

Textual modeling languages

Oszkár Semeráth, Kristóf Marussy

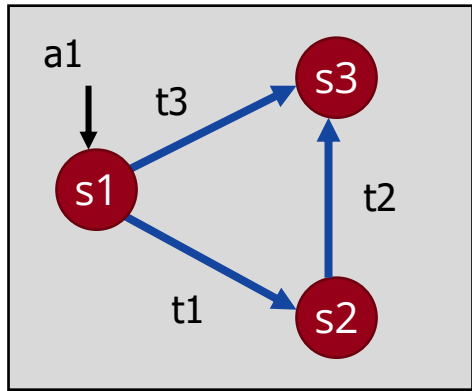


**Critical Systems
Research Group**

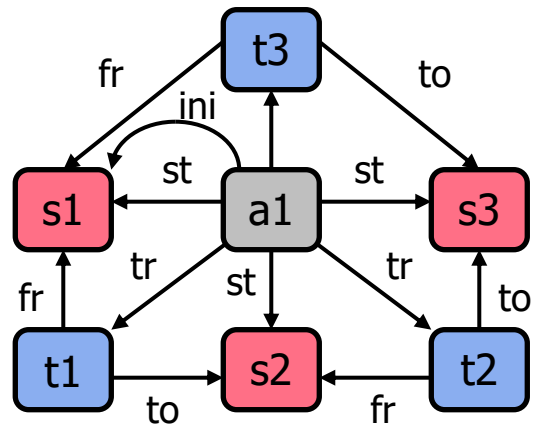
Last lecture: Concrete and Abstract Syntax

- **Question:** How to capture models?
- **Answer:** graph-based structures!

Concrete syntax



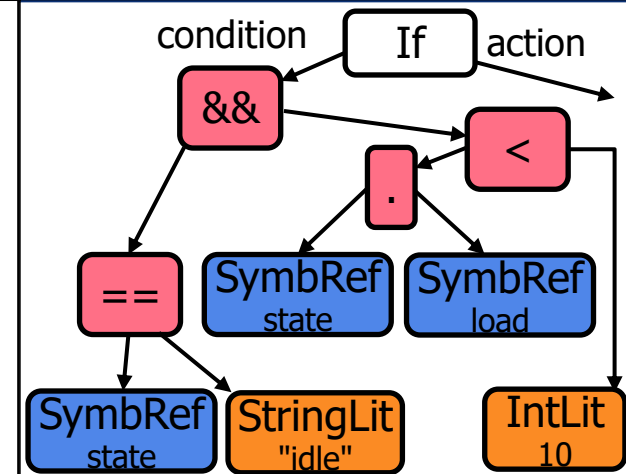
Abstract syntax



Concrete syntax

```
if (  
  state ==  
  "idle" &&  
  this.load < 10)  
  ...
```

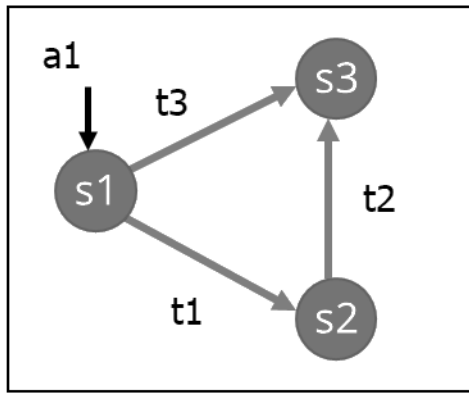
Abstract syntax



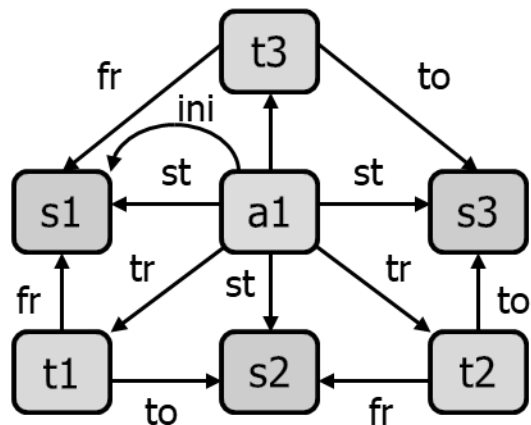
Last lecture: Concrete and Abstract Syntax

- **Question:** How to capture models?
- **Answer:** graph-based structures!

Concrete syntax



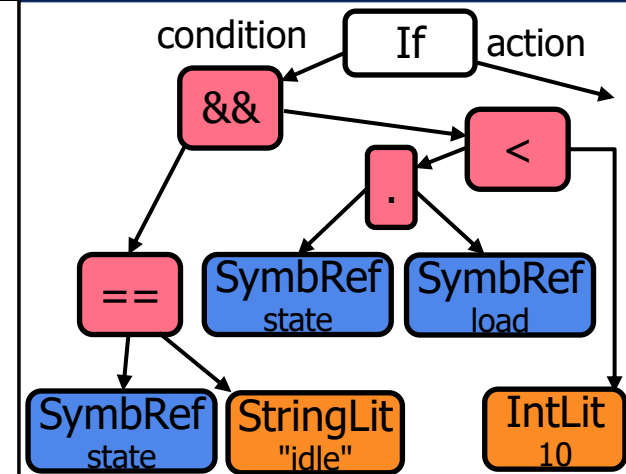
Abstract syntax



Concrete syntax

```
if (  
  state ==  
  "idle" &&  
  this.load < 10)  
  ...
```

Abstract syntax



- **Today's Topic:** Textual syntax and textual editors

Grammars, Context-Free Parsing

Theoretical Background

Conceptual overview

1. There are textual documents that needs to be parsed.
2. Some documents are valid, some are invalid. We want to efficiently capture which textual documents are valid.
→ **Grammars**
3. We want to create editors for valid document.
→ **Parsers**
4. For each document, we want to construct a model.
→ **Abstract syntax graphs**

Terminology

- Alphabet (Σ): Set of possible symbols in the document, e.g.,
 - $\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$ a numerical field
 - ASCII characters textual documents
 - $\{a, b\}$, then we can use only those two letters
- Words (Σ^*): Sequence of symbols, e.g.,
 - “0012”, “222”, “1”, or “ ” also called ε
 - “*class person { }*”
 - worlds like “*babababbbba*”
- Language (L): a subset of all words ($L \subseteq \Sigma^*$)
 - We do not like “0012”, but we like “222”
 - We do like “*class person { }*” but not “*class person { { }*”
- **How to capture a language?**

Formal Grammar

- Formal grammar

$$G = (N, T, P, S)$$

- N : nonterminal symbols
- T : terminal symbols (alphabet)
- P : production rules
- S : start symbol ($S \in N$)

Example: $G = (N, T, P, S)$

- $\text{Num} \rightarrow \text{Digit Num}$

← Traditionally first rule is Start symbol

- $\text{Num} \rightarrow \text{Digit}$

- $\text{Digit} \rightarrow 0 \mid 1 \mid 2 \dots \mid 9$

Formal Grammar

- Notation:

- A, B, C : nonterminals in N ,
- a, b, c : terminals in T
- $\alpha, \beta, \gamma \in (T \cup N)^*$

- Regular rules:

- $B \rightarrow a$
- $B \rightarrow aC$

EASY

- Context-free (CF) rules:

- $B \rightarrow \alpha$
- $B \rightarrow \varepsilon$ (empty)

Theoretical background
+ Good technologies

- Context-dependent rules:

- $\alpha \rightarrow \beta$

Same power as a computer

- There are languages that cannot be captured by rules



Derivation and Language

- **Derivation step** using grammar $G = (T, N, P, S)$

- $\alpha A \gamma \rightarrow \alpha \beta \gamma$
- applying production rule: $A \rightarrow \beta$
- $\alpha A \gamma, \alpha \beta \gamma$: sentential forms

- **Derivation over G:** $S \rightarrow^* w$ where

- S : start symbol
- \rightarrow^* transitive closure (apply as long as possible)
- $w \in T^*$: sentence, i.e. string of terminals only

- **Language generated by G**

- $L(G) = \{w \in T^* \mid \text{there exists a derivation } S \rightarrow^* w \text{ of } G\}$
- Set of sentences derivable from S

- Parsing is polynomial algorithm for regular and context-free grammars
- In general, nonterminals can be resolved in arbitrary order (non-deterministic)
- Leftmost vs. Rightmost derivation: always resolve the left/rightmost nonterminal as next step

Num

Digit Num

1 Num

1 Digit Num

1 9 Num

1 9 Digit Num

1 9 Digit Digit

1 9 Digit 6

1 9 7 6

Num \rightarrow Digit Num

Digit \rightarrow 1

Num \rightarrow Digit Num

Digit \rightarrow 9

Num \rightarrow Digit Num

Num \rightarrow Digit

Digit \rightarrow 6

Digit \rightarrow 7

Binary Operations

Example: $G = (N, T, P, S)$

$\text{Exp} \rightarrow \text{Num}$

$\text{Exp} \rightarrow \text{Exp} \text{ "+" } \text{Exp}$

$\text{Exp} \rightarrow \text{Exp} \text{ "-" } \text{Exp}$

$\text{Exp} \rightarrow \text{Exp} \text{ "*" } \text{Exp}$

$\text{Exp} \rightarrow \text{Exp} \text{ "/" } \text{Exp}$

$\text{Exp} \rightarrow \text{"(" Exp ")"}$

Two Derivations of $1+2*3$:

Exp	
Exp + Exp	Exp \rightarrow Exp "+" Exp
1 + Exp	Exp \rightarrow Num
1 + Exp * Exp	Exp \rightarrow Exp "*" Exp
1 + 2 * Exp	Exp \rightarrow Num
1 + 2 * 3	Exp \rightarrow Num

Exp	
Exp * Exp	Exp \rightarrow Exp "*" Exp
Exp * 3	Exp \rightarrow Num
Exp + Exp * 3	Exp \rightarrow Exp "+" Exp
Exp + 2 * 3	Exp \rightarrow Num
1 + 2 * 3	Exp \rightarrow Num

Lexing and Parsing

Lexer

- **Input:**
 - Regular grammar / RegExp
 - Character sequence:
l,e,t, ,x,=,1,3,4,
- **Output:**
 - Token sequence: let,x,=,134,
 - Identify keywords, numbers, variables, comments

Grammar:

Num \rightarrow 0 Num | ... | 9 Num Num \rightarrow 0 | ... | 9

RegExp:

Num = Digit Digit*

Digit = (0 | 1 | ... | 9)

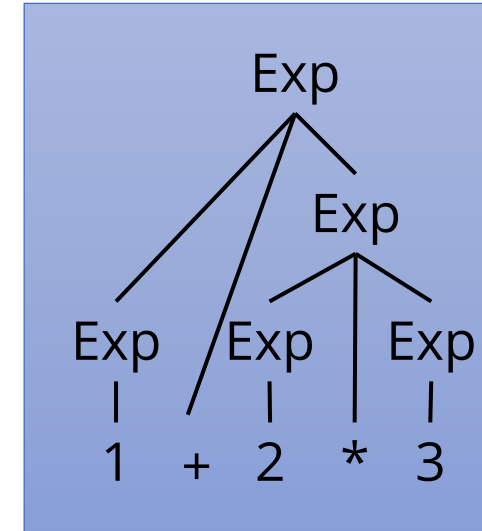
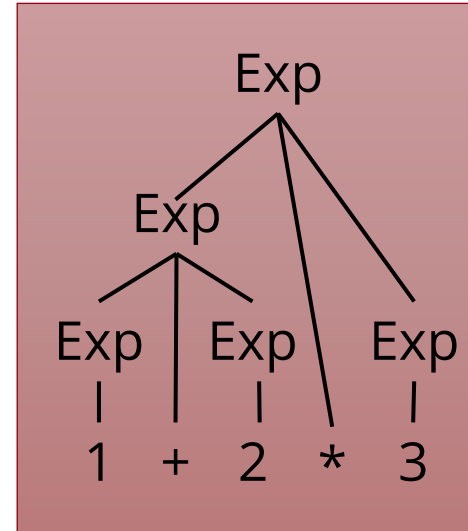
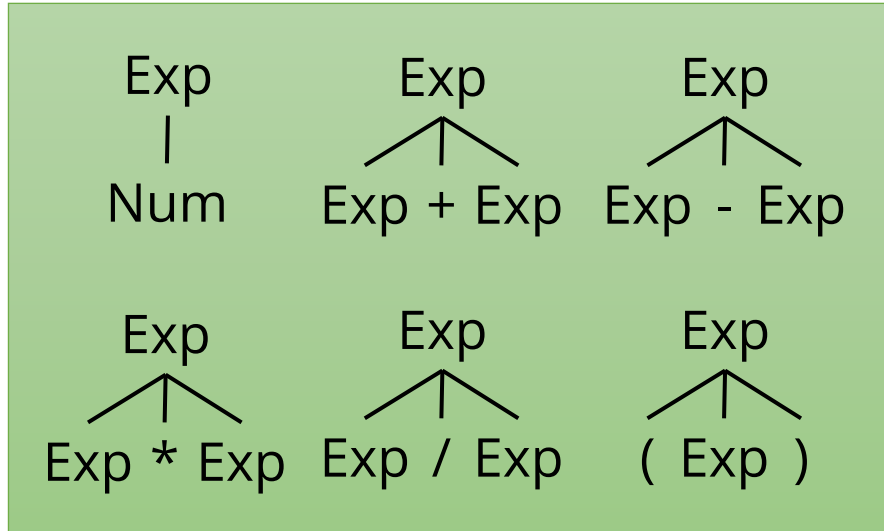
Parser

- **Input:**
 - CF grammar
 - Token sequence: let,x,=,134,
- **Output**
 - Abstract syntax tree (AST)
 - How to derive the token sequence according to grammar?

Grammar:

Exp \rightarrow Num | Exp "+" Exp
| Exp "-" Exp | Exp "*" Exp
| Exp "/" Exp | "(" Exp ")"

Parse Tree Construction



Ambiguous derivation!
How to disambiguate?

- **Parse Tree:**

- Parent node: nonterminals
- Child node: nonterminals/terminals
- Built up according to productions of the grammar

- **Ambiguous grammar:**

- Generates two distinct parse trees
- Serious problem for a parser!

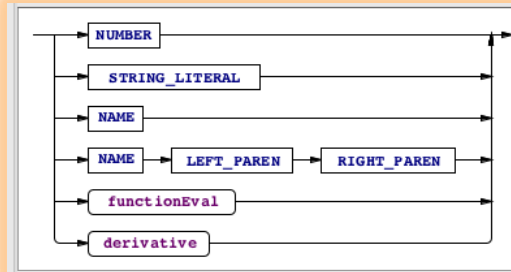
Idea: Change the grammar

Looking inside advanced IDEs

From parsers to development tools
Roundtrip property
Modern services

Parsers: The Traditional Setup

Grammar

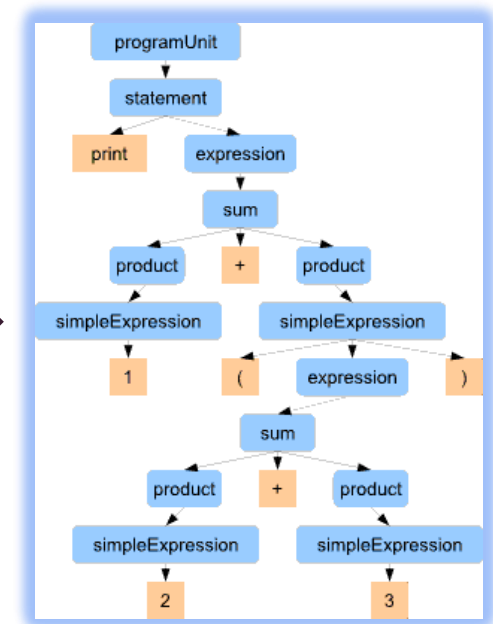


Parsing
(lexer + parser)



```
import com.lauchenauer.istockhelper.  
import com.lauchenauer.lib.ui.Vertic  
public class AboutDialog extends JDia  
protected CardLayout mLayout;  
protected JButton mCredits;  
protected JPanel mMainPanel;  
public AboutDialog(JFrame owner) {  
    super(owner);  
    setModal(true);  
    setUndecorated(true);  
    initUI();  
}  
protected void initUI() {  
    setSize(440, 600);  
    Container cont = getContentPane  
    JPanel p =
```

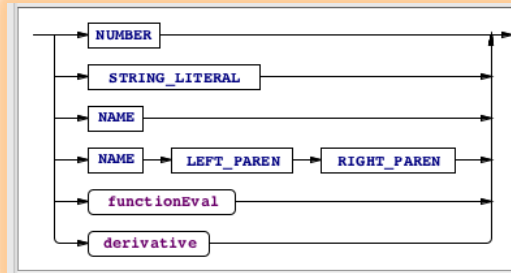
Source code of program



Abstract
syntax tree (AST)

Parsers: Setup with errors

Grammar

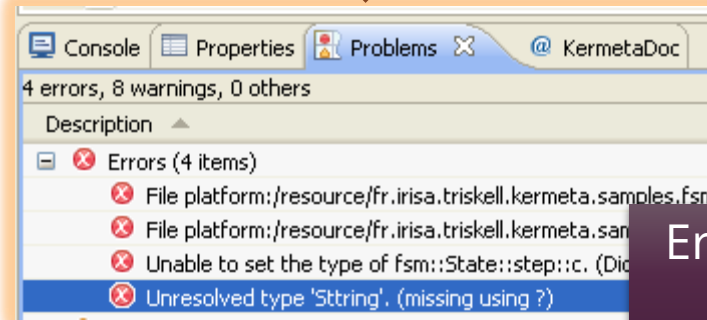


Parsing
(lexer + parser)

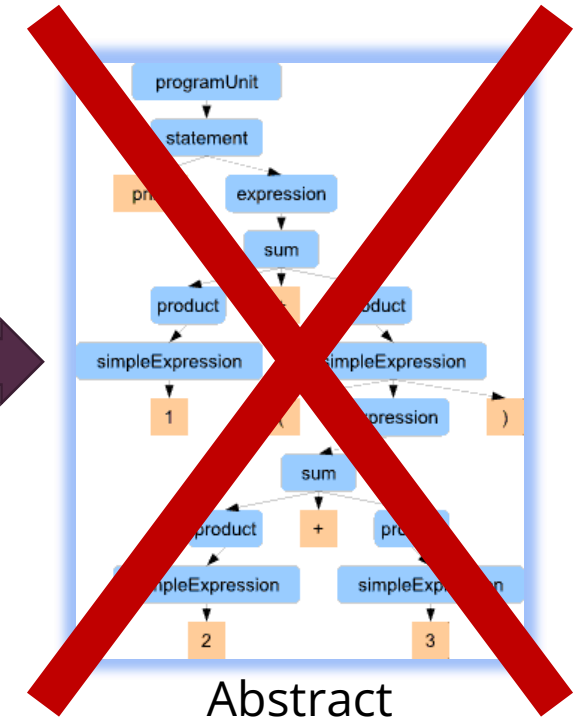
Error
report

```
import com.lauchenauer.istockhelper.  
import com.lauchenauer.lib.ui.Vertic  
public class AboutDialog extends JDia  
protected CardLayout mLayout;  
protected JButton mCredits;  
protected JPanel mMainPanel;  
public AboutDialog(JFrame owner) {  
    super(owner);  
    setModal(true);  
    setUndecorated(true);  
    initUI();  
}  
protected void initUI() {  
    setSize(440, 600);  
    Container cont = getContentPane  
    JPanel p =
```

Source code of program

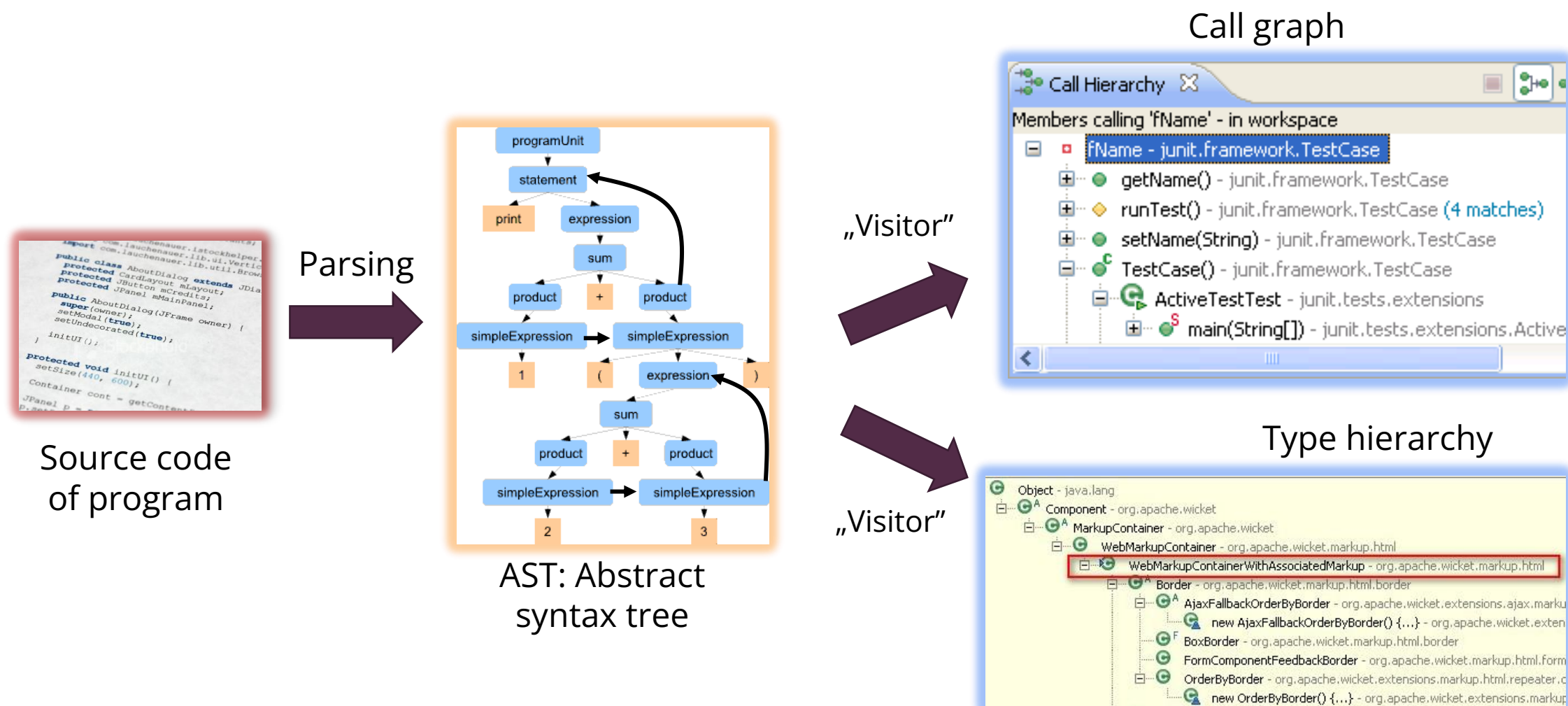


Error recovery
parsing

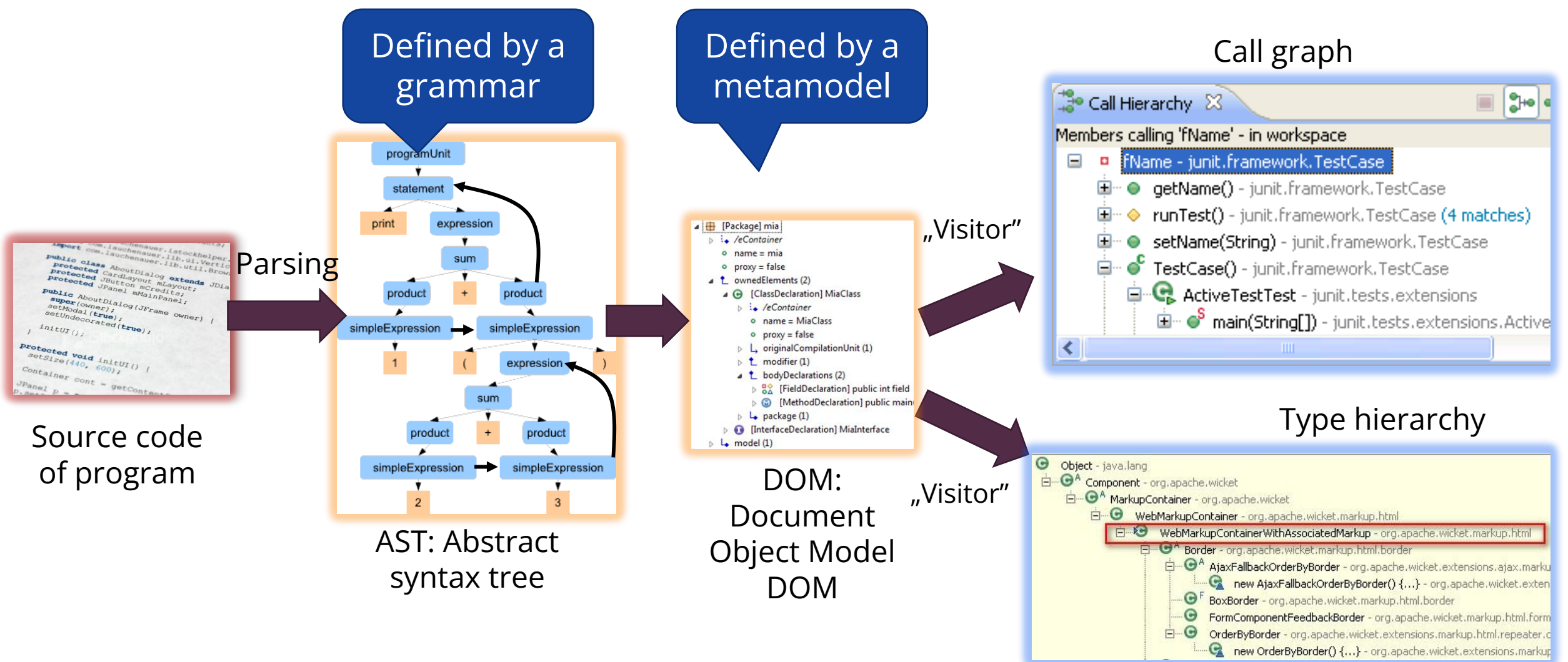


Abstract
syntax tree (AST)

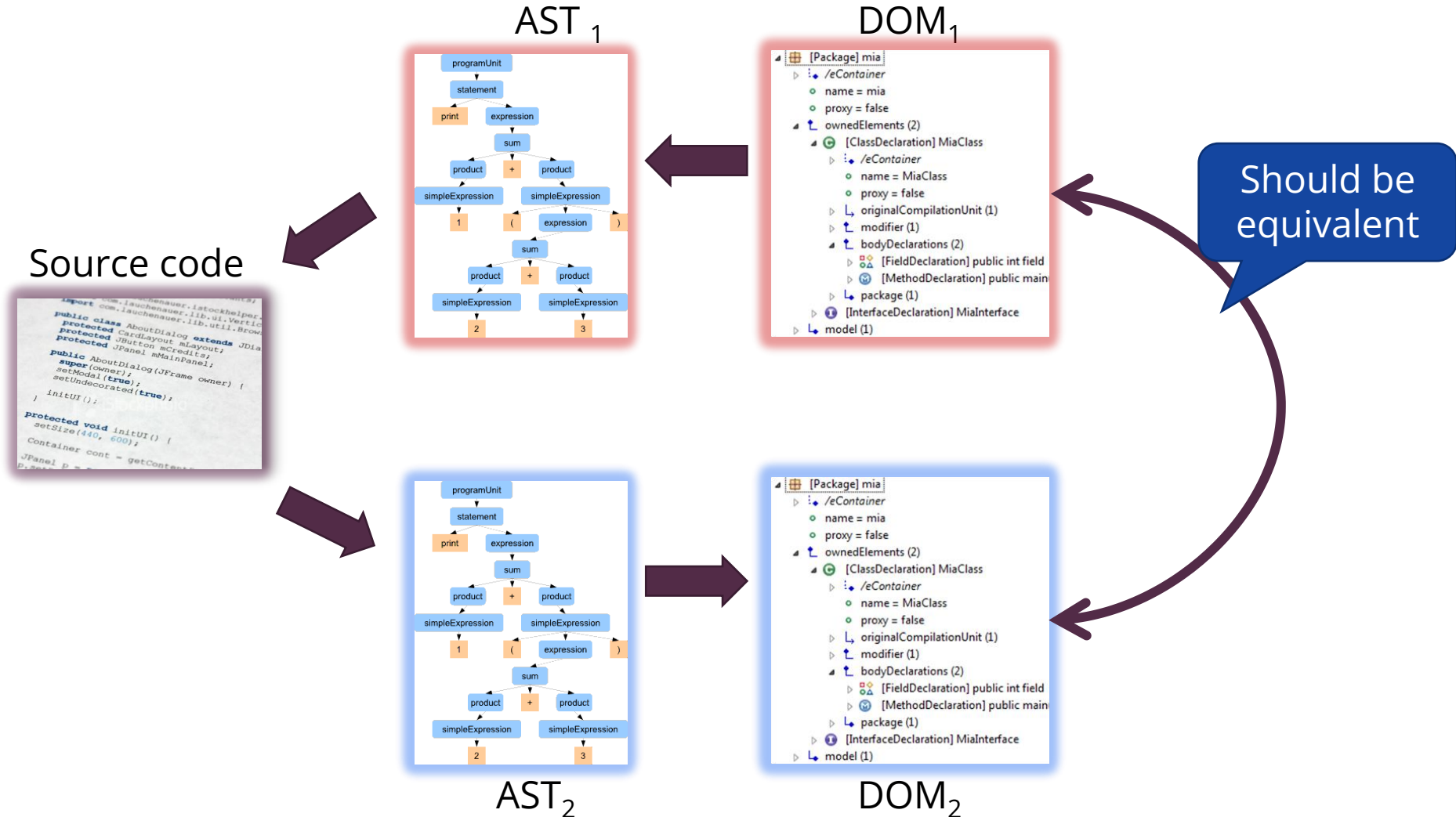
View Generation + Program Analysis



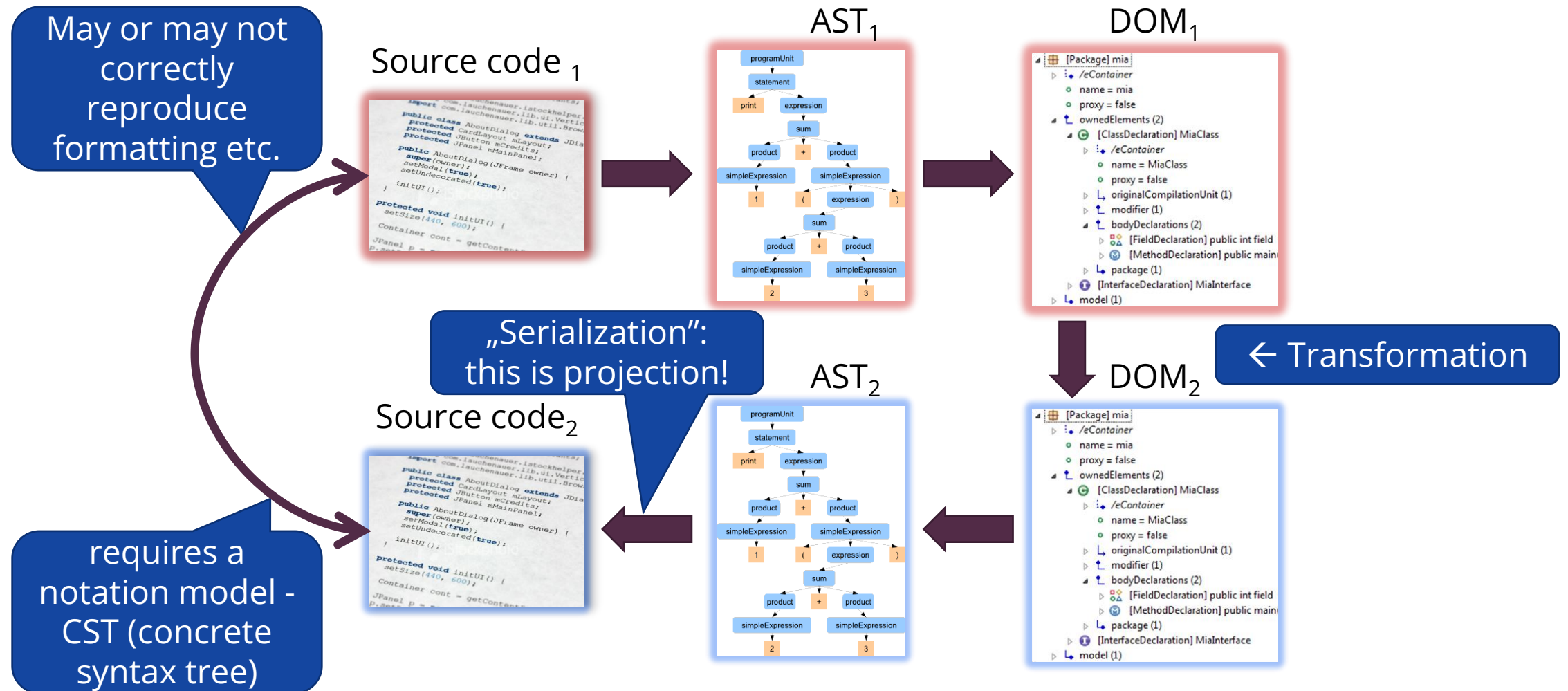
AST vs DOMs



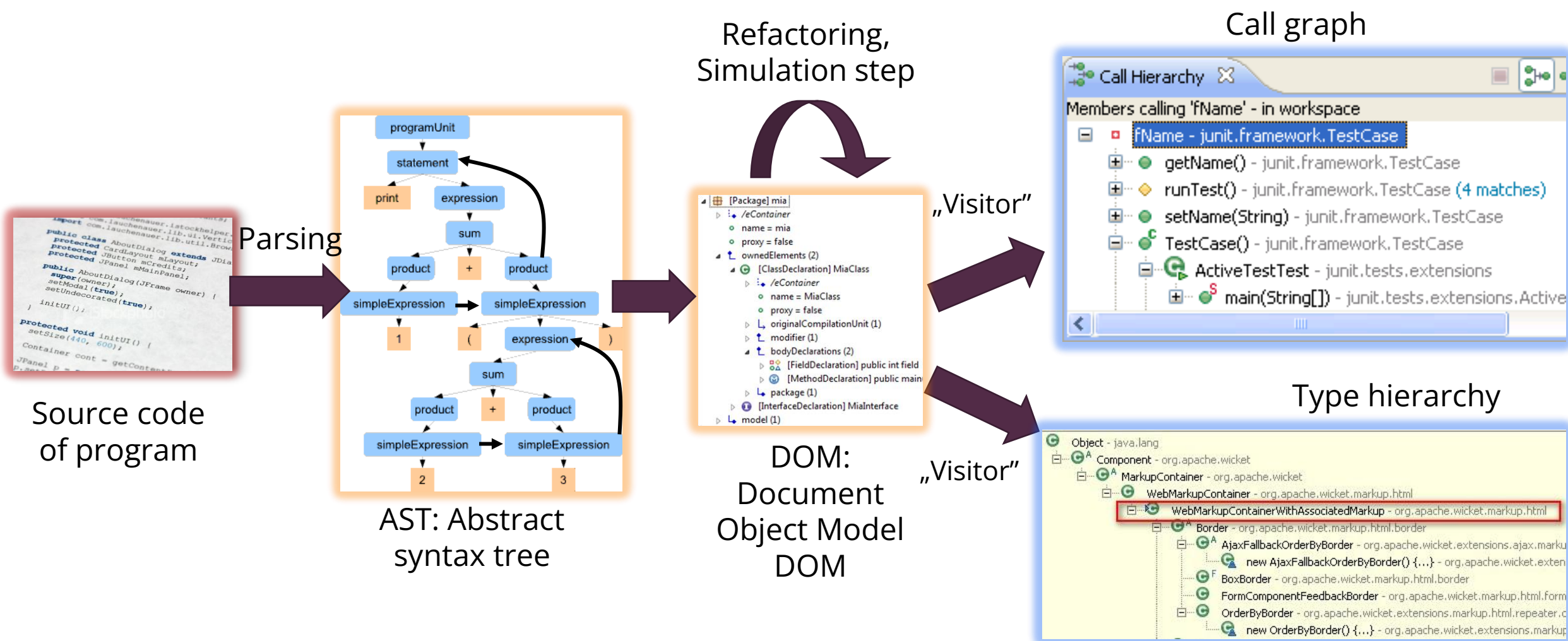
AST + DOMs: persistence roundtrip



Model editing roundtrip

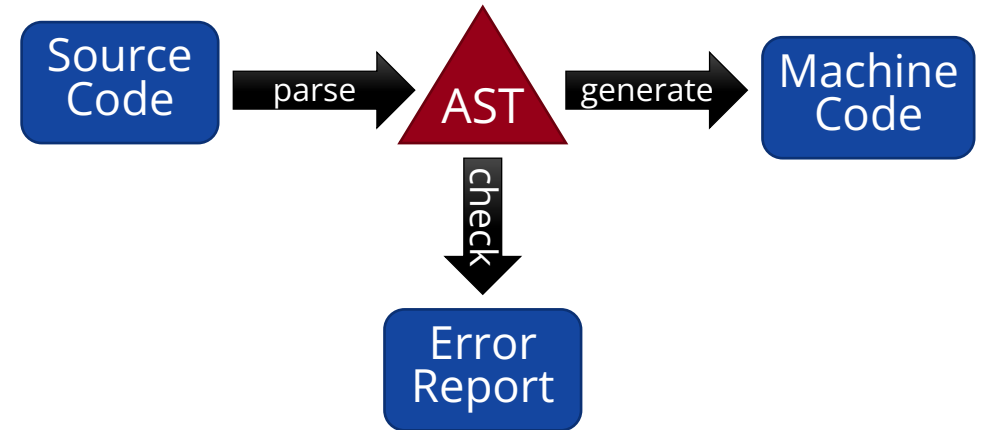


Textual DSM Languages: An Overview



Traditional Architecture of Compilers

- On demand parsing
 - Explicit user's request
 - E.g. `javac myClass.java`
- Parsing:
 - Successful: AST generated
 - Failed: Errors reported (no AST)
- Semantic checks
 - Successful: Machine code generated
 - Failed: Errors reported



Modern IDEs

- Auto-parse-and-check
 - During typing: Parse
 - Upon save: Analyze
- Parsing:
 - AST always generated
 - Error markers on failure
- Semantic analysis
 - Successful:
Machine code generated
 - Failed: Errors reported



Syntactic editor services:

- syntax checking
- syntax highlighting
- outline view
- code folding
- bracket matching...

Semantic editor services:

- error checking
- reference resolving
- hover help
- „code mining“
- code completion
- refactoring...

Source: Guido Wachsmuth (Compiler Construction at TU Delft)

Textual syntax specification

Regular grammars

Context-free grammars

Textual Domain-specific Languages

- Idea
 - Describing models as text files
- Textual development
 - Long history (30+ years)
 - Well-researched theory
 - Mature tools

(Bad) idea: Regular expressions

- Pattern matching for strings
 - Good support
 - Most programming languages
 - More or less the same syntax
 - Calculates and returns matches
- Usable as DSL parser?
 - Weak expressive power
 - Example: balanced parentheses (see pumping lemma)
 - Output is a single boolean variable (does it match?)
 - Missing: interpretability of contents, error localization
 - In some cases, not very concise...

This regular expression will only validate addresses that have had any comments stripped and replaced with whitespace



Regex: validation

- Output is a single boolean variable
 - Decides whether the string matches the language
- What is missing?

Error localization!
Debugging!

(it has a bug with "00")

Context-Free (CF) grammars

Wide-spread adoption

- Sufficient expressive power for describing models
- Intuitive semantics and usage
- Advanced parsing techniques
 - Performance
 - Error localization
 - ...
- Automated tooling for parser generation

Example: name lists CF grammar

- Terminal Symbols
 - “*Dániel*”, “*István*”, “*Zoltán*”, “*and*”, “,”
- Non-terminal Symbols
 - «*Name*», «*Sentence*», «*List*»
- Metalanguage
 - $::=$, $|$, « \rangle »

«*Sentence*» ::= «*Name*» | «*List*» *and* «*Name*»

«*List*» ::= «*List*», «*Name*» | «*Name*»

«*Name*» ::= *Zoltán* | *István* | *Dániel*

Examples: *Dániel* *Dániel, Zoltán* *Dániel, Zoltán and Dániel*

Example from: Dick Grune, Criel J.H. Jacobs: *Parsing Techniques*

The Parsing Process

Derivation trees
Lexer vs parser
Linking & scoping

Example: derivation / parse tree

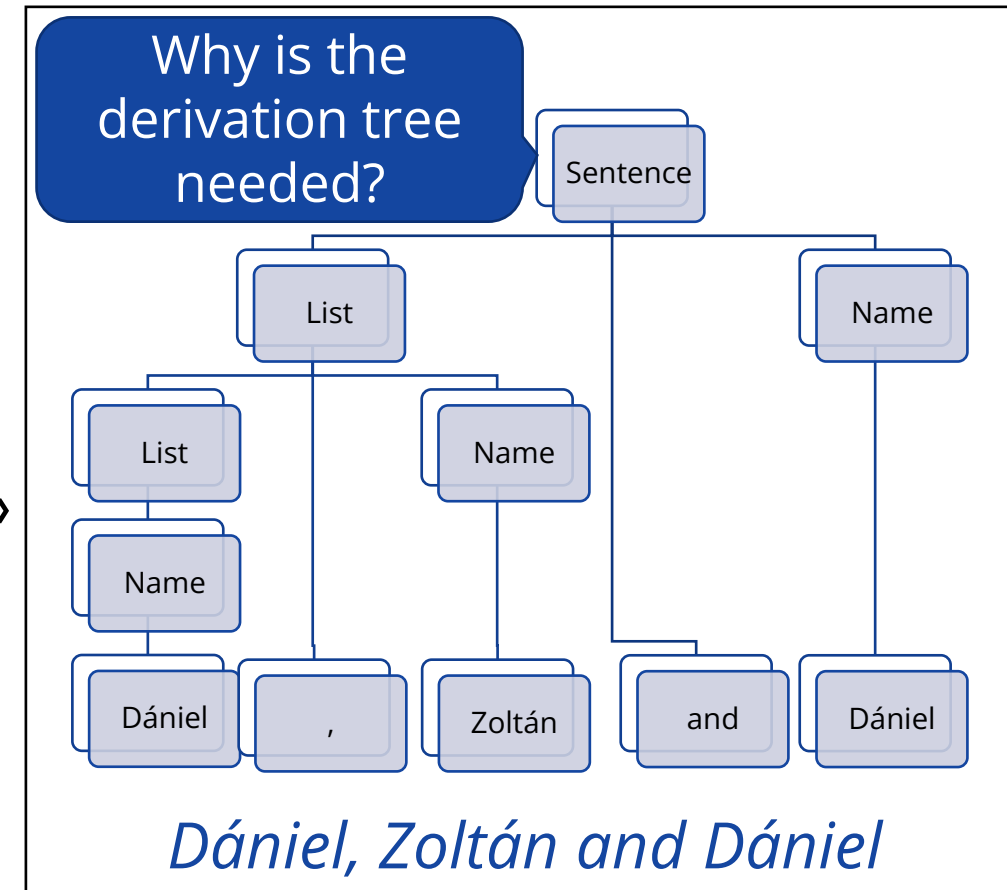
- Terminal Symbols
 - “*Dániel*”, “*István*”, “*Zoltán*”, “*and*”, “,”
- Non-terminal Symbols
 - «*Name*», «*Sentence*», «*List*»
- Metalanguage
 - $::=$, $|$, « \rangle »

«*Sentence*» $::=$ «*Name*» $|$ «*List*» *and* «*Name*»

«*List*» $::=$ «*List*», «*Name*» $|$ «*Name*»

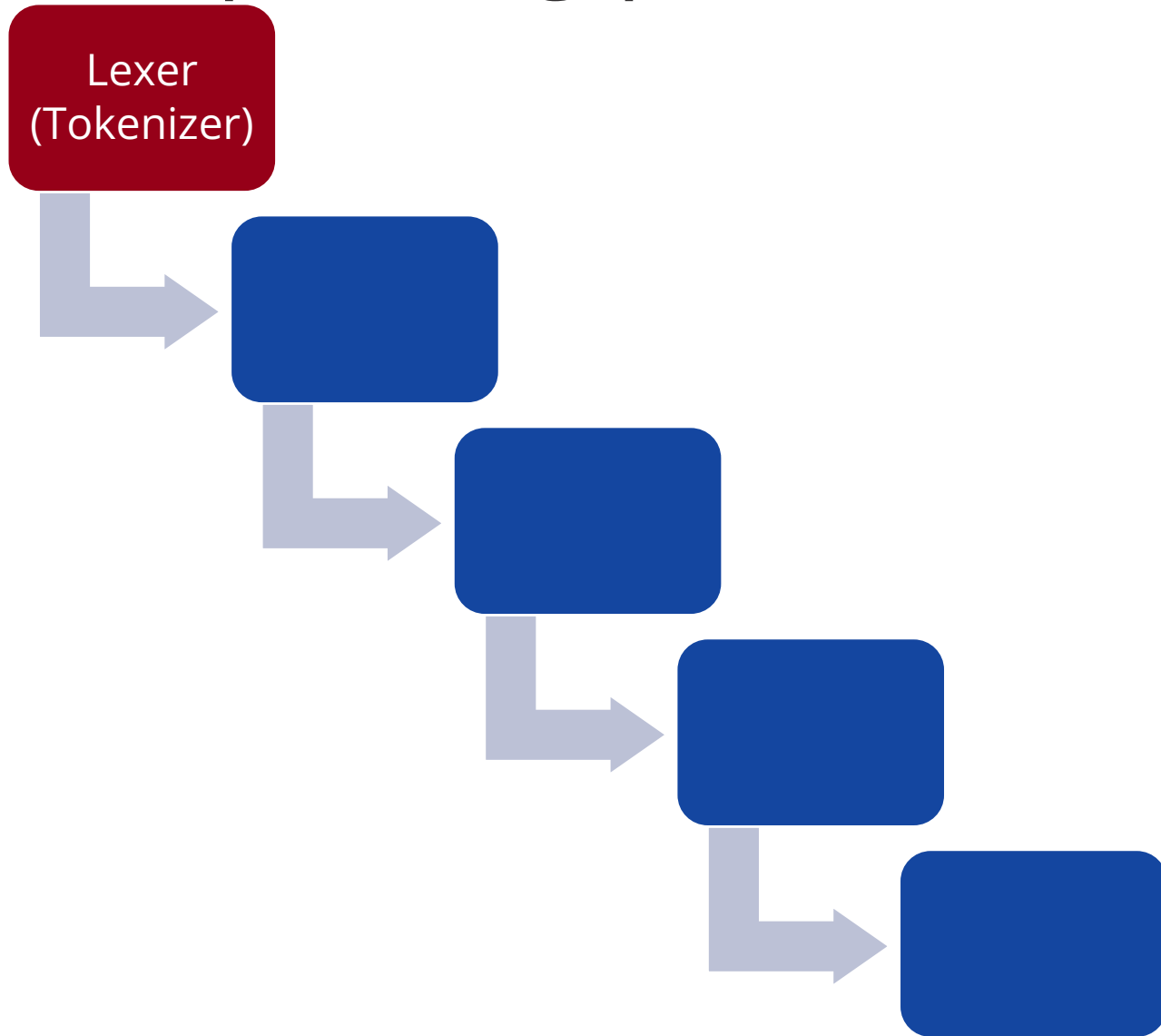
«*Name*» $::=$ *Zoltán* $|$ *István* $|$ *Dániel*

Examples: *Dániel* *Dániel, Zoltán*



Example from: Dick Grune, Ceriel J.H. Jacobs: *Parsing Techniques*

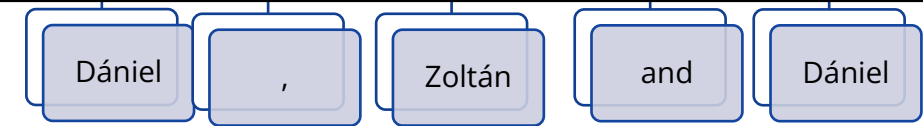
The parsing process



Lexer (Tokenizer)

- **Input:** Character stream
- **Output:** Token stream
- Tokenizing the input character stream
- Similar to the parsing problem
 - But usually simpler – Typically regular expressions
 - Only word/token identification
 - Optional task: leaving out comments

Dániel, Zoltán and Dániel

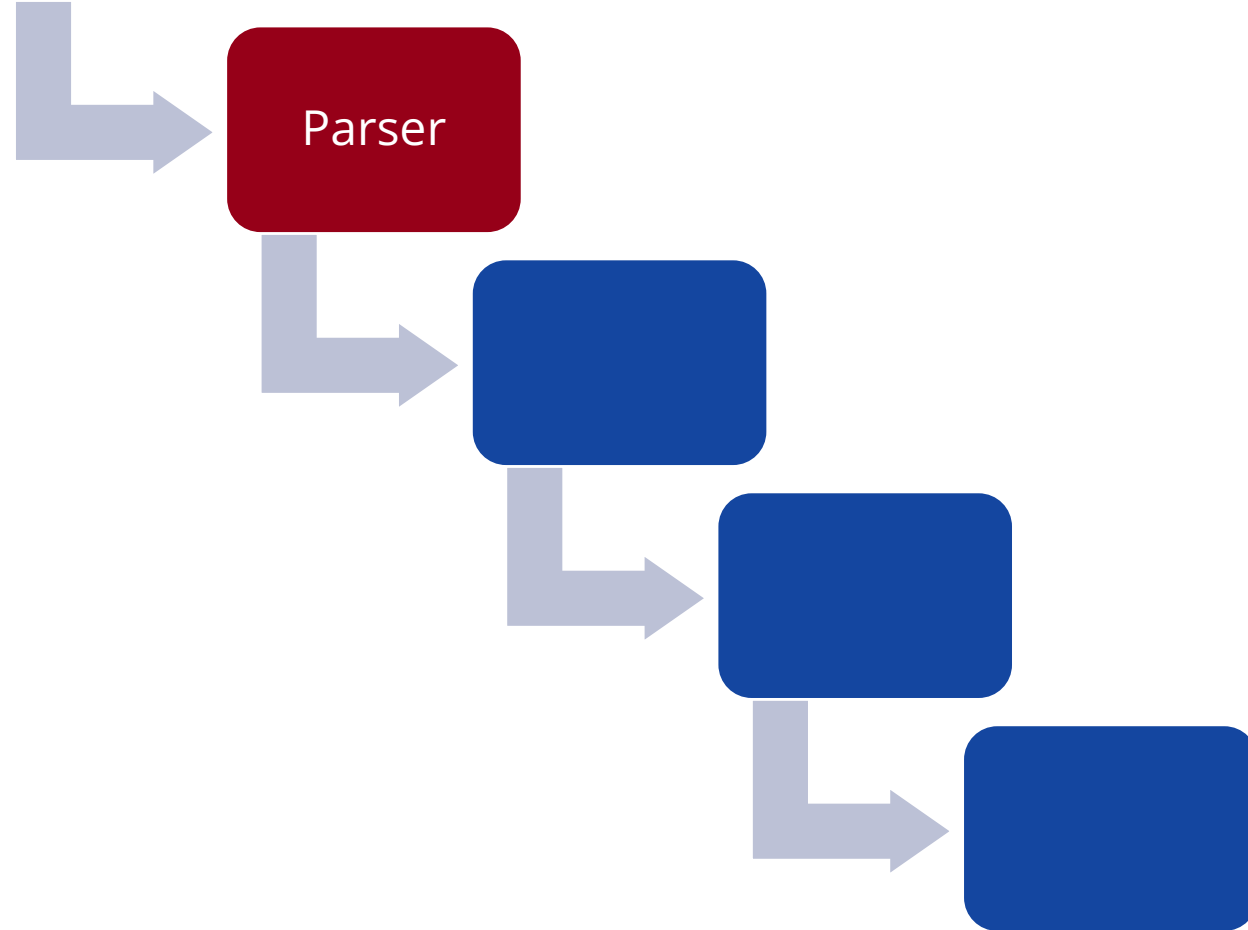


Error handling, Performance, Problem decomposition

The parsing process

Lexer
(Tokenizer)

Input: Character stream
Output: Token stream
Services: Syntax highlight



Parser

Grammar:

«Sentence» ::= «Name» | «List» and «Name»

«List» ::= «List», «Name» | «Name»

«Name» ::= Zoltán | István | Dániel

Token stream:

Dániel

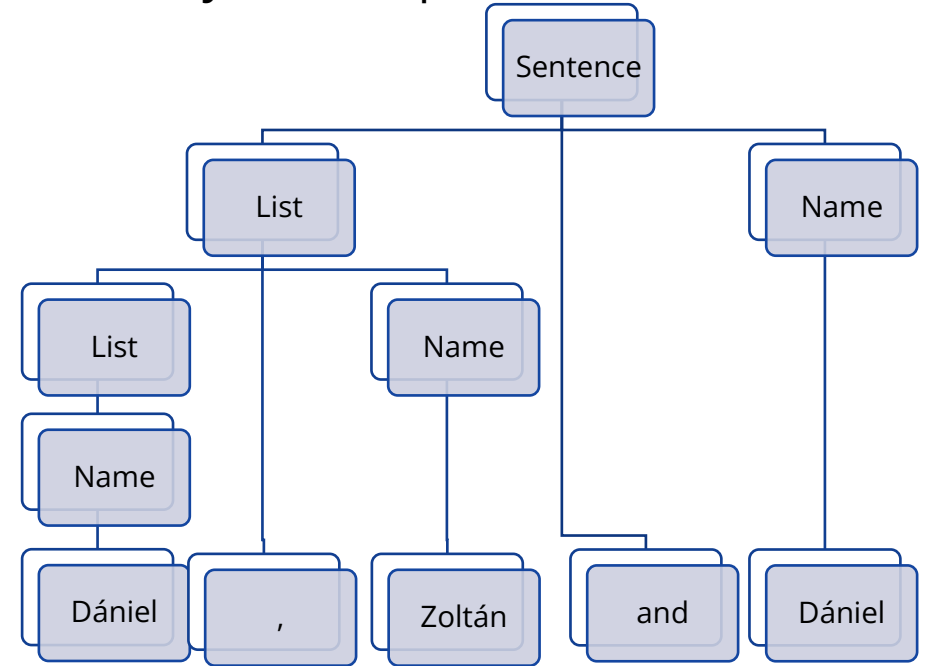
,

Zoltán

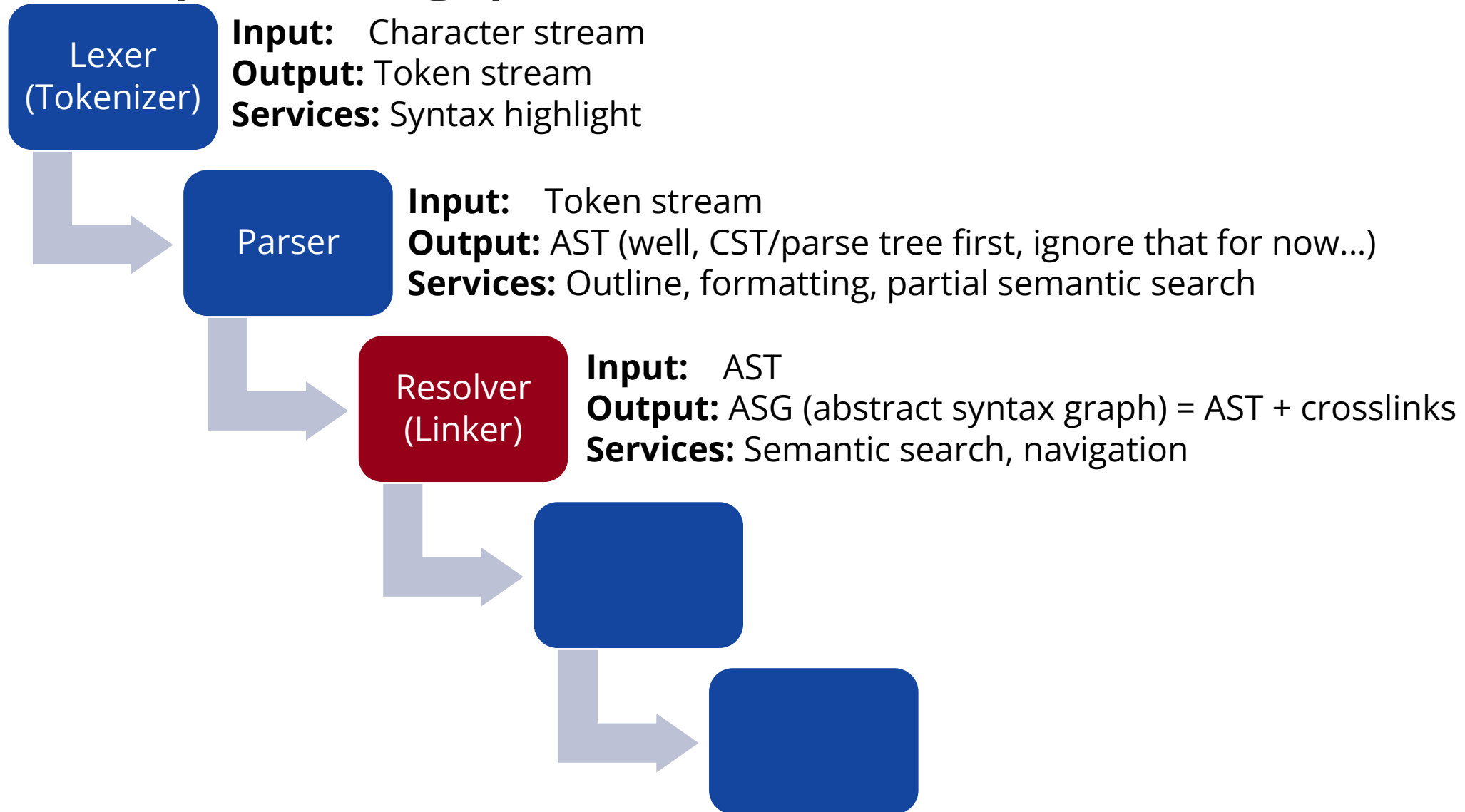
and

Dániel

Abstract Syntax Graph (ASG)



The parsing process



Variable / Identifier Handling

- Variables
 - At runtime: Value calculation/substitution
 - Editing/analysis time: Refers to other parts of the AST
- **Variable declaration:**
 - A declaration of a variable
 - Unique naming: Variable definitions must be resolvable → Extra phase
- **Variable reference:** Use of an already defined variable

`int a=3;` Declaration

`System.out.println(a);` Reference

- Parser checks: “Can a variable be named ‘a’?”
- Reference resolution: “Is a variable defined?” → Scoping problem

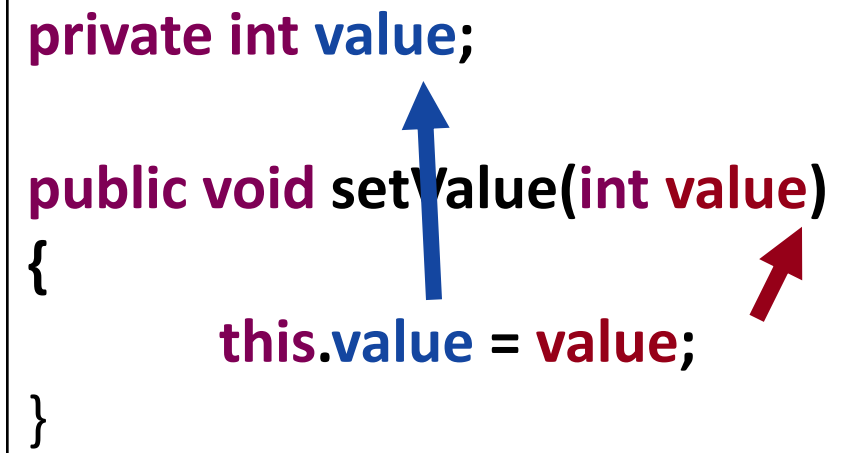
Scoping/Linking problem

- Possible approaches
 - Most specific declaration
 - **Hierarchical scopes**
 - Conflict is error
 - Qualified references
 - ...
- Linking may be *lazy*

Scope: the set of elements (variable names, identifiers) that can be referenced at a given point. Used:

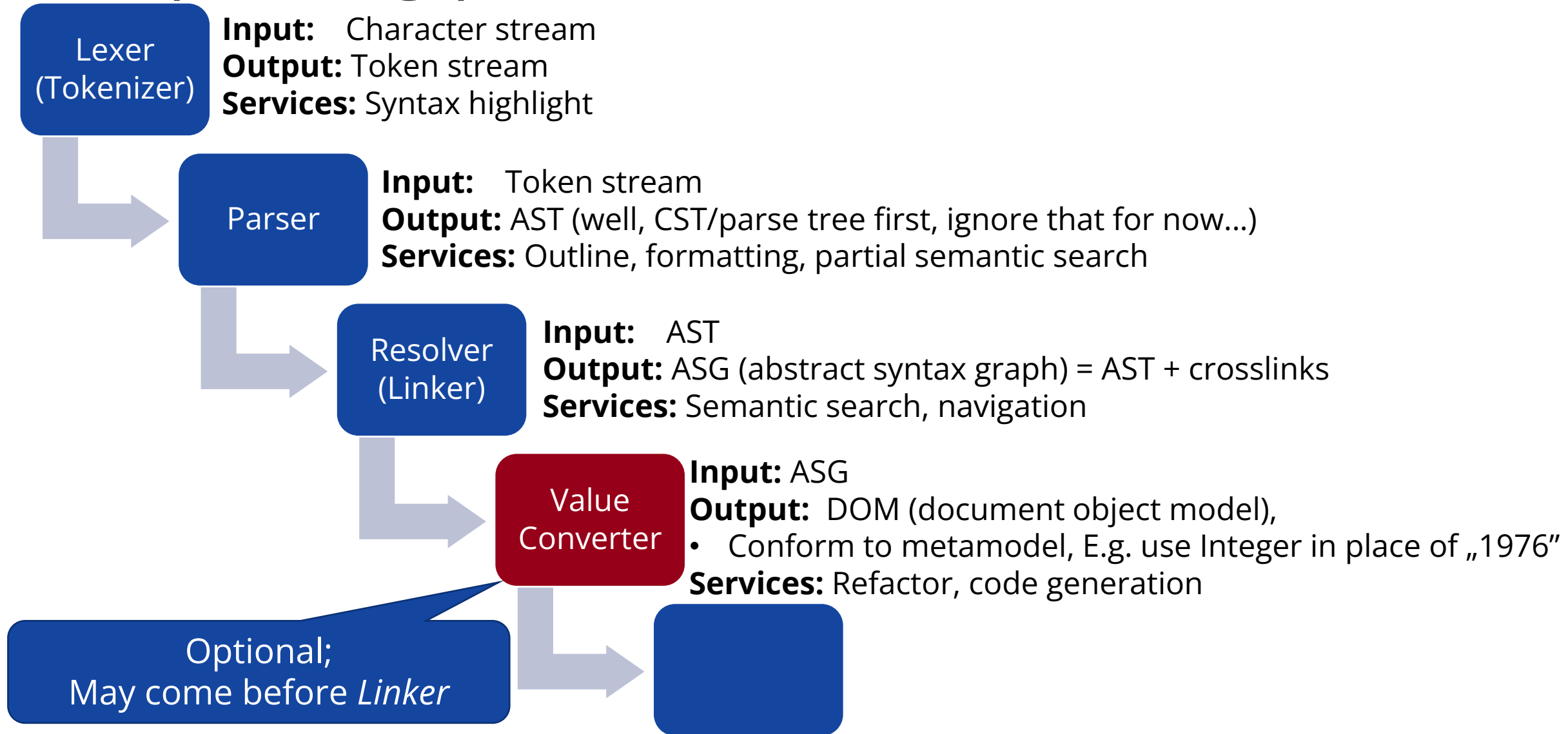
- During parsing (resolving/linking),
- During editing (content assist)

```
private int value;  
  
public void setValue(int value)  
{  
    this.value = value;  
}
```



Which variable declaration is referred by 'value'?

The parsing process



The parsing process

