are allocated from a JVM thread's private stack when a method is invoked; each frame has its own array of **local variables** and **operand stack**
 - local variables are indexed
 - a *single local variable* can hold a value of type `boolean`, `byte`, `char`, `short`, `int`, `float`, `reference`, or `return address` (computational type category 1)
 - a *pair of local variables* can hold a value of type `long` or `double` (computational type category 2)
 - the **operand stack** is empty at creation time; an entry can hold any value
 - the **local variables** contains the parameters (including the implicit `this` parameter in local variable 0)

4

Structure of the Java Virtual Machine (JVM)



5

The size of the operand stack and the size of the locals depends on the called method.

Category two values require two slots on the operand stack / two locals.

Structure of the Java Virtual Machine (JVM)

- Special Methods
 - the name of instance initialization methods (Java constructors) is "**<init>**"
 - the name of the class or interface initialization method (Java static initializer) is "**<clinit>**"
- Exceptions are instance of the class **Throwable** or one of its subclasses; exceptions are thrown if:
 - an **athrow** instruction was executed
 - an abnormal execution condition occurred (e.g., division by zero)

6

Structure of the Java Virtual Machine (JVM)

- Instruction Set Summary
 - an instruction consists of a *one-byte opcode* specifying the operation and zero or more operands (arguments to the operation)
 - most instructions *encode type information in their name*; in particular those operating on primitive types (e.g., **iadd**, **fadd**, **dadd**)
 - some are generic and are only restricted by the computational type category of the values (e.g. **swap**, **dup2**)

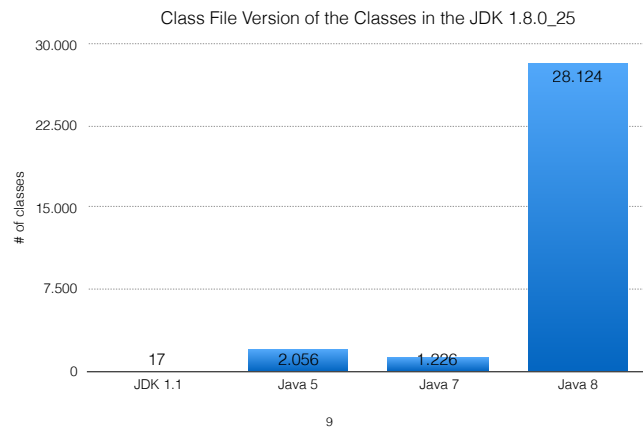
7

Structure of the Java Virtual Machine (JVM)

- Categories of Instructions
 - Load and store instructions (e.g., **aload_0**, **istore(x)**)
(Except of the load and store instructions, the only other instruction that manipulates a local variable is **inc**.)
 - Arithmetic instructions (e.g., **iadd**, **iushr**)
 - (Primitive/Base) Type conversion instructions (e.g., **i2d**, **l2d**, **l2i**)
 - Object/Array creation and manipulation (e.g., **new**, **checkcast**)
 - (Generic) Operand Stack Management Instructions (e.g., **dup**)
 - Control Transfer Instructions (e.g., **itlt**, **if_icmplt**, **goto**, **jsr**, **ret**)
(Some are further modified using the **wide** modifier.)
obsolete since Java 5
 - Method Invocation and Return instructions (e.g., **invokespecial**, **return**)
 - Throwing Exceptions (**athrow**)
 - Synchronization (**monitorenter**, **monitorexit**)

8

Obsolete Code in the JDK



Structure of the Java Virtual Machine (JVM)

- The maximum length of a method is 65536.
(This is a frequent issue with generated code.)
- A method can have only 65536 local variables.
(The maximum number of locals in the JDK is 142.)
(The maximum number of locals in OPAL is/was 1136.)
- The maximum stack size of a single method is 65536.
(The maximum stack size of any method in the JDK is 42.)

10

Java Bytecode

PC	Instruction
0	iload_0
1	iload_1
2	if_icmple →7
5	iload_0
6	ireturn
7	iload_1
8	ireturn

```
static int max(int i, int j) {
    if (i > j)
        return i;
    else
        return j;
}
```

11

The parameters are stored in the locals when the method is invoked. Recall that double and long values need two slots.
(In case of instance methods, the implicit this parameter is stored in local variable 0.)

Java Bytecode

```
static int max(int i, int j) {
    if (i > j)
        return i;
    else
        return j;
}
```

PC	Instruction	Operand Stack	Registers/Local Vars.
0	iload_0		0: an int 1: an int
1	iload_1	an int	0: an int 1: an int
2	if_icmple →7	an int an int	0: an int 1: an int
5	iload_0		0: int ∈ [-2147483647, MAX] 1: int ∈ [MIN, 2147483646]
6	ireturn	int ∈ [-2147483647, MAX]	0: int ∈ [-2147483647, MAX] 1: int ∈ [MIN, 2147483646]
7	iload_1		0: an int 1: an int
8	ireturn	an int	0: an int 1: an int

12

Notice the compilation of the if statement. It is common that in the bytecode the if operator is the inverse one, because if the condition evaluates to true, we then perform the jump (to the else branch) while in the source code, we simply fall through in case the condition evaluates to true.

Java Bytecode - Loops

```
static int factorial(int n) {
    int r = 1;
    while (n > 0) {
        r *= n;
        n--;
    }
    return r;
}
```

PC	Instruction
0	iconst_1
1	istore_1
2	goto →12
5	iload_1
6	iload_0
7	imul
8	istore_1
9	iinc (lv=0, val=-1)
12	iload_0
13	ifgt →5
16	iload_1
17	ireturn

13

When we have a loop, we will always find an unconditional goto in the bytecode. Here, the if condition is actually the same as in the source code!

By moving the test to the end, we have no goto instruction while executing the nth iteration ($n > 1$).

Java Bytecode - Object Creation

```
static int numberOfDigits(int i) {
    assert (i > 0);
    return ((int) Math.floor(Math.log10(i))) + 1;
}
```

PC	Instruction	Code
0	getstatic TACDemo { boolean \$assertionsDisabled }	assert(i > 0)
3	ifne →18	
6	iload_0	
7	ifgt →18	
10	new java.lang.AssertionError	
13	dup	
14	invokespecial java.lang.AssertionError { void <init> () }	
17	athrow	Math.log10(i)
18	iload_0	
19	i2d	
20	invokestatic java.lang.Math { double log10 (double) }	
23	invokestatic java.lang.Math { double floor (double) }	Math.floor(...)
26	d2i	(int) "Typecast"
27	iconst_1	1
28	iadd	+
29	ireturn	return

14

Initialization of an object is done in two steps. First allocation of the memory (using new) and then calling the constructor (<init>).

Assertions are compiled to a field access instruction to check if assertions are enabled. Afterwards (conditionally), the assertion is checked.

Java Bytecode - Basic Exception Handling

▼ Method Body (Size: 31 bytes, Max Stack: 2, Max Locals: 3)

PC	Line	Instruction	Exceptions
0	109	aload_1	
1		checkcast	
4		astore_2	
5	110	aload_2	
6		invokeinterface	
11		lshl	
14	111	iconst_m1	
15		aload_2	
16		invokeinterface	
21		idiv	
22		ireturn	
23	113	iconst_1	
24		ireturn	
25	114	astore_2	
26	115	iconst_0	
27		ireturn	
28	116	astore_2	
29	117	aload_2	
30		athrow	

int tryCatch(Object o) {
 try {
 List<?> l = (List<?>) o;
 if (l.size() == 0)
 return -1 / l.size();
 else
 return 1;
 } catch (ClassCastException cce) {
 return 0;
 } catch (Error | RuntimeException e) {
 throw e;
 }
}

▼ Exception Table:

1.	try [0-22]	catch 25	java.lang.ClassCastException
2.	try [0-22]	catch 28	java.lang.Error
3.	try [0-22]	catch 28	java.lang.RuntimeException

inclusive →
 exclusive → 15

At runtime the first handler that can handle an exception will be invoked.

Java Bytecode - Finally

▼ Method Body (Size: 39 bytes, Max Stack: 1, Max Locals: 5)

PC	Line	Instruction	Exceptions
0	122	iconst_0	
1		istore_2	
2	124	aload_1	
3		instanceof	
6		ifeq	
9	125	aload_1	
10		checkcast	
13		invokeinterface	
18		istore_2	
19	126	goto	
22	127	iconst_m1	
23		istore_2	
24	129	iload_2	
25		istore	
27	131	iinc	
30	129	iload	
32		ireturn	
33	130	astore_3	
34	131	iinc	
37	132	aload_3	
38		athrow	

int tryFinally(Object o) {
 int result = 0;
 try {
 if (o instanceof List<?>) {
 result = ((List<?>) o).size();
 } else {
 result = -1;
 }
 }
 return result;
} finally {
 result += 1;
}

▼ Exception Table:

1.	try [2-27]	catch 33	Any
----	------------	----------	-----

16

The finally block is generally included twice - independent of its size! Once, for the case if no exception is thrown and once for the case when an exception is thrown.


```

static int readFirstByte(String path) throws Exception {
    try (FileReader r = new FileReader(path)) {
        return r.read();
    }
}

```

▼ Method Body (Size: 59 bytes, Max Stack: 3, Max Locals: 4)

PC	Line	Instruction	Exceptions
0	138	aconst_null	
1		astore_1	
2		aconst_null	
3		astore_2	
4		new	3 Any
7		dup	
8		aload_0	
9		invokespecial	java.io.FileReader { void <init> (java.lang.String) }
12		astore_3	
13	139	aload_3	
14		invokevirtual	java.io.FileReader { int read () }
17	140	aload_3	
18		ifnull	26
21		aload_3	
22		invokevirtual	java.io.FileReader { void close () }
25	139	return	
26		astore_1	2 Any
27	140	aload_3	
28		ifnull	35
31		aload_3	
32		invokevirtual	java.io.FileReader { void close () }
35		aload_1	
36		athrow	
37		astore_2	
38		aload_1	
39		thrownull	47
42		aload_2	
43		astore_1	
44		goto	57
47		aload_1	
48		aload_2	
49		if_acmpneq	57
52		aload_1	
53		aload_2	
54		invokevirtual	java.lang.Throwable { void addSuppressed (java.lang.Throwable) }
57		aload_1	
58		athrow	

17

The stream is closed in the regular case, but also if an exception occurs - even if the exception occurs while closing the stream!

A try-with-resources statement can handle the case where the resource is “null”.

```

static int readFirstByte(String path) throws Exception {
    try (FileReader r = new FileReader(path)) {
        return r.read();
    }
}

```

PC	Line	Instruction	Exceptions
0	138	aconst_null	
1		astore_1	
2		aconst_null	
3		astore_2	
4		new	3 Any
7		dup	
8		aload_0	
9		invokespecial	java.io.FileReader { void <init> (java.lang.String) }
12		astore_3	
13	139	aload_3	
14		invokevirtual	java.io.FileReader { int read () }
17	140	aload_3	
18		ifnull	26
21		aload_3	
22		invokevirtual	java.io.FileReader { void close () }
25	139	return	
26		astore_1	2 Any
27	140	aload_3	
28		ifnull	35
31		aload_3	
32		invokevirtual	java.io.FileReader { void close () }
35		aload_1	
36		athrow	
37		astore_2	
38		aload_1	
39		thrownull	47

Exception handlers:

try [13-17) catch 26 Any

try [25-26) catch 26 Any

try [4-37) catch 37 Any

Java Bytecode - Synchronization

```
public class TACDemo {
    private static volatile TACDemo instance;

    static TACDemo getInstance() {
        TACDemo instance = TACDemo.instance;
        // thread-safe double checked locking
        if (instance == null) {
            synchronized (TACDemo.class) {
                instance = TACDemo.instance;
                if (instance == null) {
                    instance = new TACDemo();
                    TACDemo.instance = instance;
                }
            }
        }
        return instance;
    }
}
```

19

PC	Instruction	Exception Handlers
0	getstatic TACDemo { TACDemo instance }	
3	astore_0	
4	aload_0	
5	ifnonnull →41	
8	ldc TACDemo.class	
10	dup	
11	astore_1	
12	monitorenter	
13	getstatic TACDemo { TACDemo instance }	1: Any
16	astore_0	
17	aload_0	
18	ifnonnull →33	
21	new TACDemo	
24	dup	
25	invokespecial TACDemo { void <init> () }	
28	astore_0	
29	aload_0	
30	putstatic TACDemo { TACDemo instance }	
33	aload_1	
34	monitorexit	
35	goto →41	
38	aload_1	2: Any
39	monitorexit	
40	athrow	
41	aload_0	
42	areturn	

synchronized blocks are translated to monitorenter/monitorexit(+) instructions.
synchronized blocks will always result in a try-finally handler to ensure that the lock is released even in case of an exception!

Java Bytecode - Invokedynamic

```
static <T> List<T> sortIt(List<T> l) {
    l.sort(
        (T a, T b) -> { return a.hashCode() - b.hashCode(); }
    );
    return l;
}
```

20

Java Bytecode - Invokedynamic

```
static <T> List<T> sortIt(List<T> l) {
    l.sort((T a, T b) -> { return a.hashCode() - b.hashCode(); });
    return l;
}
```

static java.util.List sortIt(java.util.List)

Signature : <T:Ljava/lang/Object;>(Ljava/util/List<TT>;)Ljava/util/List<TT>;

▼ Method Body (Size: 13 bytes, Max Stack: 2, Max Locals: 1)

PC	Line	Instruction
0	164	aload_0
1		invokedynamic (Bootstrap Method Attribute[0], java.util.Comparator compare ())
6		invokeinterface java.util.List { void sort (java.util.Comparator)
11	167	aload_0
12		areturn

▼ LineNumberTable

start_pc: 0, line_number: 164
start_pc: 11, line_number: 167

▼ LocalVariableTable

pc=[0 → 13] / lv=0 ⇒ java.util.List l

▼ LocalVariableTypeTable

pc=[0 → 13] / lv=0 ⇒ l : Ljava/util/List<TT>;

Lambda expression in Java source code are compile using invoke dynamic instructions.

Some information is optional; e.g., the line number table and also the local variable (type) tables. However, the signature is practically not optional for libraries.

Java Bytecode - Invokedynamic

```
static <T> List<T> sortIt(List<T> l) {
    l.sort((T a, T b) -> { return a.hashCode() - b.hashCode(); });
    return l;
}
```

static java.util.List sortIt(java.util.List)

Signature : <T:Ljava/lang/Object;>(Ljava/util/List<TT>;)Ljava/util/List<TT>;

▼ Method Body (Size: 13 bytes, Max Stack: 2, Max Locals: 1)

PC	Line	Instruction
0	164	aload_0
1		invokedynamic (Bootstrap Method Attribute[0], java.util.Comparator compare ())

MethodHandle(kind=REF_invokeStatic invokestatic C.m:(A*)T, java.lang.invoke.LambdaMetafactory (java.lang.invoke.CallSite metafactory (java.lang.invoke.MethodHandles\$Lookup, java.lang.String, java.lang.invoke.MethodType, java.lang.invoke.MethodType, java.lang.invoke.MethodHandle, java.lang.invoke.MethodType)))

Parameters:

- MethodType(int (java.lang.Object, java.lang.Object))
- MethodHandle(

kind=REF_invokeStatic invokestatic C.m:(A*)T,
 TACDemo { int lambda\$0 (java.lang.Object, java.lang.Object) }
- MethodType(int (java.lang.Object, java.lang.Object))

The Bootstrap Method Attribute [0].

APSA-02-Java Bytecode.key - 28. April 2016

Java Bytecode - Invokedynamic

```
static <T> List<T> sortIt(List<T> l) {
    l.sort((a, b) -> { return a.hashCode() - b.hashCode(); });
    return l;
}
```

static java.util.List sortIt(java.util.List)

Signature : <T>Ljava/lang/Object;-(Ljava/util/List<T>;)Ljava/util/List<T>;

▼ Method Body (Size: 13 bytes, Max Stack: 2, Max Locals: 1)

PC	Line	Instruction
0	164	aload_0
1		invokedynamic (Bootstrap Method Attribute[0], java.util.Comparator compare ())
6		
11	167	
12		

▼ LineNumberTable

start_pc: 0, line_number: 164

▼ LocalVariableTable

pc=[0 → 13] / n=0 ⇒ java.lang.Object a

pc=[0 → 13] / n=1 ⇒ java.lang.Object b

▼ LocalVariableTypeTable

pc=[0 → 13] / n=0 ⇒ a: TT;

pc=[0 → 13] / n=1 ⇒ b: TT;

The “lambda method”.

Analyzing Java Bytecode works (very) well, because
Java compilers deliberately don't optimize.

```
static long optimizableExpression(double d, int i, long l) {
    return (long)((d * d) + i) * 0L;
}
```

PC	Instruction
0	dload_0
1	dload_0
2	dmul
3	iload_2
4	i2d
5	dadd
6	d2l
7	lconst_0
8	lmul
9	lreturn

Example

The result of the expression is always “0” and this could be easily detected, but the compiler won't optimize it. However, the JVM will (most likely) do it.

Analyzing Java Bytecode works (very) well, because
Java compilers deliberately don't optimize.

```
static void optimizableIndexInc(int[] is, int i) {  
    is[i++] = 0;  
    is[i++] = 1;  
}
```

PC	Instruction
0	aload_0
1	iload_1
2	iinc (lv=1,val=1)
5	iconst_0
6	istore
7	aload_0
8	iload_1
9	iinc (lv=1,val=1)
12	iconst_1
13	istore
14	return

25

Example

Here, the second iinc instruction is useless and this could be easily detected, but it is still not optimized.

Analyzing Java Bytecode works (very) well, because
Java compilers deliberately don't optimize.

```
int always9() {  
    int i = 3;  
    int j = 3;  
    return i * j;  
}
```

PC	Instruction
0	iconst_3
1	istore_1
2	iconst_3
3	istore_2
4	iload_1
5	iload_2
6	imul
7	ireturn

26

Example

Analyzing Java Bytecode works (very) well, because
Java compilers deliberately don't optimize.

**Java compilers perform constant propagation for final local
variables and final fields (primitive values and Strings).**
Expressions are evaluated if all parameters are primitive constants.

```
int always6() {  
    final int i = 3;  
    final int j = 2;  
    return i * j;  
}
```

PC	Instruction
0	iconst_3
1	istore_1
2	iconst_2
3	istore_2
4	bipush 6
6	ireturn

```
int always8() {  
    return 4 * 2;  
}
```

PC	Instruction
0	bipush 8
2	ireturn

```
double eXpi() { return Math.E * Math.PI; }
```

PC	Instruction
0	ldc2_w 8.539734222673566d
3	dreturn

27

Example

Summary

Sources of Potentially Dead Code Created by Java Compilers

- Finally blocks are generally included twice.
- Switches always have default branches.
- Constant expressions are evaluated.
- Constant propagation for final (local) variables/final fields is performed. (This includes primitive types and "Strings").

28

Other Peculiarities

- Types are represented using binary notation. In binary notation packages are separated using "/": e.g., java/lang/Object.
- The JVM has no "negate" instruction. A negation in Java (!b) is compiled to an if instruction followed by a push of the corresponding value.
- The JVM has no direct support for shortcut-evaluation (&&, ||).
- (Reliably) identifying anonymous inner classes is broken for older class files.
- The switch instructions are four byte aligned.
- The instruction set is not orthogonal; i.e., to achieve a certain effect many instructions exist.
- The catch block is not immediately available.

29

Additional differences compared to Java are: String concatenation is compiled using StringBuilders and Varargs are always stored in arrays.

Three-Address Code

Three-Address Code (TAC)

- Three-address code is a sequence of statements with the general form:
`x = y op z`
- where x,y and z are (local variable) names, constants (in case of y and z) or compiler-generated temporaries
- The name was chosen, because “most” statements use three addresses: two for the *operators* and *one to store the result*

31

General Types of Three-Address Statements

- **Assignment** statements `x = y bin_op z` or `x = unary_op z`
- **Copy** statements `x = y`
- **Unconditional jumps**: `goto l` (and `jsr l`, `ret` in case of Java bytecode)
- **Conditional jumps**: `if (x rel_op y) goto l` (else fall through), `switch`
- **Method call and return**: `invoke(m, params)`, `return x`
- **Array access**: `a[i]` or `a[i] = x`
- *More IR specific types.*

32

Converting Java Bytecode to Three-Address Code

- Core Idea:

- Compute for each instruction the current stack layout by following the control flow; i.e., compute the types of values found on the stack before the instruction is evaluated.
(The JVM specification guarantees that the operand stack always has the same layout independent of the taken path.)
- Assign each local variable to a variable where the name is based on the local variable index.
E.g., an `inc(lv=1, val=1)` instruction is transformed into the three address code: `r_1 = r_1 + 1`
- Assign each variable on the operand stack to a corresponding local variable with an index based on the position on the stack.
E.g., if the operand stack is empty and we push the constant 1, then the three address code would be: `op_0 = 1`; if we would then push another value 2 then the code would be: `op_1 = 2` and an addition of the two values would be: `op_0 = op_0 + op_1`

33

Converting Java Bytecode to Three-Address Code

```
static int numberOfDigits(int i) {
    return ((int) Math.floor(Math.log10(i))) + 1;
}
```

PC	Instruction	Stack Layout (before execution)	Three-Address Code
	<i>Initialization</i>		<code>r_0 = i; // parameter</code>
0	<code>iload_0</code>		<code>op_0 = r_0;</code>
1	<code>i2d</code>	0: Integer Value	<code>op_0 = (double) op_0;</code>
2	<code>invokestatic java.lang.Math.log10 (double):double</code>	0: Double Value	<code>op_0 = Math.Log10(op_0);</code>
5	<code>invokestatic java.lang.Math.floor(double):double</code>	0: Double Value	<code>op_0 = Math.floor(op_0);</code>
8	<code>d2i</code>	0: Double Value	<code>op_0 = (int) op_0;</code>
9	<code>iconst_1</code>	0: Integer Value	<code>op_1 = 1;</code>
10	<code>iadd</code>	1: Integer Value 0: Integer Value	<code>op_0 = op_0 + op_1;</code>
11	<code>ireturn</code>	0: Integer Value	<code>return op_0;</code>

34

Java Bytecode vs. Three-Address Code

Java	Bytecode	Three Address Code																
<pre>static int max(int i, int j) { if (i > j) return i; else return j; }</pre>	<table border="1"> <thead> <tr> <th>PC</th> <th>Instruction</th> </tr> </thead> <tbody> <tr><td>0</td><td>iload_0</td></tr> <tr><td>1</td><td>iload_1</td></tr> <tr><td>2</td><td>if_icmple →7</td></tr> <tr><td>5</td><td>iload_0</td></tr> <tr><td>6</td><td>ireturn</td></tr> <tr><td>7</td><td>iload_1</td></tr> <tr><td>8</td><td>ireturn</td></tr> </tbody> </table>	PC	Instruction	0	iload_0	1	iload_1	2	if_icmple →7	5	iload_0	6	ireturn	7	iload_1	8	ireturn	<pre>0: r_0 = i; 1: r_1 = j; 2: op_0 = r_0; 3: op_1 = r_1; 4: if(op_0 <= op_1) goto 7; 5: op_0 = r_0; 6: return op_0; 7: op_0 = r_1; 8: return op_0;</pre>
PC	Instruction																	
0	iload_0																	
1	iload_1																	
2	if_icmple →7																	
5	iload_0																	
6	ireturn																	
7	iload_1																	
8	ireturn																	
	<p>After (Peephole) Optimizations</p> <p>↓</p> <pre>0: if(i <= j) goto 2; 1: return i; 2: return j;</pre>																	

35

Peephole optimizations use a “sliding window” over the cfg’s basic blocks (discussed later) to perform, e.g., the following optimizations:

- elimination of redundant loads and stores
- constant folding
- constant propagation
- common subexpression elimination
- copy propagation
- strength reduction ($x * 2 \Rightarrow x + x$; $x / 2 = x \gg 1$)
- elimination of useless instructions ($y = x * 0 \Rightarrow y = 0$; **may cause false negatives!**)
- *exploitation of the instruction set (not really possible for Java bytecode)*

Java Bytecode vs. Three-Address Code

Java	Bytecode	Three Address Code																										
<pre>static int factorial(int n) { int r = 1; while (n > 0) { r *= n; n--; } return r; }</pre>	<table border="1"> <thead> <tr> <th>PC</th> <th>Instruction</th> </tr> </thead> <tbody> <tr><td>0</td><td>iconst_1</td></tr> <tr><td>1</td><td>istore_1</td></tr> <tr><td>2</td><td>goto →12</td></tr> <tr><td>5</td><td>iload_1</td></tr> <tr><td>6</td><td>iload_0</td></tr> <tr><td>7</td><td>imul</td></tr> <tr><td>8</td><td>istore_1</td></tr> <tr><td>9</td><td>iinc 0, -1</td></tr> <tr><td>12</td><td>iload_0</td></tr> <tr><td>13</td><td>ifgt →5</td></tr> <tr><td>16</td><td>iload_1</td></tr> <tr><td>17</td><td>ireturn</td></tr> </tbody> </table>	PC	Instruction	0	iconst_1	1	istore_1	2	goto →12	5	iload_1	6	iload_0	7	imul	8	istore_1	9	iinc 0, -1	12	iload_0	13	ifgt →5	16	iload_1	17	ireturn	<pre>0: r_0 = n; 1: op_0 = 1; 2: r_1 = op_0; 3: goto 9; 4: op_0 = r_1; 5: op_1 = r_0; 6: op_0 = op_0 * op_1; 7: r_1 = op_0; 8: r_0 = r_0 + -1; 9: op_0 = r_0; 10: if(op_0 > 0) goto 4; 11: op_0 = r_1; 12: return op_0;</pre> <p>After (Peephole) Optimizations</p> <p>↓</p> <pre>0: r_0 = n; 1: r_1 = 1; 2: goto 5; 3: r_1 = r_1 * r_0; 4: r_0 = r_0 + -1; 5: if(r_0 > 0) goto 3; 6: return r_1;</pre>
PC	Instruction																											
0	iconst_1																											
1	istore_1																											
2	goto →12																											
5	iload_1																											
6	iload_0																											
7	imul																											
8	istore_1																											
9	iinc 0, -1																											
12	iload_0																											
13	ifgt →5																											
16	iload_1																											
17	ireturn																											

36

A Three-Address Code for Java Bytecode

Basic Statements

```
ASSIGNMENT(  
    pc: PC,  
    targetVar: VAR,  
    expr: VAL_EXPR  
)  
  
GOTO(pc: PC, target: Int)  
  
JUMP_TO_SUBROUTINE(pc: PC, target: Int)  
RET(pc: PC, returnAddressVar: VAR)  
  
NOP(pc: PC)  
  
IF(  
    pc: PC,  
    left: VAL_EXPR,  
    condition: RelationalOperator,  
    right: VAL_EXPR,  
    target: Int  
)  
  
SWITCH(  
    pc: PC,  
    defaultTarget: PC,  
    index: VAL_EXPR,  
    npairs: IndexedSeq[(Int, PC)]  
)  
  
RETURN_VALUE(pc: PC, expr: VAL_EXPR)  
RETURN(pc: PC)  
  
ARRAY_STORE(  
    pc: PC,  
    arrayRef: VAR,  
    index: VAL_EXPR,  
    value: VAL_EXPR  
)  
  
THROW(pc: PC, exception: VAR)  
  
MONITOR_ENTER(pc: PC, objRef: VAR)  
MONITOR_EXIT(pc: PC, objRef: VAR)
```

37

A Three-Address Code for Java Bytecode

Field Access and Method Call Statements

```
PUT_STATIC(  
    pc: PC,  
    declaringClass: ObjectType,  
    name: String,  
    value: VAL_EXPR  
)  
  
PUT_FIELD(  
    pc: PC,  
    declaringClass: ObjectType,  
    name: String,  
    objRef: VAR,  
    value: VAL_EXPR  
)  
  
NON_VIRTUAL_METHOD_CALL(  
    pc: PC,  
    declaringClass: ReferenceType,  
    name: String,  
    descriptor: MethodDescriptor,  
    receiver: VAR,  
    params: List[VAL_EXPR]  
)  
  
VIRTUAL_METHOD_CALL(  
    pc: PC,  
    declaringClass: ReferenceType,  
    name: String,  
    descriptor: MethodDescriptor,  
    receiver: VAR,  
    params: List[VAL_EXPR]  
)  
  
STATIC_METHOD_CALL(  
    pc: PC,  
    declaringClass: ReferenceType,  
    name: String,  
    descriptor: MethodDescriptor,  
    params: List[VAL_EXPR]  
)
```

38

A Three-Address Code for Java Bytecode

Expressions

```
INSTANCEOF(pc: PC, value: VAR, t: ReferenceType)
CHECKCAST(pc: PC, value: VAR, t: ReferenceType)
NEW(pc: PC, tpe: ObjectType)
NEWARRAY(
  pc: PC,
  counts: List[VALEXPR],
  tpe: ArrayType
)
ARRAYLOAD(pc: PC, index: Var, arrayRef: Var)
ARRAYLENGTH(pc: PC, arrayRef: Var)
VALEXPR
PARAM(cTpe: ComputationalType, name: String)
SIMPLEVAR(id: Int, cTpe: ComputationalType)
INTCONST(pc: PC, value: Int)
LONGCONST(pc: PC, value: Long)
FLOATCONST(pc: PC, value: Float)
DOUBLECONST(pc: PC, value: Double)
STRINGCONST(pc: PC, value: String)
CLASSCONST(pc: PC, value: ReferenceType)
NULLEXPR(pc: PC)

COMPARE(
  pc: PC,
  left: VALEXPR,
  condition: RelationalOperator,
  right: VALEXPR
)
BINARYEXPR(
  pc: PC,
  cTpe: ComputationalType,
  op: BinaryArithmeticOperator,
  left: VALEXPR, right: VALEXPR
)
PREFIXEXPR(
  pc: PC,
  cTpe: ComputationalType,
  op: UnaryArithmeticOperator,
  operand: VALEXPR
)
PRIMITIVEYPECASTEXPR(
  pc: PC,
  targetType: BaseType,
  operand: VALEXPR
)
```

39

A Three-Address Code for Java Bytecode

Expressions

```
GETFIELD(
  pc: PC,
  declaringClass: ObjectType,
  name: String,
  objRef: VALEXPR
)
GETSTATIC(
  pc: PC,
  declaringClass: ObjectType,
  name: String
)
INVOKEDYNAMIC(
  pc: PC,
  bootstrapMethod: BootstrapMethod,
  name: String,
  descriptor: MethodDescriptor,
  params: List[EXPR]
)
METHODTYPECONST(pc: PC, desc: MethodDescriptor)
METHODHANDLECONST(pc: PC, desc: MethodHandle)
NONVIRTUALFUNCTIONCALL(
  pc: PC,
  declaringClass: ReferenceType,
  name: String,
  descriptor: MethodDescriptor,
  receiver: EXPR,
  params: List[EXPR]
)
VIRTUALFUNCTIONCALL(
  pc: PC,
  declaringClass: ReferenceType,
  name: String,
  descriptor: MethodDescriptor,
  receiver: EXPR,
  params: List[EXPR]
)
STATICFUNCTIONCALL(
  pc: PC,
  declaringClass: ReferenceType,
  name: String,
  descriptor: MethodDescriptor,
  params: List[EXPR]
)
```

40

Control-Flow Graph

- The control-flow graph (CFG) represents the control flow of a single method.
- Each node represents a **basic block**. A basic block is a maximal-length sequence of statements without jumps in and out (and no exceptions are thrown by intermediate instructions).
- The arcs represent the inter-node control flow.

41

Control-Flow Graph

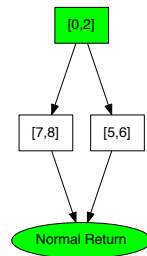
Java

```
static int max(int i, int j) {  
    if (i > j)  
        return i;  
    else  
        return j;  
}
```

Bytecode

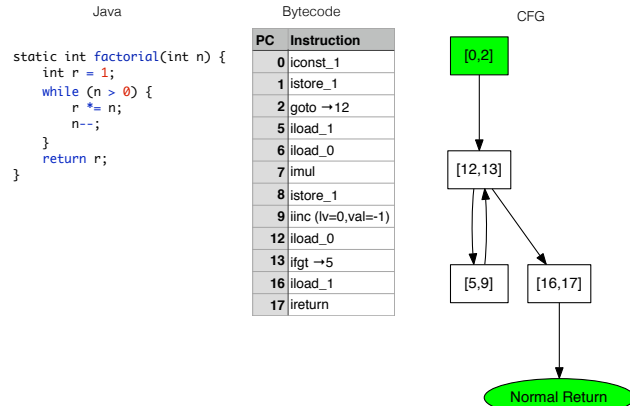
PC	Instruction
0	iload_0
1	iload_1
2	if_icmple →7
5	iload_0
6	ireturn
7	iload_1
8	ireturn

CFG



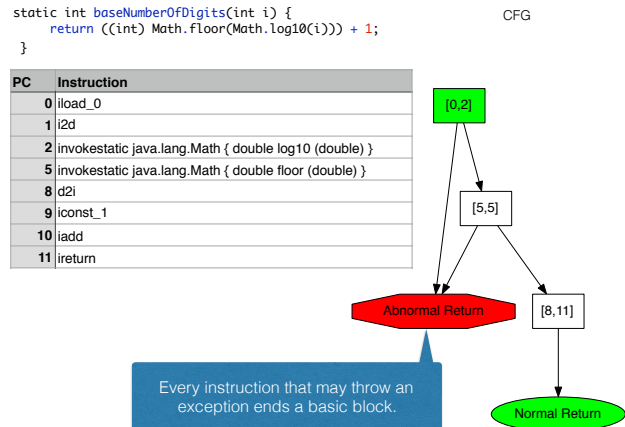
42

Loops - Control-Flow Graph



43

Java Bytecode



44

In general, the basic blocks in Java bytecode are rather small; many bytecode instruction can throw an exception.

Java Bytecode - Finally

```
int tryFinally(Object o) {
    int result = 0;
    try {
        if (o instanceof List<?>) {
            result = ((List<?>) o).size();
        } else {
            result = -1;
        }
        return result;
    } finally {
        result += 1;
    }
}
```

▼ Method Body (Size: 39 bytes, Max Stack: 1, Max Locals: 5)

PC	Line	Instruction	Exceptions
0	122	iconst_0	
1		istore_2	
2	124	aload_1	
3		instanceof java.util.List	1: Any
6		ifeq 22	
9	125	aload_1	
10		checkcast java.util.List	
13		invokeinterface java.util.List int size ()	
18		istore_2	
19	126	goto 24	
22	127	iconst_m1	
23		istore_2	
24	129	iload_2	
25		istore 4	
27	131	iinc 2 1	
30	129	iload 4	
32		ireturn	
33	130	astore_3	
34	131	iinc 2 1	
37	132	aload_3	
38		athrow	

▼ Exception Table:

1. try [2-27] catch 33 Any

45

When we reach an exception handler the operand stack always only contains the thrown exception.

Java Bytecode - Finally

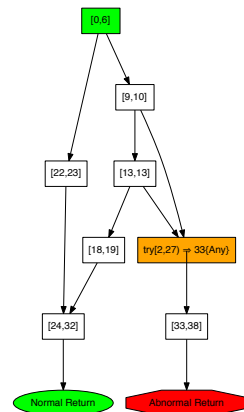
▼ Method Body (Size: 39 bytes, Max Stack: 1, Max Locals: 5)

PC	Line	Instruction	Exceptions
0	122	iconst_0	
1		istore_2	
2	124	aload_1	
3		instanceof java.util.List	1: Any
6		ifeq 22	
9	125	aload_1	
10		checkcast java.util.List	
13		invokeinterface java.util.List int size ()	
18		istore_2	
19	126	goto 24	
22	127	iconst_m1	
23		istore_2	
24	129	iload_2	
25		istore 4	
27	131	iinc 2 1	
30	129	iload 4	
32		ireturn	
33	130	astore_3	
34	131	iinc 2 1	
37	132	aload_3	
38		athrow	

▼ Exception Table:

1. try [2-27] catch 33 Any

46



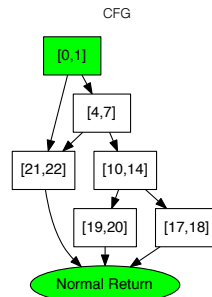
The method call (List.size) now goes to the exception handler.

The CFG is an over approximation; it represents all paths.

```
static boolean checkIt(int a, int b) {
    if (b < 0) {
        if (a < 100) {
            if (b > 1000)
                return true;
            else
                return false;
        }
    }
    return true;
}
```

Bytecode

PC	Instruction
0	iload_1
1	ifge21
4	iload_0
5	bipush 100
7	if_icmpge →21
10	iload_1
11	sipush 1000
14	if_icmple →19
17	iconst_1
18	ireturn
19	iconst_0
20	ireturn
21	iconst_1
22	ireturn



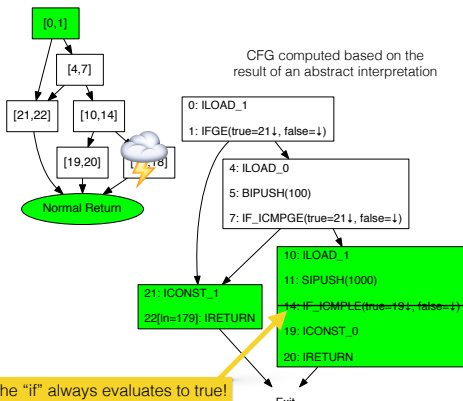
When we perform a decent data-flow analysis, we can determine that the third if condition will always evaluate to false and, hence, the edge to the basic block starting with instruction 17 is dead.

The CFG is an over approximation.

Bytecode

PC	Instruction
0	iload_1
1	ifge21
4	iload_0
5	bipush 100
7	if_icmpge →21
10	iload_1
11	sipush 1000
14	if_icmple →19
17	iconst_1
18	ireturn
19	iconst_0
20	ireturn
21	iconst_1
22	ireturn

Conservative CFG



The "if" always evaluates to true!

We can now conclude that the instructions 17 and 18 are dead.

(Hence, every node dominates itself.)

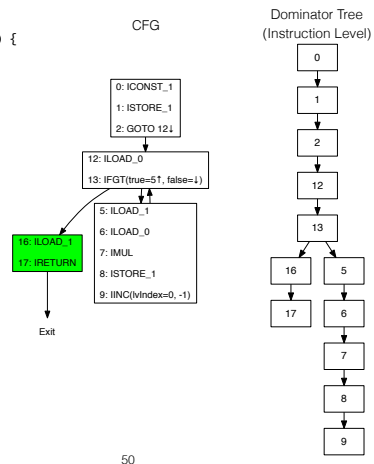
- A node *d* of a (control-) flow graph *dominates* node *n*, if every path from the *initial node* of the flow graph to *n* goes through *d*.
(Every node dominates itself.)
- In a dominator tree the initial node is the *unique* root node and each node *d* only dominates its descendants.
- Dominator information is useful in identifying loops;
by identifying the header and the back edge.

49

Loops - CFG and Dominator Tree

```
Java
static int factorial(int n) {
    int r = 1;
    while (n > 0) {
        r *= n;
        n--;
    }
    return r;
}
```

PC	Instruction
0	iconst_1
1	istore_1
2	goto →12
5	iload_1
6	iload_0
7	imul
8	istore_1
9	iinc (lv=0, val=-1)
12	iload_0
13	ifgt →5
16	iload_1
17	ireturn



Given the dominator tree it is then possible to compute various other helpful data structure: Reverse Dominator Tree, Dominance Frontier, Control Dependence Graph...

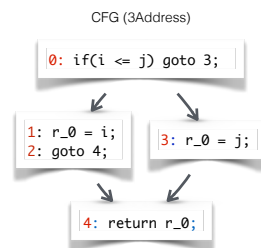
Post-Dominator Tree

- The post-dominator tree is the dominator tree computed using the *reverse* control flow graph.
- The reverse control flow graph is the control flow graph where all edges are reversed and — *if necessary* — a new artificial unique root node is added.
- The post-dominator tree is used, e.g., to determine control dependency.

51

Post-Dominator Tree

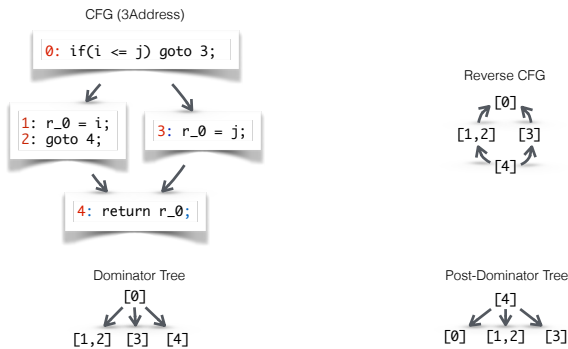
```
Java
static int max(int i, int j) {
    int max;
    if (i > j)
        max = i;
    else
        max = j;
    return max;
}
```



52

(Hence, every node dominates itself.)

Post-Dominator Tree



53

(Hence, every node dominates itself.)

Control-Dependence

- An instruction/statement is control dependent on a predicate if the value of the predicate controls the execution of the instruction.

- Let G be a control flow graph; Let X and Y be nodes in G ; Y is control dependent on X iff
 - there exists a directed path P from X to Y with any Z in $P \setminus \{X, Y\}$ post-dominated by Y
 - X is not post-dominated by Y

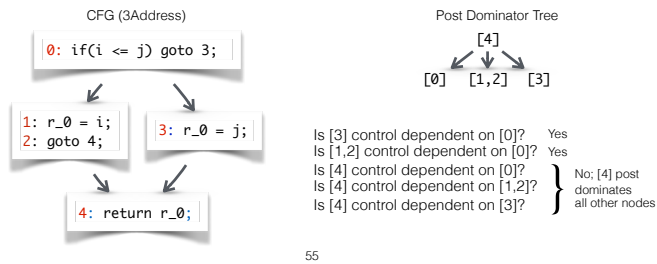
54

(Paper: Jeanne Ferrante, Karl J. Ottenstein, and Joe D. Warren. 1987. The program dependence graph and its use in optimization. ACM Trans. Program. Lang. Syst. 9, 3 (July 1987), 319-349. DOI=<http://dx.doi.org/10.1145/24039.24041>)

First rule in other words: starting with X it is not possible to bypass Y . Second rule in other words: X will not be executed anyway.

Control-Dependence

- An instruction/statement is control dependent on a predicate if the value of the predicate controls the execution of the instruction.



This is a common scenario; e.g., used by java.lang.String.

Identification of Lazily Initialized Fields

(,e.g. to facilitate the identification of immutable classes)
Control-Dependence

```

class S(final val i: Int, final val j: Int) {
  private[this] var hash: Int = _

  override def hashCode(): Int = {
    if (hash == 0) {
      hash = i*31+j
    }
    hash
  }
}
  
```

hash is updated if and only if hash still has the default value; after that hash is never updated again (unless the computation's result was the default value).

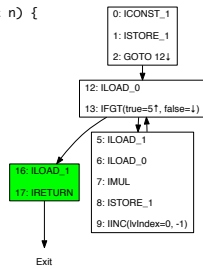
The update is thread-safe.

56

Def-Use Dependencies

Java CFG Def-Use Information

```
static int factorial(int n) {
    int r = 1;
    while (n > 0) {
        r *= n;
        n--;
    }
    return r;
}
```



PC	Instruction	Used By	Uses
0	iconst_1	{7,17}	
1	istore_1	N/A	
2	goto →12		
5	iload_1		
6	iload_0		
7	imul	{7,17}	1. Operand (-1,9) 2. Operand (0,7)
8	istore_1		
9	inc (lv=0, val=-1)	{7,9,13}	
12	iload_0		
13	ifgt →5		{-1,9}
16	iload_1		
17	ireturn		{0,7}

Here "-1" is used to indicate a usage of the first parameter.

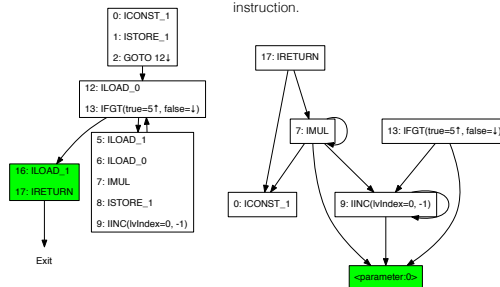
57

Loops - Def-Use Dependencies

Java CFG Implicit Def-Use dependencies at the level of instructions.

```
static int factorial(int n) {
    int r = 1;
    while (n > 0) {
        r *= n;
        n--;
    }
    return r;
}
```

PC	Instruction
0	iconst_1
1	istore_1
2	goto →12
5	iload_1
6	iload_0
7	imul
8	istore_1
9	inc (lv=0, val=-1)
12	iload_0
13	ifgt →5
16	iload_1
17	ireturn



58

SSA Code (STATIC SINGLE ASSIGNMENT (FORM))

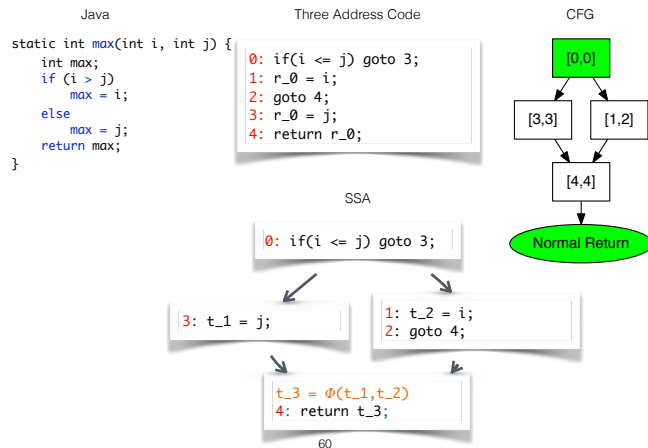
- Motivation: Many analyses require def-use information. I.e., the information where a used local variable is defined or vice versa.
- In SSA form each variable *has only one static definition-site in the program text*.
- When two control-flow paths merge, a selector function ϕ is used that initializes the variable based on the control flow that was taken.
- SSA form facilitates data-flow analyses; e.g., facilitates the elimination of redundant loads and computations across basic block boundaries.

59

Here, “static” refers to the definition site - it may be in a loop and may be executed multiple times.

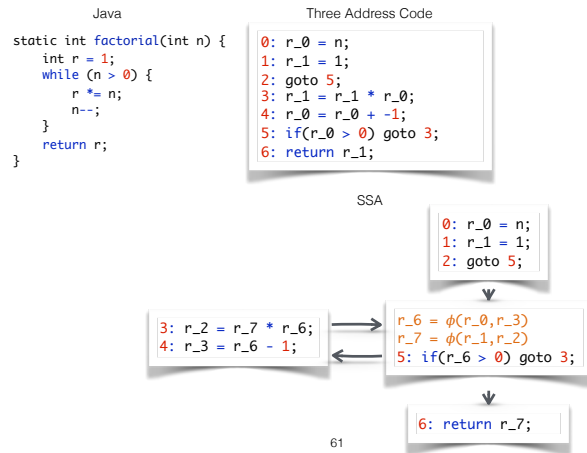
SSA is an improvement of def-use chains.

Java Bytecode vs. Three-Address Code vs. SSA



60

Java Bytecode vs. Three-Address Code vs. SSA



Class Hierarchy

- A core data-structure which encodes the project's class hierarchy and which offers query functionality to get a type's supertypes and subtypes.
- In Java `java.lang.Object` is the super type of all types (including interface types). I.e., every interface (in the byte code explicitly) inherits from `Object`.
- In practice it is - when you want to analyze libraries - *basically always* the case that the class hierarchy contains "holes" (is not upwards closed).



62