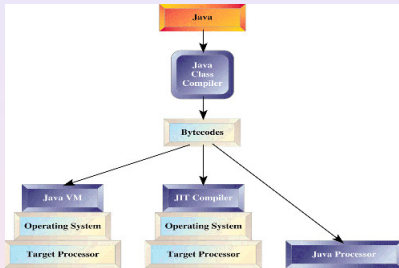
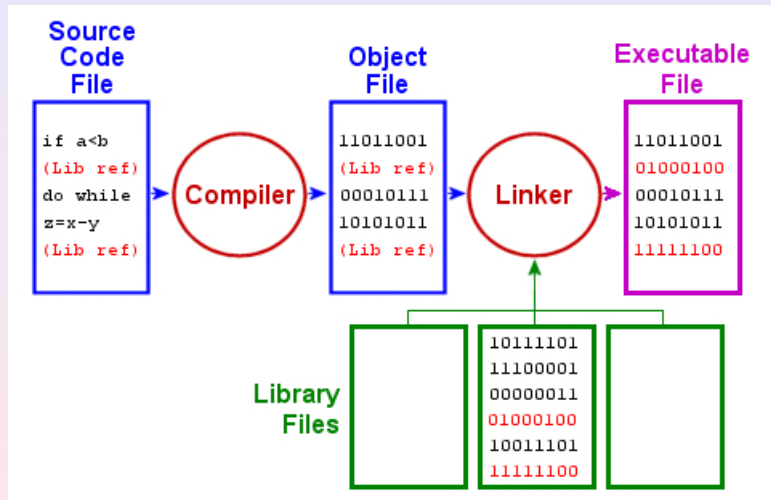
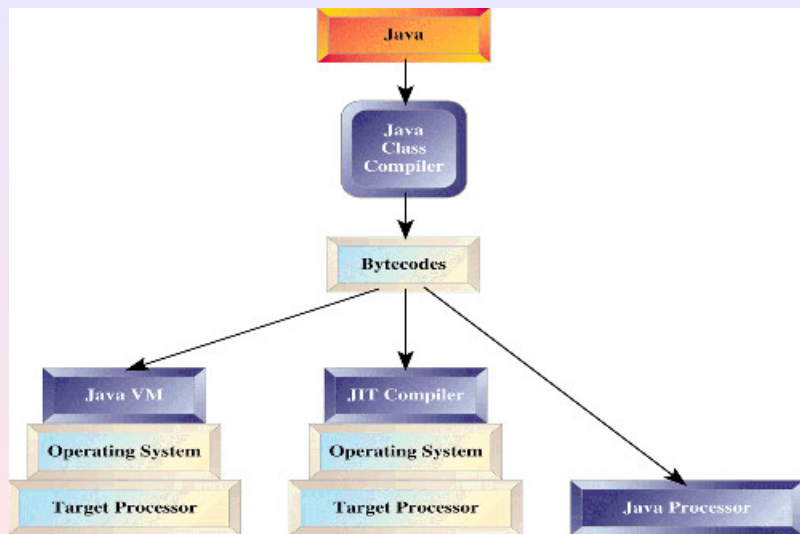

Short Notes on Compiler Construction



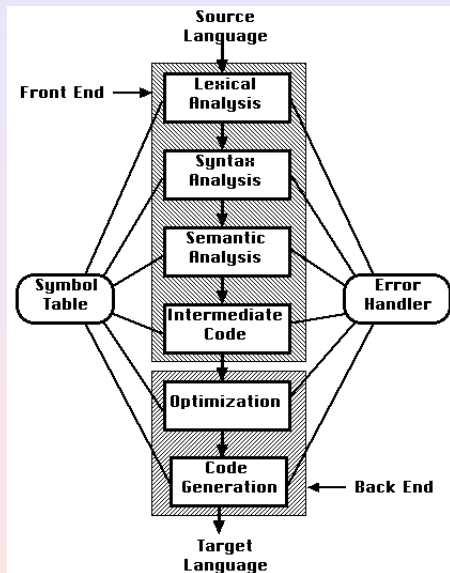
Compilation into native code



Compilation through an intermediate bytecode



Compilation phases



My open source Kitten compiler:
<https://github.com/HotMoka/Kitten>

Lexical analysis

- **input:** textual source code
- **output:** sequence of tokens
- **error:** unknown token

Input: Kitten source code

```
class Led {  
    field boolean state  
  
    constructor() {}  
  
    method void on()  
        this.state := true  
  
    method void off()  
        this.state := false  
  
    method boolean isOn()  
        return this.state  
  
    method boolean isOff()  
        return !this.state  
}
```

Output: sequence of tokens

CLASS from 0 to 4
ID(Led) from 6 to 8
LBRACE from 10 to 10
FIELD from 14 to 18
BOOLEAN from 20 to 26
ID(state) from 28 to 32
CONSTRUCTOR from 37 to 47
LPAREN from 48 to 48
RPAREN from 49 to 49
LBRACE from 51 to 51
RBRACE from 52 to 52
METHOD from 57 to 62
VOID from 64 to 67
ID(on) from 69 to 70
...

Old friends: Regular expressions

An *alphabet* Λ is a finite collection of symbols (*characters*)

Regular Expressions \mathcal{R} over Λ

- $\emptyset \in \mathcal{R}$ (empty set)
- $\varepsilon \in \mathcal{R}$ (empty string)
- $\Lambda \subseteq \mathcal{R}$ (single characters)
- if $r_1, r_2 \in \mathcal{R}$ then $r_1 r_2 \in \mathcal{R}$ (sequence)
- if $r_1, r_2 \in \mathcal{R}$ then $r_1 | r_2 \in \mathcal{R}$ (union)
- if $r \in \mathcal{R}$ then $r^* \in \mathcal{R}$ (iteration)

A regular expression denotes a language

- $\mathcal{L}(\emptyset) = \emptyset$
- $\mathcal{L}(\varepsilon) = \{\varepsilon\}$
- $\mathcal{L}(a) = \{a\}$ for every $a \in \Lambda$
- $\mathcal{L}(r_1 r_2) = \{s_1 s_2 \mid s_1 \in \mathcal{L}(r_1) \text{ and } s_2 \in \mathcal{L}(r_2)\}$
- $\mathcal{L}(r_1 | r_2) = \mathcal{L}(r_1) \cup \mathcal{L}(r_2)$
- $\mathcal{L}(r^*) = \{s_1 \cdots s_n \mid n \geq 0 \text{ and } s_i \in \mathcal{L}(r) \text{ for all } 0 \leq i \leq n\}$.

For instance

- `then` denotes the keyword `{then}`
- `[a-zA-Z][a-zA-Z0-9_]` denotes the identifiers

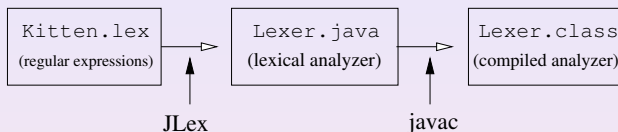
What you can not do with regular expressions

- you cannot define recursive languages
- you cannot count (match parentheses)

C, Java, Kitten source code is defined in a recursive way and is correct only if parentheses match

Building a lexical analyzer with JLex

<https://www.cs.princeton.edu/~appel/modern/java/JLex/>



then	{return tok(THEN, null);}
+	{return tok(PLUS, null);}
:=	{return tok(ASSIGN, null);}
[a-zA-Z][a-zA-Z0-9_]*	{return tok(ID, yytext());}
[0-9]+	{return tok(INTEGER,
	new Integer(yytext()));}
[0-9]*.[0-9]+	{return tok(FLOATING,
	new Float(yytext()));}

- 1 regular expressions \Rightarrow NFA
- 2 NFA \Rightarrow DFA
- 3 DFA \Rightarrow `Lexer.java`

Disambiguation rules

- longest match
- rule priority

Syntactical analysis

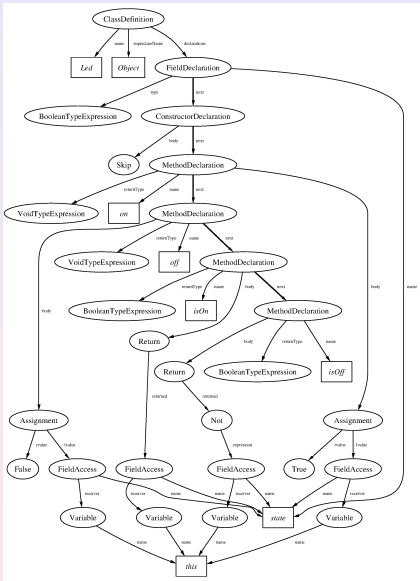
Syntactical analysis

- **input:** sequence of tokens
- **output:** abstract syntax tree
- **error:** syntax error

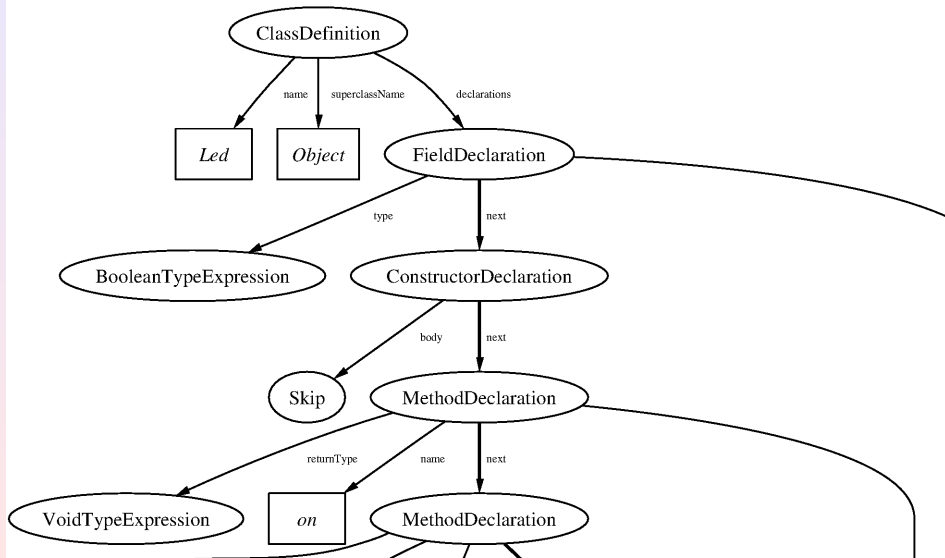
Input: sequence of tokens (yes, you have seen this already)

CLASS from 0 to 4
ID(Led) from 6 to 8
LBRACE from 10 to 10
FIELD from 14 to 18
BOOLEAN from 20 to 26
ID(state) from 28 to 32
CONSTRUCTOR from 37 to 47
LPAREN from 48 to 48
RPAREN from 49 to 49
LBRACE from 51 to 51
RBRACE from 52 to 52
METHOD from 57 to 62
VOID from 64 to 67
ID(on) from 69 to 70
...

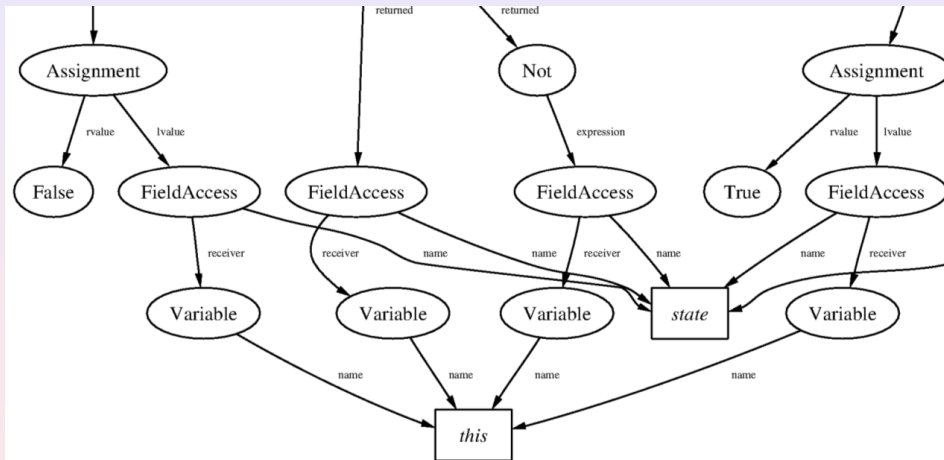
Abstract syntax *tree*



A closer look 1/2



A closer look 2/2



Context-free grammars

A *context-free grammar* over an alphabet Λ is a quadruple $\langle T, N, I, P \rangle$ where

- $T \subseteq \Lambda$ is a set of *terminals*
- N is a set of *non-terminals*
- $I \in N$ is the *starting non-terminal*
- P is a set of *productions*, that is, arrows such as $L \rightarrow \gamma$ where $L \in N$ and $\gamma \in (T \cup N)^*$.

For instance:

$$T = \{a, b\}$$

$$N = \{I\}$$

$$P = \{I \rightarrow \varepsilon, I \rightarrow a/b\}$$

Derivations

Given $G = \langle T, N, I, P \rangle$ we say that β is *derived in G in a step* from α iff there is $L \rightarrow \gamma \in P$ such that $\alpha = \eta L \delta$ and $\beta = \eta \gamma \delta$. We write it as $\alpha \Rightarrow \beta$. A *derivation* for G is a sequence of steps $\alpha \Rightarrow \beta_1 \Rightarrow \beta_2 \dots$

For instance:

$$ab/b \Rightarrow aba/bb$$

$$aba/bb \Rightarrow ababb$$

$$I \Rightarrow a/b$$

$$I \Rightarrow^* I$$

$$ab/b \Rightarrow^* ababb.$$

Language of a grammar

Given a grammar G on an alphabet Λ , its language $L(G)$ is

$$L(G) = \{\alpha \text{ ground} \mid I \Rightarrow^* \alpha\}.$$

For instance, for the previous example of grammar:

$$L(G) = \{a^n b^n \mid n \geq 0\}.$$

Parse trees

A *parse tree* for $G = \langle T, N, I, P \rangle$ is a tree such that

- 1 its nodes are labeled with an element from N or from T or with ε ;
- 2 its root is labeled with I
- 3 its leaves are labeled with elements from T or with ε
- 4 for every node labeled with L and its children labeled with e_1, \dots, e_n (left-to-right) we have $L \rightarrow e_1 \cdots e_n \in P$.

The frontier of the tree is the string derived by the tree from its root.

A parse tree stands for more derivations, by abstracting away the order or replacement of the non-terminals

Ambiguity 1/2

Consider the grammar

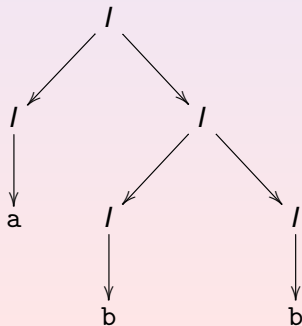
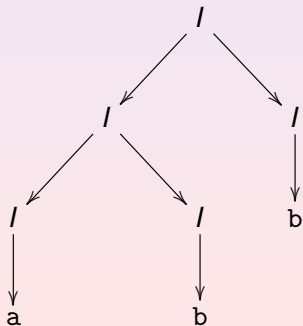
$$I \rightarrow \varepsilon$$

$$I \rightarrow a$$

$$I \rightarrow b$$

$$I \rightarrow II$$

It admits two parse trees for the same word abb:



Why ambiguity is bad

- Ambiguity means that the same word can be given two structures, that is, potentially, two distinct interpretations
- With an ambiguous grammar, a compiler does not know which is the right interpretation of a program
- In a perfect world, grammars should be non-ambiguous but this often entails that they become complex and unnatural
- We will give a concrete example later

Recursive descent parsing

$I \rightarrow com \$$

$com \rightarrow exp \text{ ASSIGN INTEGER}$

$exp \rightarrow \text{INTEGER}$

$exp \rightarrow \text{ID}$

$exp \rightarrow \text{MINUS } exp$

We can implement it with a recursive descent parser:

```
public void parse() { parseI(); }  
private void parseI() { parseCom(); eat(EOF); }  
private void parseCom() { parseExp(); eat(ASSIGN); eat(INTEGER); }  
private void parseExp() {  
    switch (lookahead) {  
        case ID: eat(ID); break;  
        case INTEGER: eat(INTEGER); break;  
        case MINUS: eat(MINUS); parseExp();  
        default: syntax_error(lookahead);  
    }  
}
```


When does recursive descent parsing fail?

Recursive descent parsing fails

- when the first token is not enough to distinguish two productions for the same non-terminal
- when the grammar is left-recursive

For instance it fails for:

$$I \rightarrow A\$$$

$$A \rightarrow a$$

$$A \rightarrow Aa$$

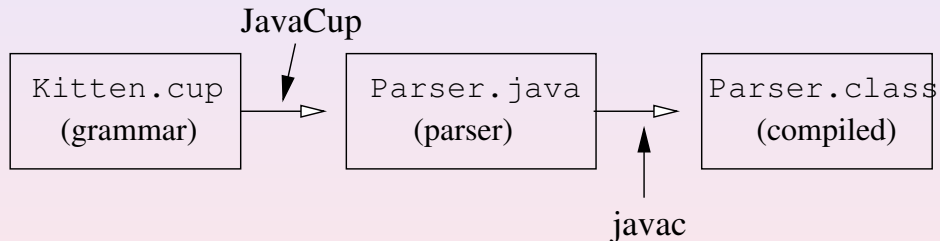
In practice, it fails for the grammars of all sensible programming languages

LR parsing uses a stack automaton to remember what has already been seen during parsing and which are the possible productions that we can use in the future, depending on the input token that will get processed by the compiler:

- details are long and boring: see any compilers book
- the construction of the parser is automatic from the grammar
- a special case of *LR* parsing is the standard approach for compiler construction
- it will anyway fail for ambiguous grammars

Automatic parser generation with JavaCup

<http://www2.cs.tum.edu/projects/cup/>



The grammar of Kitten for JavaCup: types

```
type ::=
    ID
  | BOOLEAN
  | INT
  | FLOAT
  | ARRAY OF type ;
```

```
typeplus ::=
    type
  | VOID ;
```

The grammar of Kitten for JavaCup: leftvalues and expressions

```
lvalue ::=
    ID                      // variable
  | exp DOT ID              // a field of an object
  | exp LBRACK exp RBRACK ; // an element of an array

exp ::=
    lvalue                  // a leftvalue
  | TRUE                    // literals...
  | FALSE
  | INTEGER
  | FLOATING
  | STRING
  | NIL
  | NEW ID LPAREN expseq RPAREN // object creation
  | NEW type LBRACK exp RBRACK  // array creation
```

The grammar of Kitten for JavaCup: more expressions

```
| exp AS type           // cast into type
| exp PLUS exp          // arithmetic....
| exp MINUS exp         // AMBIGUITY
| exp TIMES exp
| exp DIVIDE exp
| MINUS exp             // unary minus
| exp LT exp            // comparisons...
| exp LE exp            // AMBIGUITY
| exp GT exp
| exp GE exp
| exp EQ exp
| exp NEQ exp
| exp AND exp           // logical operations...
| exp OR exp            // AMBIGUITY
| NOT exp
| exp DOT ID LPAREN expseq RPAREN // method call
| LPAREN exp RPAREN ;    // parentheses
```

The grammar of Kitten for JavaCup: resolving the ambiguity

Ambiguity can be resolved with a non-ambiguous grammar (complex, innatural) or with *precedence and associativity* rules:

```
precedence left AND, OR;  
precedence left NOT;  
precedence nonassoc EQ, NEQ, LT, LE, GT, GE;  
precedence left PLUS, MINUS;  
precedence left TIMES, DIVIDE;  
precedence left UMINUS;
```

The grammar of Kitten for JavaCup: commands

```
command ::=  
    lvalue ASSIGN exp // assignment  
| type ID ASSIGN exp // variable declaration  
| RETURN              // return from void method  
| RETURN exp          // return from non-void method  
| IF LPAREN exp RPAREN THEN command // conditionals...  
| IF LPAREN exp RPAREN THEN command ELSE command  
| WHILE LPAREN exp RPAREN command // while loop  
| FOR LPAREN command SEMICOLON exp  
    SEMICOLON command RPAREN command // for loop  
| LBRACE statements RBRACE // local scope  
| exp DOT ID LPAREN expseq RPAREN; // method call, again!
```

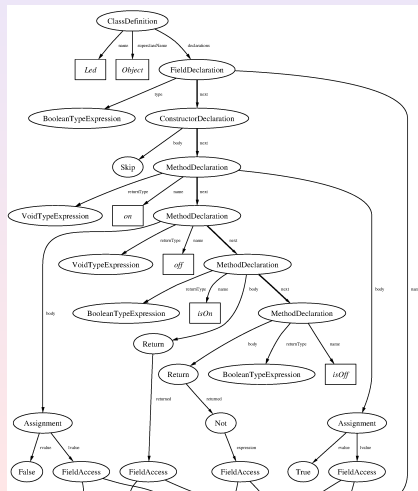

The grammar of Kitten for JavaCup: classes

```
class ::=
    CLASS ID LBRACE class_members RBRACE
  | CLASS ID EXTENDS ID LBRACE class_members RBRACE ;

class_members ::=
  | FIELD type ID class_members
  | CONSTRUCTOR LPAREN formals RPAREN command
    class_members
  | METHOD typeplus ID LPAREN formals RPAREN command
    class_members ;
```

Building the abstract syntax

Given a grammar and a program, JavaCup checks that the program agrees with the grammar, otherwise it issues a syntax error. However, we want more: we want to build the *abstract syntax tree* of the program:



Semantical actions

JavaCup allows one to specify *semantical actions* that are fired when a derivation is used in a parse tree. We will use this possibility to build the syntax, recursively

Semantical actions for types

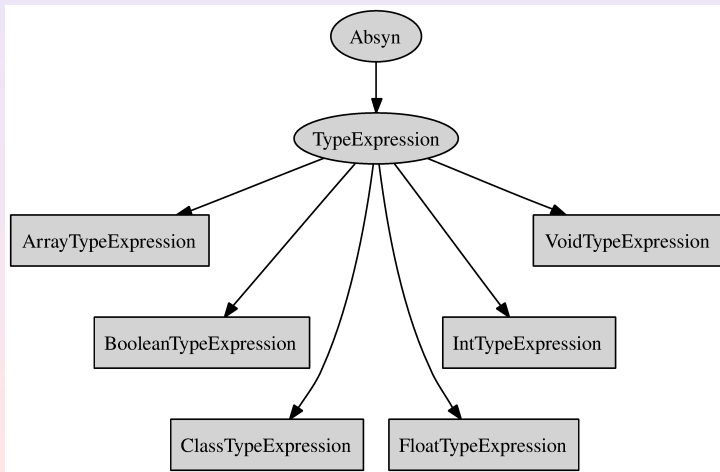
```
type ::=
    ID
  | BOOLEAN
  | INT
  | FLOAT
  | ARRAY OF type ;
```

become

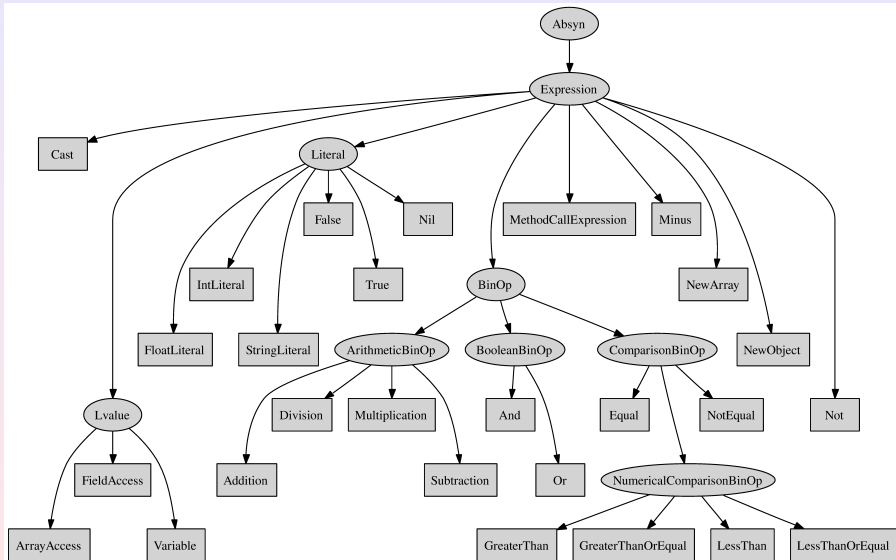
```
type ::=
    ID:id      {: RESULT = new ClassTypeExpression(idleft, id);
  | BOOLEAN:b  {: RESULT = new BooleanTypeExpression(bleft); :}
  | INT:i      {: RESULT = new IntTypeExpression(ileft); :}
  | FLOAT:f    {: RESULT = new FloatTypeExpression(fleft); :}
  | ARRAY:a OF type:t
    {: RESULT = new ArrayTypeExpression(aleft, t); :} ;
```

Classes for the abstract syntax: types

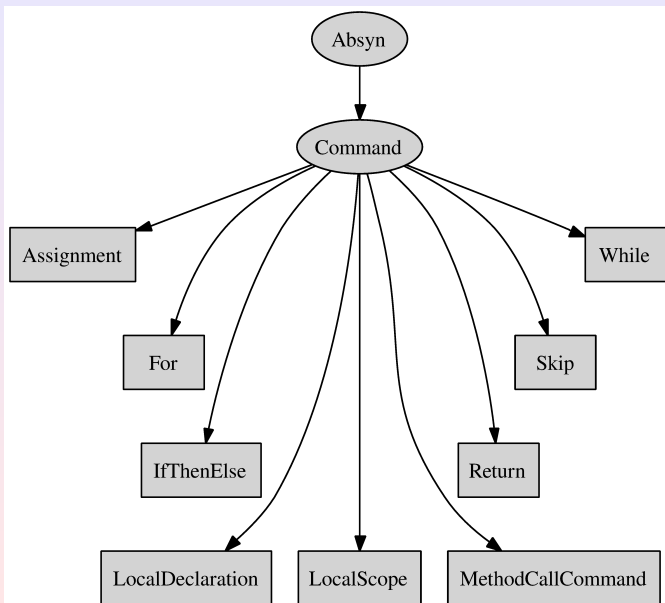
The classes for the abstract syntax are organized into a hierarchy, so that it will be simpler to share code among them in the future



Classes for the abstract syntax: expressions



Classes for the abstract syntax: commands



Semantical analysis

- **input:** abstract syntax tree
- **output:** annotated abstract syntax tree
- **error:** type error, unreachable code. . .

Why do we need a semantical analysis phase?

Grammars do not check for

- type consistency
- undeclared variables
- deadcode
- using the return value of a void method
- ...

The simplest implementation of those checks/inferences is through inductive definitions, that automatically translate into Java code that descends on the abstract syntax tree, recursively

Deadcode identification

```
while (x > 0) {  
  x := x - 1;  
  return;  
  y := y + 1  <- a warning should be issued here  
}
```

$command \vdash^{cdc} Boolean$

$command \vdash^{cdc} true$ means that the execution of *command* definitely ends with the last command(s) of *command*.

Deadcode rules

$$\frac{}{\text{Return}(x) \vdash^{\text{cdc}} \text{true}}$$

$$\frac{}{\text{Assignment}(\text{lvalue}, \text{rvalue}) \vdash^{\text{cdc}} \text{false}}$$

$$\frac{\text{then} \vdash^{\text{cdc}} b_1 \quad \text{else} \vdash^{\text{cdc}} b_2}{\text{IfThenElse}(\text{condition}, \text{then}, \text{else}) \vdash^{\text{cdc}} b_1 \wedge b_2}$$

$$\frac{\text{if } com_1 \vdash^{\text{cdc}} \text{true} \text{ issue a warning} \quad com_2 \vdash^{\text{cdc}} b}{com_1; com_2 \vdash^{\text{cdc}} b}$$

$$\frac{\text{body} \vdash^{\text{cdc}} b}{\text{While}(\text{condition}, \text{body}) \vdash^{\text{cdc}} \text{false}}$$

Type inference and checking 1/4

$$\boxed{\rho \vdash \text{expression} : \text{type}}$$

$$\frac{\rho(\text{name}) \text{ is defined}}{\rho \vdash \text{Variable}(\text{name}) : \rho(\text{name})}$$

$$\frac{\begin{array}{l} \rho \vdash \text{receiver} : \kappa \quad \kappa \in \text{ClassType} \\ \text{field} = \kappa.\text{fieldLookup}(\text{name}) \quad \text{field} \neq \text{null} \end{array}}{\rho \vdash \text{FieldAccess}(\text{receiver}, \text{name}) : \text{field.getType}()}$$

$$\frac{\rho \vdash \text{array} : t \quad t \in \text{ArrayType} \quad \rho \vdash \text{index} : \text{INT}}{\rho \vdash \text{ArrayAccess}(\text{array}, \text{index}) : t.\text{getElementsType}()}$$

Type inference and checking 2/4

$$\overline{\rho \vdash \text{False}() : \text{BOOLEAN}}$$

$$\overline{\rho \vdash \text{True}() : \text{BOOLEAN}}$$

$$\overline{\rho \vdash \text{Nil}() : \text{NIL}}$$

$$\overline{\rho \vdash \text{IntLiteral}() : \text{INT}}$$

$$\overline{\rho \vdash \text{FloatLiteral}() : \text{FLOAT}}$$

$$\overline{\rho \vdash \text{StringLiteral}(\text{value}) : \text{ClassType.mk}(\text{STRING})}$$

$$\frac{\rho \vdash \text{expression} : \text{BOOLEAN}}{\rho \vdash \text{Not}(\text{expression}) : \text{BOOLEAN}}$$

$$\frac{\rho \vdash \text{expression} : t \quad t \leq \text{FLOAT}}{\rho \vdash \text{Minus}(\text{expression}) : t}$$

$$\frac{\text{intoType} = \tau[\![\text{type}]\!] \quad \rho \vdash \text{expression} : \text{fromType} \quad \text{intoType} < \text{fromType}}{\rho \vdash \text{Cast}(\text{type}, \text{expression}) : \text{intoType}}$$

Type inference and checking 3/4

$$\frac{\text{ClassType.mk}(\text{className}) = \kappa \quad \rho \vdash \text{actuals} : \vec{\tau} \quad \kappa.\text{constructorsLookup}(\vec{\tau}) = \{\text{constructor}\}}{\rho \vdash \text{NewObject}(\text{className}, \text{actuals}) : \kappa}$$

$$\frac{\rho \vdash \text{elementsType} : t \quad \rho \vdash \text{size} : \text{INT}}{\rho \vdash \text{NewArray}(\text{elementsType}, \text{size}) : \text{ArrayType.mk}(t)}$$

$$\frac{\rho \vdash \text{receiver} : \kappa \quad \kappa \in \text{ClassType} \quad \rho \vdash \text{actuals} : \vec{\tau} \quad \kappa.\text{methodsLookup}(\text{name}, \vec{\tau}) = \{\text{method}\} \quad r = \text{method.getReturnType}()}{\rho \vdash \text{MethodCallExpression}(\text{receiver}, \text{name}, \text{actuals}) : r}$$

$$\frac{\rho \vdash \text{left} : t_l \quad \rho \vdash \text{right} : t_r \quad t_l \leq \text{FLOAT} \quad t_r \leq \text{FLOAT}}{\rho \vdash \text{ArithmeticBinOp}(\text{left}, \text{right}) : t_l.\text{leastCommonSupertype}(t_r)}$$

$$\frac{\rho \vdash \text{left} : t_l \quad \rho \vdash \text{right} : t_r \quad (\text{either } t_l \leq t_r \text{ or } t_r \leq t_l)}{\rho \vdash \text{Equal}(\text{left}, \text{right}) : \text{BOOLEAN}}$$

Type inference and checking 4/4

$$\frac{t = \tau[\![type]\!] \quad \rho \vdash \text{initialiser} : i \quad i \leq t}{\rho \vdash \text{LocalDeclaration}(\text{type}, \text{name}, \text{initialiser}) : \rho[\text{name} \mapsto t]}$$

$$\frac{\rho \vdash \text{lvalue} : t_l \quad \rho \vdash \text{rvalue} : t_r \quad t_r \leq t_l}{\rho \vdash \text{Assignment}(\text{lvalue}, \text{rvalue}) : \rho}$$

$$\frac{\rho \vdash \text{condition} : \text{Type.BOOLEAN} \quad \rho \vdash \text{then} : \rho' \quad \rho \vdash \text{else} : \rho''}{\rho \vdash \text{IfThenElse}(\text{condition}, \text{then}, \text{else}) : \rho}$$

$$\frac{\text{expression} \neq \text{null} \quad \rho \vdash \text{expression} : t \quad \text{it is in a method that returns } r \geq t}{\rho \vdash \text{Return}(\text{expression}) : \rho}$$

$$\frac{\rho \vdash \text{condition} : \text{Type.BOOLEAN} \quad \rho \vdash \text{body} : \rho'}{\rho \vdash \text{While}(\text{condition}, \text{body}) : \rho}$$

Intermediate code generation

Generation of the intermediate Kitten bytecode

- **input:** annotated abstract syntax tree
- **output:** Kitten bytecode

Why not Java bytecode instead?

- JB is too low level
- JB is only implicitly typed
- JB uses integers for Booleans
- JB has many optimized variants of the same instruction
- ...

Intermediate code generation

Intermediate Kitten bytecode

Led():

return void

on():

load 0 of type Led

const true

putfield Led.state

return void

off():

load 0 of type Led

const false

putfield Led.state

return void

isOn():

load 0 of type Led

getfield Led.state

return boolean

isOff():

load 0 of type Led

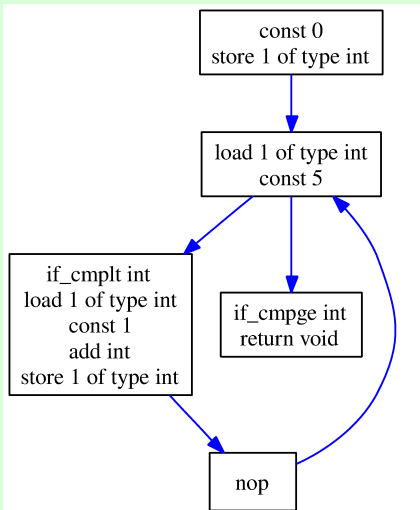
getfield Led.state

neg boolean

return boolean

Intermediate code generation

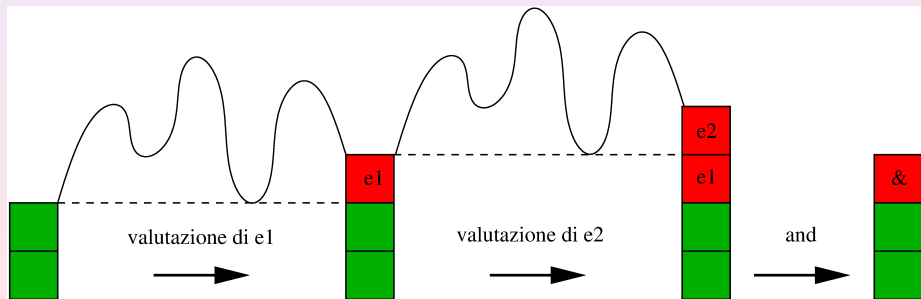
Intermediate Kitten bytecode: loops



Intermediate code generation for expressions 1/3

- 1 leave the original stack untouched
- 2 push on top the value of the expression

For instance, the evaluation of $(e_1 \text{ and } e_2)$ looks like:



Intermediate code generation for expressions 2/3

$$\gamma[\![-]\!] : \text{expression} \mapsto \text{block} \mapsto \text{block}$$

$$\gamma[\![\text{True}()]\!](\beta) = \boxed{\text{const true}} \rightarrow \beta$$

$$\gamma[\![\text{Variable}(name)]\!](\beta) = \boxed{\text{load } num \text{ of type } \tau} \rightarrow \beta$$

$$\gamma[\![\text{And}(left, right)]\!](\beta) = \gamma[\![left]\!] \left(\gamma[\![right]\!] \left(\boxed{\text{and}} \rightarrow \beta \right) \right)$$

$$\gamma[\![\text{FieldAccess}(receiver, name)]\!](\beta) = \gamma[\![receiver]\!] \left(\boxed{\text{getfield } field} \rightarrow \beta \right)$$

$$\begin{aligned} & \gamma[\![\text{ArrayAccess}(array, index)]\!](\beta) \\ &= \gamma[\![array]\!] \left(\gamma[\![index]\!] \left(\boxed{\text{arrayload from array of } \tau} \rightarrow \beta \right) \right) \end{aligned}$$

Intermediate code generation for expressions 3/3

$$\gamma[\![\text{Cast}(\text{type}, \text{expression})]\!](\beta) = \gamma[\![\text{expression}]\!]\left(\boxed{\text{cast from } \tau' \text{ into } \tau} \rightarrow \beta\right)$$

$$\gamma[\![\text{Addition}(\text{left}, \text{right})]\!](\beta) = \gamma^\tau[\![\text{left}]\!]\left(\gamma^\tau[\![\text{right}]\!]\left(\boxed{\text{add } \tau} \rightarrow \beta\right)\right)$$

$$\begin{aligned} \gamma[\![\text{MethodCallExpression}(\text{receiver}, \text{name}, \text{actuals})]\!] \\ = \gamma[\![\text{receiver}]\!]\left(\gamma^{\vec{t}}[\![\text{actuals}]\!]\left(\boxed{\text{virtualcall method}} \rightarrow \beta\right)\right) \end{aligned}$$

$$\begin{aligned} \gamma[\![\text{NewObject}(\text{className}, \text{actuals})]\!](\beta) \\ = \boxed{\begin{array}{l} \text{new } \kappa \\ \text{dup } \kappa \end{array}} \rightarrow \gamma^{\vec{t}}[\![\text{actuals}]\!]\left(\boxed{\text{constructorcall con}} \rightarrow \beta\right) \end{aligned}$$

Intermediate code generation for Boolean expressions

$$\gamma^{test}[\![exp]\!](\beta_{true})(\beta_{false}) = \gamma[\![exp]\!] \left(\boxed{\text{nop}} \langle \begin{array}{l} \boxed{\text{if_true}} \rightarrow \beta_{true} \\ \boxed{\text{if_false}} \rightarrow \beta_{false} \end{array} \right)$$

Passive compilation of leftvalues

$$\boxed{\gamma^{passive} \llbracket lvalue, rvalue \rrbracket (\beta)}$$

Variable(*name*)

$$\gamma^{\tau} \llbracket rvalue \rrbracket \left(\boxed{\text{store num of type } \tau} \rightarrow \beta \right)$$

FieldAccess(*receiver*, *name*)

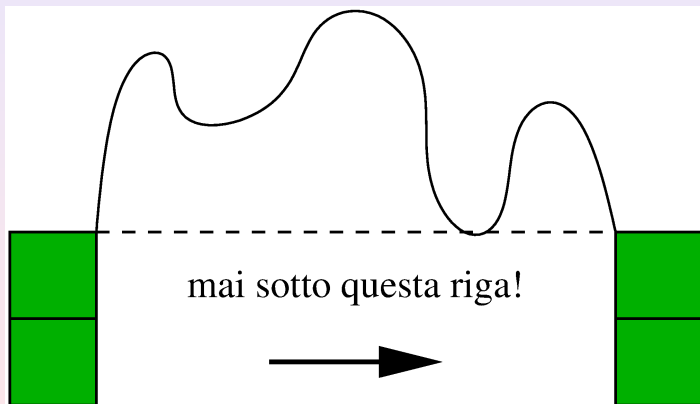
$$\gamma \llbracket receiver \rrbracket \left(\gamma^{\tau} \llbracket rvalue \rrbracket \left(\boxed{\text{putfield field}} \rightarrow \beta \right) \right)$$

ArrayAccess(*array*, *index*)

$$\gamma \llbracket array \rrbracket \left(\gamma \llbracket index \rrbracket \left(\gamma^{\tau} \llbracket rvalue \rrbracket \left(\boxed{\begin{array}{c} \text{arraystore into} \\ \text{array of } \tau \end{array}} \rightarrow \beta \right) \right) \right)$$

Intermediate code generation for the commands 1/2

- 1 leave the original stack untouched



Intermediate code generation for the commands 2/2

$$\gamma[\![\text{Skip}()]\!](\beta) = \beta \quad \gamma[\![\text{LocalScope}(body)]\!](\beta) = \gamma[\![body]\!](\beta)$$

$$\gamma[\![\text{IfThenElse}(cond, then, else)]\!](\beta) = \gamma^{test}[\![cond]\!](\gamma[\![then]\!](\beta))(\gamma[\![else]\!](\beta))$$

$$\gamma[\![\text{Assignment}(lvalue, rvalue)]\!](\beta) = \gamma^{passive}[\![lvalue, rvalue]\!](\beta)$$

$$\gamma[\![\text{While}(cond, body)]\!](\beta) = \underbrace{\boxed{\text{nop}}}_{pivot} \rightarrow \gamma^{test}[\![cond]\!](\gamma[\![body]\!](pivot))(\beta)$$

Java bytecode generation

Generation of the object Java bytecode

- **input:** Kitten bytecode
- **output:** Java bytecode

- ➊ each Kitten bytecode is translated into one or more Java bytecodes, exploiting potential optimized variants
- ➋ each block of Kitten bytecodes becomes a sequential snippet of Java bytecodes
- ➌ such snippets are flattened into a linear sequence of Java bytecodes, by adding goto's
- ➍ everything is finally packaged into class files
- ➎ a Java bytecode manipulation library such as BCEL does most of the work: <http://commons.apache.org/proper/commons-bcel>

Security guarantees of the Java bytecode

A controlled low-level language

- code cannot be modified at runtime \Rightarrow **no metamorphic code**
- at each given program point, number and static types of stack elements and local variables are fixed and statically known, for all possible execution paths \Rightarrow **strong types, you cannot blow up the stack**
- casts are checked \Rightarrow **you cannot pretend that a frog is a prince**
- initialized local variables cannot be used \Rightarrow **you cannot exploit stale values**
- new objects cannot be used before calling one of their constructors \Rightarrow **you cannot forge raw data**
- all jumps go to static targets, inside the same method \Rightarrow **you cannot dynamically build your jump targets**

Identical constraints hold for the Dalvik Android bytecode