

COMPILER CONSTRUCTION

# Parsing & Lexing

Hannes Saffrich

University of Freiburg  
Department of Computer Science  
Programming Languages

26. April 2024

# Parser

- ▶ A parser checks if a word is part of a language
- ▶ An alphabet  $\Sigma$  is a set
- ▶ The elements  $a \in \Sigma$  are called letters or *terminal symbols*
- ▶ A word  $w$  is a list of letters, i.e.  $w \in \Sigma^*$
- ▶ A language  $\mathcal{L}$  is a set of words, i.e.  $\mathcal{L} \subseteq \Sigma^*$
- ▶ A word  $w$  is part of the language  $\mathcal{L}$  if  $w \in \mathcal{L}$
- ▶ Typically, a language is described by a grammar
- ▶ Typically, a parser produces additional information:
  - ▶ if a word is in the language, it produces one or more syntax trees
  - ▶ if a word is not in the language, it produces an error message describing why it is not in the language

# Lexer

- ▶ A lexer is an optional preprocessing step for a parser
- ▶ It translates words from one alphabet to words of another alphabet

$$\Sigma^* \rightarrow \Delta^*$$

- ▶ Typically, the input words are actual strings, i.e.  $\Sigma^* = str$
- ▶ The output letters  $t \in \Delta$  are called *Tokens*
- ▶ A lexer serves two purposes:
  - ▶ Allow the parser to be based on a more readable grammar, e.g. by removing whitespace and comments, and treating numbers or variable names as single letters
  - ▶ Increase performance as some parsing tasks can be done more efficiently in a lexer
- ▶ Typically, a lexer is described by a mapping from regular expressions to tokens

# Running Example: A Circuit Description Language

$\langle prog \rangle ::= (\langle expr \rangle ;)^*$

$\langle expr \rangle ::= \text{True}$

|  $\text{False}$

|  $\langle var \rangle$

|  $!\langle expr \rangle$

|  $\langle expr \rangle \& \langle expr \rangle$

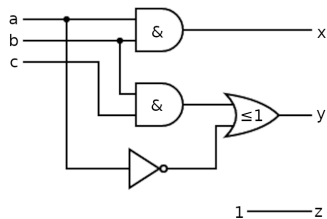
|  $\langle expr \rangle | \langle expr \rangle$

Example Program:

$a \& b;$

$!a | b \& c;$

$\text{True};$



# Lexer Example

Input String

```
"  
a & b;  
!a | b & c;  
True;  
"
```

Output Tokens

```
[  
  Id('a'), And(), Id('b'), Semicolon(),  
  Not(), Id('a'), Or(), Id('b'), And(),  
    Id('c'), Semicolon(),  
  True(), Semicolon(),  
]
```

## Lexer Description (Semi-Formal)

```
Id      = "[a-zA-Z][a-zA-Z0-9]*"  
True    = "True"  
False   = "False"  
And     = "&"  
Or      = "\\|"  
Not     = "!"  
Semicolon = ";"  
Ignore1 = "[ \\t\\n]+"Ignore2 = "#[^\n]*\n"
```

# Lexer Implementation

- ▶ Let's write a lexer!
- ▶ ... or to be more precise: a lexer library!

## Grammars: Definition

- ▶ Typically, languages are specified as grammars
- ▶ Formally, a grammar is a tuple  $(N, \Sigma, S, P)$  where
  - ▶  $N$  is a finite set of non-terminal symbols
  - ▶  $\Sigma$  is a set of terminal symbols
  - ▶  $S \in N$  is the start symbol
  - ▶  $P \subseteq (N \cup \Sigma)^* \times (N \cup \Sigma)^*$  are the production rules

- ▶ Example:

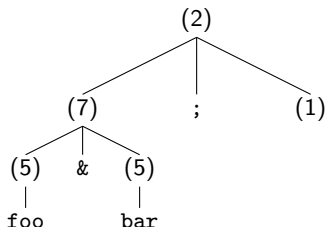
- ▶  $N = \{\langle prog \rangle, \langle expr \rangle, \langle var \rangle\}$
- ▶  $\Sigma = \{\text{True, False, ;, \&, |, a, b, ...}\}$
- ▶  $S = \langle prog \rangle$  is the start symbol
- ▶  $P$  contains the following rules:

$\langle prog \rangle \rightarrow$	$\langle expr \rangle \rightarrow !\langle expr \rangle$
$\langle prog \rangle \rightarrow \langle expr \rangle ; \langle prog \rangle$	$\langle expr \rangle \rightarrow \langle expr \rangle \& \langle expr \rangle$
$\langle expr \rangle \rightarrow \text{True}$	$\langle expr \rangle \rightarrow \langle expr \rangle   \langle expr \rangle$
$\langle expr \rangle \rightarrow \text{False}$	$\langle var \rangle \rightarrow \text{Var}(\_)$
$\langle expr \rangle \rightarrow \langle var \rangle$	

# Grammars: Language

- ▶ The language of a grammar is the set of words that can be derived from the start symbol by applying production rules

- ▶ Example:  $\langle prog \rangle \Rightarrow \langle expr \rangle ; \langle prog \rangle$   
 $\Rightarrow \langle expr \rangle ;$   
 $\Rightarrow \langle expr \rangle \& \langle expr \rangle ;$   
 $\Rightarrow \langle var \rangle \& \langle expr \rangle ;$   
 $\Rightarrow \text{foo} \& \langle expr \rangle ;$   
 $\Rightarrow \text{foo} \& \langle var \rangle ;$   
 $\Rightarrow \text{foo} \& \text{bar} ;$



- ▶ Rules:

$\langle prog \rangle \rightarrow$	(1)	$\langle expr \rangle \rightarrow ! \langle expr \rangle$	(6)
$\langle prog \rangle \rightarrow \langle expr \rangle ; \langle prog \rangle$	(2)	$\langle expr \rangle \rightarrow \langle expr \rangle \& \langle expr \rangle$	(7)
$\langle expr \rangle \rightarrow \text{True}$	(3)	$\langle expr \rangle \rightarrow \langle expr \rangle   \langle expr \rangle$	(8)
$\langle expr \rangle \rightarrow \text{False}$	(4)	$\langle var \rangle \rightarrow \text{Var}(\_)$	(9)
$\langle expr \rangle \rightarrow \langle var \rangle$	(5)		



## Grammars: EBNF Notation

- ▶ The *Extended Backus–Naur form (EBNF)* is a more concise notation for writing down the production rules of a grammar
- ▶ Multiple rules with the same left side are combined by writing the right side of the rule as a regular expression over the alphabet  $N \cup \Sigma$
- ▶ Our example grammar in EBNF notation:

$$\langle prog \rangle ::= (\langle expr \rangle ;)^*$$
$$\langle expr \rangle ::= \text{True} \mid \text{False} \mid \langle var \rangle \mid !\langle expr \rangle \mid \langle expr \rangle \& \langle expr \rangle \mid \langle expr \rangle \mid \langle expr \rangle$$
$$\langle var \rangle ::= \text{Var}(\_)$$

- ▶ Our example grammar in rule notation:

$$\langle prog \rangle \rightarrow$$
$$\langle expr \rangle \rightarrow !\langle expr \rangle$$
$$\langle prog \rangle \rightarrow \langle expr \rangle ; \langle prog \rangle$$
$$\langle expr \rangle \rightarrow \langle expr \rangle \& \langle expr \rangle$$
$$\langle expr \rangle \rightarrow \text{True}$$
$$\langle expr \rangle \rightarrow \langle expr \rangle \mid \langle expr \rangle$$
$$\langle expr \rangle \rightarrow \text{False}$$
$$\langle var \rangle \rightarrow \text{Var}(\_)$$
$$\langle expr \rangle \rightarrow \langle var \rangle$$

## Grammars: Ambiguity

- ▶ A grammar is ambiguous, if a word can be derived in multiple ways
- ▶ Our example grammar is actually ambiguous...

$$\langle prog \rangle ::= (\langle expr \rangle ;)^*$$
$$\langle expr \rangle ::= \text{True} \mid \text{False} \mid \langle var \rangle \mid !\langle expr \rangle \mid \langle expr \rangle \& \langle expr \rangle \mid \langle expr \rangle \mid \langle expr \rangle$$
$$\langle var \rangle ::= \text{Var}(\_)$$

- ▶ The word  $x \ \& \ y \ \& \ z ;$  can be derived as  $(x \ \& \ y) \ \& \ z ;$  and  $x \ \& \ (y \ \& \ z) ;$
- ▶ Even worse, the word  $x \ \& \ y \ \mid \ z ;$  can be derived as  $(x \ \& \ y) \ \mid \ z ;$  and  $x \ \& \ (y \ \mid \ z) ;$  where only the first one is valid (“and” binds stronger than “or”)
- ▶ This is no problem if the grammar only used to describe abstract syntax tree data types, but it matters for parsing
- ▶ We can fix this by slightly adjusting the grammar to capture the *concrete syntax*

## Grammars: Ambiguity (Fix 1)

- ▶ Grammar for abstract syntax (ambiguous):

$$\langle prog \rangle ::= (\langle expr \rangle ;)^*$$
$$\langle expr \rangle ::= \text{True} \mid \text{False} \mid \langle var \rangle \mid !\langle expr \rangle \mid \langle expr \rangle \& \langle expr \rangle \mid \langle expr \rangle \mid \langle expr \rangle$$
$$\langle var \rangle ::= \text{Var}(\_)$$

- ▶ Grammar for concrete syntax (unambiguous):

$$\langle prog \rangle ::= (\langle expr \rangle ;)^*$$
$$\langle expr \rangle ::= \text{True} \mid \text{False} \mid \langle var \rangle \mid (!\langle expr \rangle) \mid ((\langle expr \rangle \& \langle expr \rangle)) \mid ((\langle expr \rangle \mid \langle expr \rangle))$$
$$\langle var \rangle ::= \text{Var}(\_)$$

- ▶ This forbids `!a & b | c;` and requires us to write `(((!a) & b) | c);`

# Grammars: Ambiguity

- ▶ What we actually want is:
  - ▶  $\&$  binds stronger than  $|$ , e.g.  $x \& y | z$  means  $(x \& y) | z$
  - ▶  $!$  binds stronger than  $\&$ , e.g.  $!x \& y$  means  $(!x) \& y$
  - ▶  $\&$  is left-associative, e.g.  $x \& y \& z$  means  $(x \& y) \& z$
  - ▶  $|$  is left-associative, e.g.  $x | y | z$  means  $(x | y) | z$
  - ▶  $!$  is right-associative, e.g.  $!!x$  means  $!( ! x)$
- ▶ As  $\&$  and  $|$  are associative operators, making them left-associative is an arbitrary choice, and we could just as well make them right-associative

## Grammars: Ambiguity (Fix 2)

- ▶ Grammar for abstract syntax (ambiguous):

$$\langle prog \rangle ::= (\langle expr \rangle ;)^*$$
$$\langle expr \rangle ::= \text{True} \mid \text{False} \mid \langle var \rangle \mid !\langle expr \rangle \mid \langle expr \rangle \& \langle expr \rangle \mid \langle expr \rangle \mid \langle expr \rangle$$
$$\langle var \rangle ::= \text{Var}(\_)$$

- ▶ Grammar for concrete syntax (unambiguous):

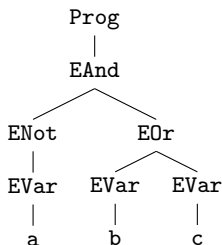
$$\langle prog \rangle ::= (\langle expr \rangle ;)^*$$
$$\langle expr \rangle ::= \langle expr \rangle \mid \langle expr1 \rangle \mid \langle expr1 \rangle$$
$$\langle expr1 \rangle ::= \langle expr1 \rangle \& \langle expr2 \rangle \mid \langle expr2 \rangle$$
$$\langle expr2 \rangle ::= !\langle expr2 \rangle \mid \langle expr3 \rangle$$
$$\langle expr3 \rangle ::= \text{True} \mid \text{False} \mid \langle var \rangle \mid (\langle expr \rangle)$$
$$\langle var \rangle ::= \text{Var}(\_)$$

- ▶ This forces `!a & b | c;` to be parsed as `(((!a) & b) | c);`

# Grammars: Abstract Syntax vs Concrete Syntax

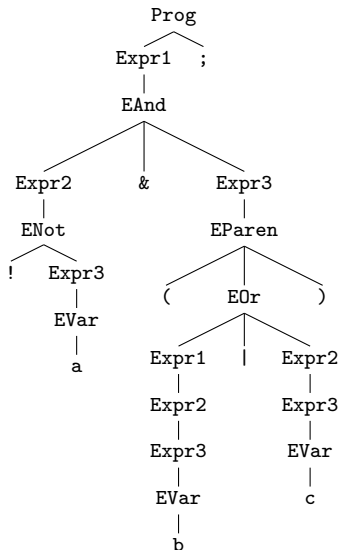
Example word: !a & (b | c);

Abstract Syntax Tree:



In most scenarios, we want to parse according to the concrete syntax, but have the parser generate an abstract syntax tree → attribute grammars

Concrete Syntax Tree:



# Grammars: Classification

- ▶ The *Chomsky Hierarchy* classifies languages by the kind of production rules that are required to describe them:

- ▶ *Recursively enumerable* languages can be described by rules of the form:

$$\alpha \rightarrow \beta \quad \forall \alpha, \beta \in (N \cup \Sigma)^*, \alpha \text{ not empty}$$

- ▶ Parsing is semi-decidable
- ▶ *Context-sensitive* languages can be described by rules of the form:

$$\alpha A \gamma \rightarrow \alpha \beta \gamma \quad \forall A \in N, \alpha, \beta, \gamma \in (N \cup \Sigma)^*, \beta \text{ not empty}$$

- ▶ Parsing is PSPACE-complete
- ▶ *Context-free* languages can be described by rules of the form:

$$A \rightarrow \alpha \quad \forall A \in N, \alpha \in (N \cup \Sigma)^*$$

- ▶ Parsing with  $O(n^3)$  worst-case time complexity
- ▶ *Regular languages* can be described by rules of the form:

$$A \rightarrow a \quad A \rightarrow aB \quad \forall A, B \in N, a \in \Sigma$$

- ▶ Parsing with  $O(n)$  worst-case time complexity

## Grammars: Classification

- ▶ We are only concerned with regular and context free languages
  - ▶ Regular languages for lexing
  - ▶ Context-free languages for parsing
- ▶ The majority of programming languages can be described by a context-free grammar\*
- ▶ There are many subclasses of context-free languages, which are more restrictive, but have better time and space complexity
  - ▶ e.g. the language classes LL, LR, LALR can be parsed in linear time, but are still powerful enough to parse Java
  - ▶ writing grammars in these subclasses can be more annoying, because less shapes of production rules are allowed

\*: Some require minor preprocessing, like Python (“semantic whitespace”)



# The Earley Parser

- ▶ Can parse arbitrary context-free grammars
- ▶ Worst-case time complexity:
  - ▶  $O(n^3)$  for ambiguous grammars (we don't care about those)
  - ▶  $O(n^2)$  for unambiguous grammars
- ▶ Relatively simple (170 lines of Python)
- ▶ Can be reused as it is basically a grammar interpreter
- ▶ Nice for prototyping as it supports general context-free grammars
- ▶ Probably a bit slow for single files with millions of tokens
- ▶ Produces multiple syntax trees when used with ambiguous grammars

# The Earley Parser: Basic Principle

- ▶ Loops once over each token of the input word
- ▶ Tracks which production rules made how much progress in a chart
- ▶ A *chart* maps each token index to a set of dotted rules
- ▶ A *dotted rule* is a production rule, which
  - ▶ has a marker (dot) in its right side that denotes how much of that rule was matched by the input so far
  - ▶ is annotated with the token index at which it was first added to the chart
- ▶ `chart[0]` is initialized with dotted rules derivable from the start symbol
- ▶ For each token `t` we compute `chart[i]` from `chart[i-1]` by checking, which rule expects a `t` after the dot and moving their dot one symbol to the right [some details omitted]
- ▶ A word is accepted, if the last entry of `chart` contains a rule with the start symbol on the left side and a dot at the end of its right side

# The Earley Parser: Example (1 / 7)

- ▶ Input word: "True & False ;"
- ▶ Step 1: Add rules with the start symbol to the chart:

$$\text{chart}[0] = \{ (0, \langle \text{prog} \rangle \rightarrow \bullet), \\ (0, \langle \text{prog} \rangle \rightarrow \bullet \langle \text{expr} \rangle ; \langle \text{prog} \rangle) \}$$

- ▶ Step 2: For each dotted rule in  $\text{chart}[0]$ , also add the rules for all non-terminals which come immediately after a dot:

$$\text{chart}[0] = \text{chart}[0] \cup \{ (0, \langle \text{expr} \rangle \rightarrow \bullet \text{True}), \\ (0, \langle \text{expr} \rangle \rightarrow \bullet \text{False}), \\ (0, \langle \text{expr} \rangle \rightarrow \bullet \langle \text{var} \rangle), \\ (0, \langle \text{expr} \rangle \rightarrow \bullet ! \langle \text{expr} \rangle), \\ (0, \langle \text{expr} \rangle \rightarrow \bullet \langle \text{expr} \rangle \& \langle \text{expr} \rangle), \\ (0, \langle \text{expr} \rangle \rightarrow \bullet \langle \text{expr} \rangle | \langle \text{expr} \rangle) \}$$

- ▶ Step 3: Repeat until the set doesn't change anymore:

$$\text{chart}[0] = \text{chart}[0] \cup \{ (0, \langle \text{var} \rangle \rightarrow \bullet \text{Var}(\_)) \}$$

## The Earley Parser: Example (2 / 7)

- ▶ Input word: "True & False ;"
- ▶ Step 4: The first input token is True, so search for dotted rules in `chart[0]`, which have an True after the dot:

$$(0, \langle expr \rangle \rightarrow \bullet True)$$

- ▶ Move the dot after the matched symbol:

$$(0, \langle expr \rangle \rightarrow True \bullet)$$

- ▶ Add the modified rules to `chart[1]`:

$$\text{chart}[1] = \{ (0, \langle expr \rangle \rightarrow True \bullet) \}$$

## The Earley Parser: Example (3 / 7)

- ▶ Input word: "True & False ;"
- ▶ Step 5: Check for completed rules in  $\text{chart}[1]$ , i.e. rules which have the dot at the end:

$$(0, \langle \text{expr} \rangle \rightarrow \text{True} \bullet)$$

- ▶ This rule was created at the beginning (0), so search for rules in  $\text{chart}[0]$ , which have  $\langle \text{expr} \rangle$  after the dot:

$$(0, \langle \text{expr} \rangle \rightarrow \bullet \langle \text{expr} \rangle \& \langle \text{expr} \rangle)$$

$$(0, \langle \text{expr} \rangle \rightarrow \bullet \langle \text{expr} \rangle | \langle \text{expr} \rangle)$$

- ▶ Move the dot after  $\langle \text{expr} \rangle$  and add them to  $\text{chart}[1]$ :

$$\text{chart}[1] = \text{chart}[1] \cup \{ (0, \langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle \bullet \& \langle \text{expr} \rangle) \\ (0, \langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle \bullet | \langle \text{expr} \rangle) \}$$

## The Earley Parser: Example (4 / 7)

- ▶ Input word: "True & False ;"
- ▶ Check if there are rules in `chart[1]` where the dot is in front of a non-terminal, and add the rules for the non-terminal also to `chart[1]` (same as for `chart[0]` in the beginning).
- ▶ There are no such rules.
- ▶ Repeat Step 5 until the `chart[1]` doesn't change anymore.
- ▶ This is already the case.

## The Earley Parser: Example (5 / 7)

- ▶ Input word: "True & False ;"
- ▶ Step 6: The next input token is & so check for rules in chart [1] with an & after the dot:

$$(0, \langle expr \rangle \rightarrow \langle expr \rangle \bullet \& \langle expr \rangle)$$

- ▶ Move the • to the right and add to chart [2]:

$$\text{chart}[2] = \{ (0, \langle expr \rangle \rightarrow \langle expr \rangle \& \bullet \langle expr \rangle) \}$$

- ▶ Step 7: Repeat Step 5 for chart [2] until it doesn't change anymore:

$$\begin{aligned} \text{chart}[2] = \text{chart}[2] \cup \{ & (2, \langle expr \rangle \rightarrow \bullet \text{True}), \\ & (2, \langle expr \rangle \rightarrow \bullet \text{False}), \\ & (2, \langle expr \rangle \rightarrow \bullet \langle var \rangle), \\ & (2, \langle expr \rangle \rightarrow \bullet ! \langle expr \rangle), \\ & (2, \langle expr \rangle \rightarrow \bullet \langle expr \rangle \& \langle expr \rangle), \\ & (2, \langle expr \rangle \rightarrow \bullet \langle expr \rangle | \langle expr \rangle), \\ & (2, \langle var \rangle \rightarrow \bullet \text{Var}(\_)) \} \end{aligned}$$

## The Earley Parser: Example (6 / 7)

- ▶ Input word: "True & False ;"
- ▶ Step 8: The next input token is False so check for rules in chart [2] with an & after the dot:

$$(2, \langle expr \rangle \rightarrow \bullet False)$$

- ▶ Move the • to the right and add to chart [3]:

$$\text{chart}[3] = \{ (2, \langle expr \rangle \rightarrow False \bullet) \}$$

- ▶ Step 9: Repeat Step 5 for chart [3] until it doesn't change anymore:

$$\begin{aligned} \text{chart}[3] = \text{chart}[3] \cup \{ & (0, \langle expr \rangle \rightarrow \langle expr \rangle \& \langle expr \rangle \bullet), \\ & (2, \langle expr \rangle \rightarrow \langle expr \rangle \bullet \& \langle expr \rangle), \\ & (2, \langle expr \rangle \rightarrow \langle expr \rangle \bullet | \langle expr \rangle), \\ & (0, \langle prog \rangle \rightarrow \langle expr \rangle \bullet ; \langle prog \rangle) \} \end{aligned}$$



## The Earley Parser: Example (7 / 7)

- ▶ Input word: "True & False ;"
- ▶ Step 10: The next input token is ; so check for rules in chart [3] with an ; after the dot:

$$(0, \langle prog \rangle \rightarrow \langle expr \rangle \bullet ; \langle prog \rangle)$$

- ▶ Move the • to the right and add to chart [4]:

$$\text{chart}[4] = \{ (0, \langle prog \rangle \rightarrow \langle expr \rangle ; \bullet \langle prog \rangle) \}$$

- ▶ Step 11: Repeat Step 5 for chart [4] until it doesn't change anymore:

$$\begin{aligned} \text{chart}[4] = \text{chart}[4] \cup \{ & (4, \langle prog \rangle \rightarrow \bullet), \\ & (4, \langle prog \rangle \rightarrow \bullet \langle expr \rangle ; \langle prog \rangle), \\ & (4, \langle expr \rangle \rightarrow \bullet \text{True}), \\ & \dots \\ & (0, \langle prog \rangle \rightarrow \langle expr \rangle ; \langle prog \rangle \bullet) \} \end{aligned}$$

- ▶ chart [4] contains a rule from the beginning (0) with the start symbol on the left side and a dot at the end of the right side, hence the input word is in the language.