



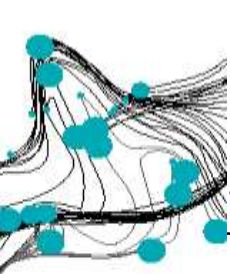
COMPILER CONSTRUCTION:

PARSE TREES & ANTLR (PRIMER)

MODULE 8: PROGRAMMING PARADIGMS

29 APRIL 2019



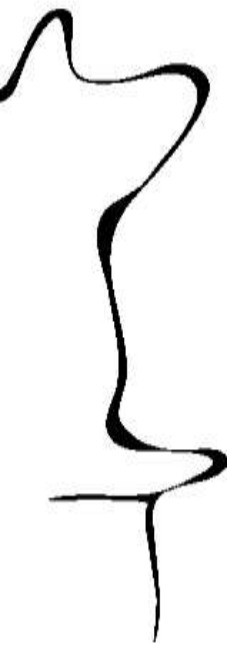


SEEN BEFORE

- Structure of compiler
- Front end
 - Scanning
 - Parsing
 - Type checking
- Optimizer
- Back end

← To be continued today

Parser: turns a sequence of tokens of language-specific token types into a parse tree



(EXTENDED) BACKUS-NAUR FORM ((E)BNF)

Context-free grammar rules

- *Original:*
 - Left hand side (LHS): Nonterminal
 - Right hand side (RHS): sequence of terminals + nonterminals
 - Multiple rules for one non-terminal, separated by |

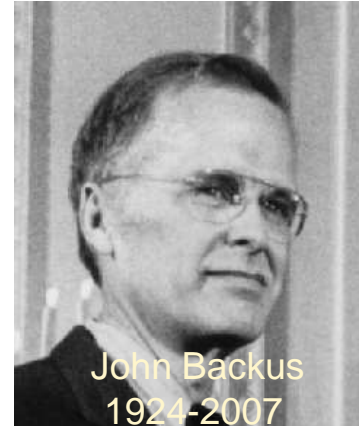
$SheepNoise \rightarrow baa\ SheepNoise$

Terminal | Nonterminal

Diagram description: The word 'baa' is circled in red and labeled 'Terminal'. The word 'SheepNoise' is also circled in red and labeled 'Nonterminal'. A vertical line with a horizontal bar at the bottom is positioned between 'baa' and 'SheepNoise', with the label 'Terminal' to its left and 'Nonterminal' to its right.

- *Extended:* use RE syntax in RHS
 - Can always be rewritten to basic BNF

$SheepNoise \rightarrow baa +$

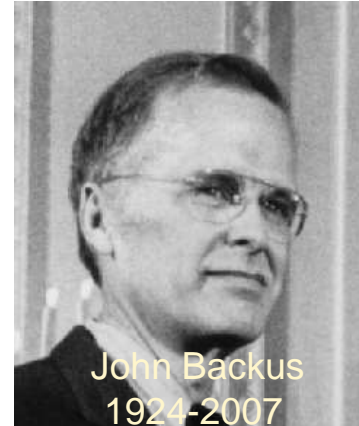


Example:

- EBNF
 $A \rightarrow B? ((a | b) C)^*$
- Corresponding BNF?

Source: John Backus: www.columbia.edu
Peter Naur: Erik tj, <http://commons.wikimedia.org>

(EXTENDED) BACKUS-NAUR FORM ((E)BNF)

$$A \rightarrow BD \mid D$$
$$D \rightarrow aC \mid bC \mid \epsilon \mid DD$$


Example:

- EBNF

$$A \rightarrow B? ((a \mid b) C)^*$$

- Corresponding BNF?

$$\begin{array}{l} A \rightarrow BD \quad D \rightarrow aCD \\ \quad \mid D \quad \quad \mid bCD \\ \quad \quad \quad \quad \mid \epsilon \end{array}$$

Source: John Backus: www.columbia.edu
Peter Naur: Eriktj, <http://commons.wikimedia.org>

DERIVATIONS AND PARSE TREES

- Derivation: sequence of rule applications
 - Starting at start symbol (duh)
 - Ending at term to be parsed (“accepted”, “generated”): *sentence*
 - Intermediate forms: *sentential form* (mixture of terminals and non-terminals)
- Choice of where to apply next rule
 - Leftmost / rightmost non-terminal
- Parse tree: node for each rule application
 - Children are the symbols of the RHS

Parse tree
Expr

1	<i>Expr</i>	→	(<i>Expr</i>)
2			<i>Expr Op</i> name
3			name
4	<i>Op</i>	→	+
5			-
6			×
7			÷

Rule	Sentential Form
	<i>Expr</i>

Rightmost Derivation of (a + b) × c

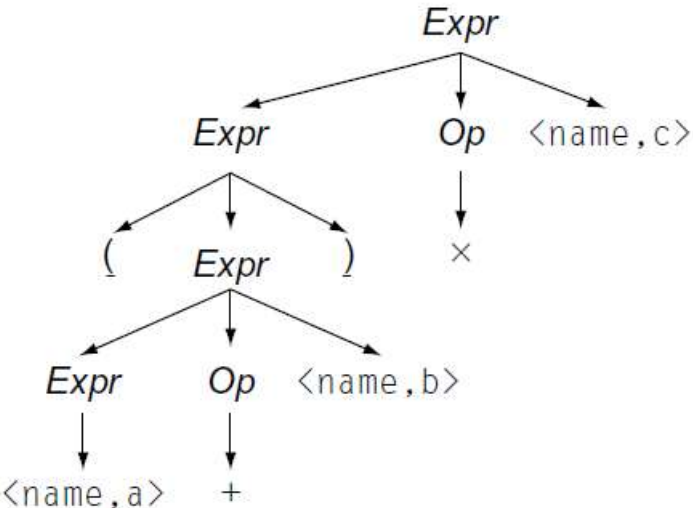
DERIVATIONS AND PARSE TREES

- Derivation: sequence of rule applications
 - Starting at start symbol (duh)
 - Ending at term to be parsed (“accepted”, “generated”): *sentence*
 - Intermediate forms: *sentential form* (mixture of terminals and non-terminals)
- Choice of where to apply next rule
 - Leftmost / rightmost non-terminal
- Parse tree: node for each rule application
 - Children are the symbols of the RHS

1	<i>Expr</i>	→	(<i>Expr</i>)
2			<i>Expr Op name</i>
3			name
4	<i>Op</i>	→	+
5			-
6			×
7			÷

Rule	Sentential Form
	<i>Expr</i>
2	<i>Expr Op name</i>
6	<i>Expr</i> × name
1	(<i>Expr</i>) × name
2	(<i>Expr Op name</i>) × name
4	(<i>Expr</i> + name) × name
3	(name + name) × name

Parse tree



Rightmost Derivation of (a + b) × c

AMBIGUITY

1) $a=2$
 2) $a=1$
 $\wedge b=1$

```

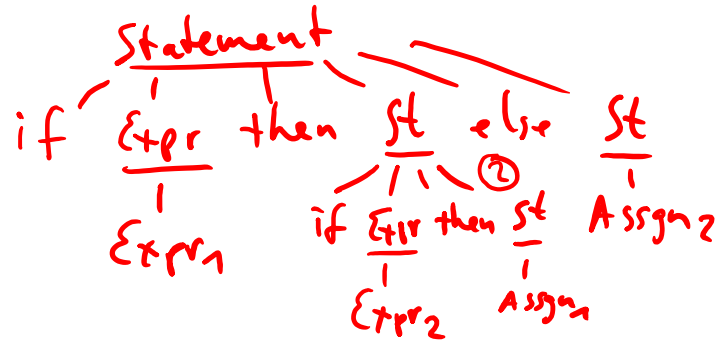
if(a==1)
{ if(b==2)
  c=2;
  else c=3; }
    
```

Let's parse
 using grammar

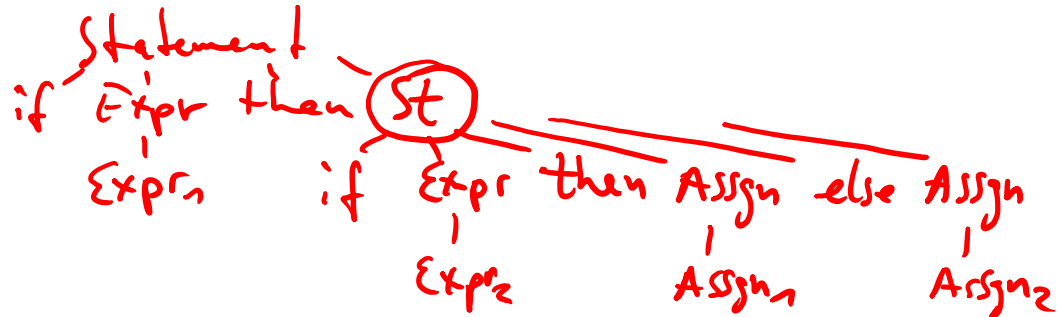
$|$ if $Expr_1$ then if $Expr_2$ then $Assignment_1$ else $Assignment_2$

1	Statement	→	if Expr then Statement else Statement	1
2			if Expr then Statement	2
3			Assignment	3
4			... other statements ...	

▪ Option 1:



▪ Option 2:



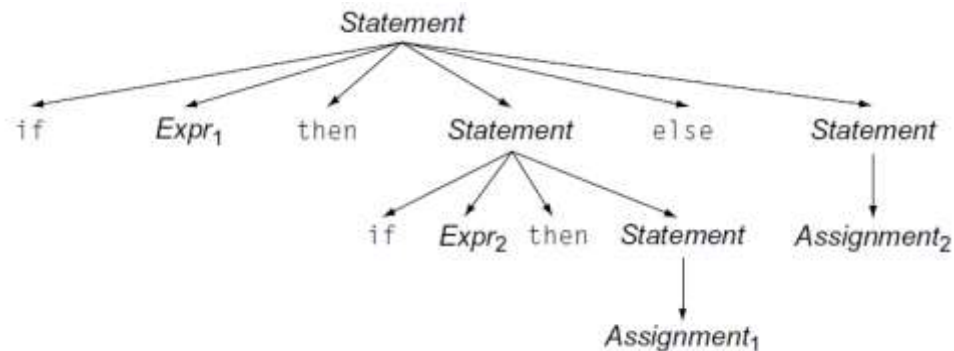
AMBIGUITY

Let's parse
using grammar

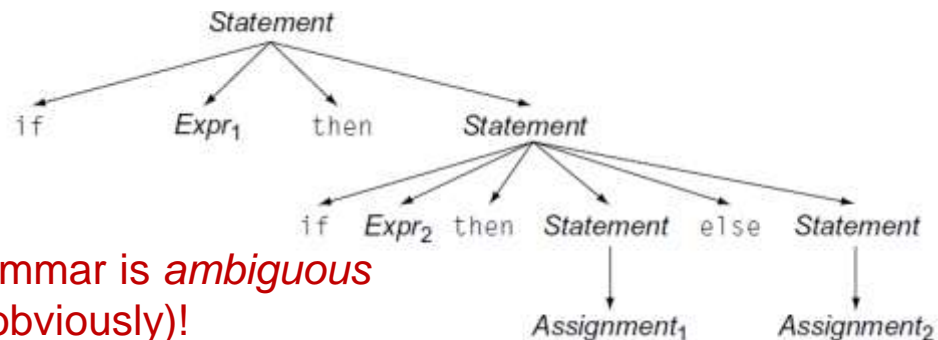
if Expr₁ then if Expr₂ then Assignment₁ else Assignment₂

- | | | | |
|---|------------------|---|--|
| 1 | <i>Statement</i> | → | <i>if Expr then Statement else Statement</i> |
| 2 | | | <i>if Expr then Statement</i> |
| 3 | | | <i>Assignment</i> |
| 4 | | | <i>... other statements ...</i> |

▪ Option 1:



▪ Option 2:



No unique answer: grammar is *ambiguous*
This is to be avoided (obviously!)
(for instance: give priority to one of the cases)

LEFT- AND RIGHT-ASSOCIATIVITY

- What is the difference between

$$\begin{array}{l} Expr \rightarrow Expr + Operand \\ | Expr - Operand \\ | Operand \end{array}$$

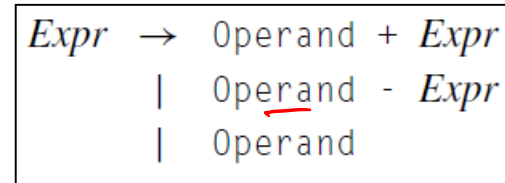
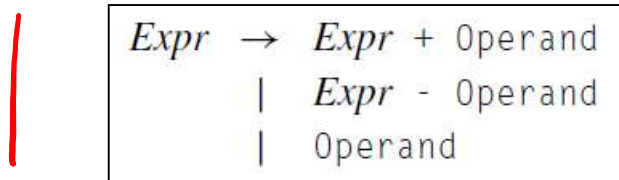
$$\begin{array}{l} Expr \rightarrow Operand + Expr \\ | Operand - Expr \\ | Operand \end{array}$$

- They accept the same language
- Consider the text: $1 - 3 + 4 - 2 = 0$ - + : left-associative
 - What should the outcome be?

LEFT- AND RIGHT-ASSOCIATIVITY

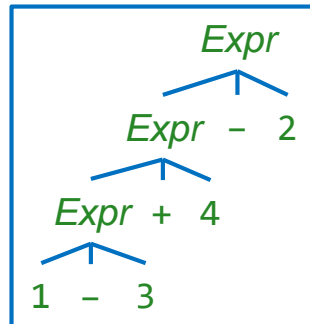
$$2^{3^4} \quad \underline{\underline{2^{(3^4)}}} \quad (2^3)^4$$

- What is the difference between

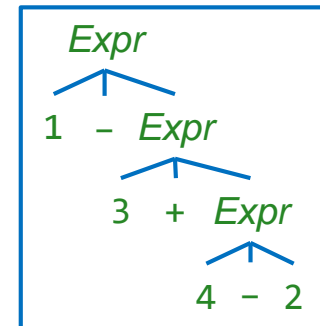


- They accept the same language
- Consider the text: $1 - 3 + 4 - 2$
 - What should the outcome be?

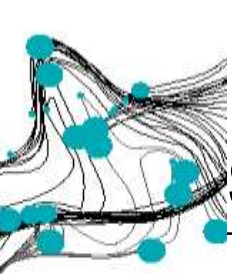
Parse tree
left hand grammar:



Parse tree
right hand grammar:



- The shape of the parse tree matters for the meaning = -4
 - For $-$ and $+$ (*left associative* operators) the left hand tree is correct
 - For $^$ (exponentiation, *right associative*) the right hand tree would be correct



SEEN BEFORE

- Structure of compiler
 - Front end
 - Scanning
 - Parsing
 - Type checking
 - Optimizer
 - Back end
- } Primer on ANTLR



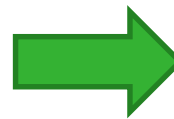
ANTLR LEXER GRAMMARS

MyFirstLanguage.g4

= scanner **lexer** grammar MyFirstLanguage;

```
KEYWORD : 'key';
```

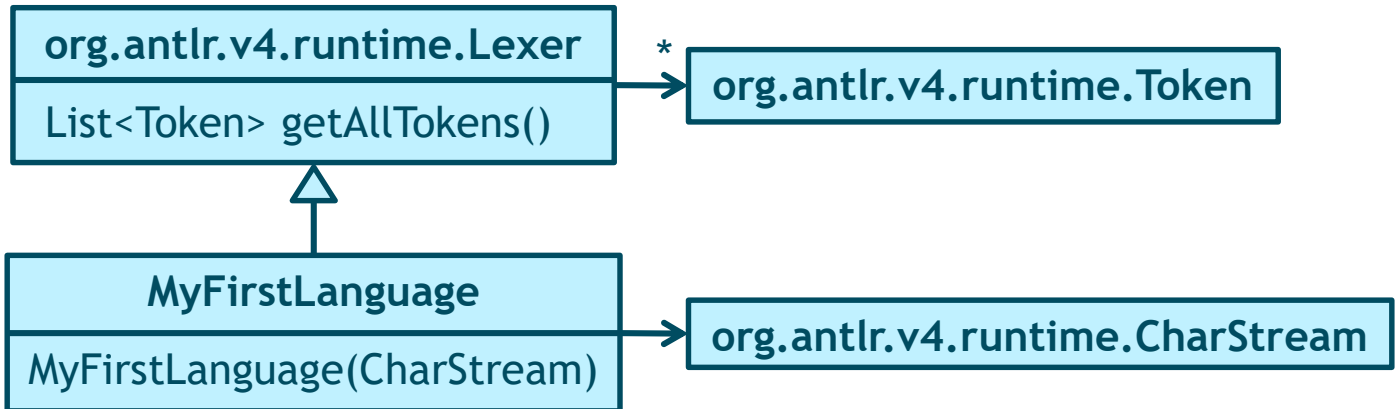
Antlr



MyFirstLanguage.tokens

MyFirstLanguage.java

Structure:



Usage:

```
CharStream stream = CharStreams.fromString(text);
Lexer lexer = new MyFirstLanguage(stream);
for (Token token : lexer.getAllTokens()) {
    System.out.print(token.toString() + " ");
}
```

PRAGMATICS OF ANTLR

- Generated Java files have to be in some package
 - Declare in Antlr grammar file:

```
lexer grammar MyFirstLanguage;  
@header{package pp.slides.cc1;}  
KEYWORD : 'key';
```

Added to header
of MyFirstLanguage.java

- Automatic generation in Eclipse plugin sometimes hard to convince
 - Context menu: “Run as -> Generate ANTLR Recognizer” *usually* works
 - Console window:

```
ANTLR Tool v4.7 (...)  
MyFirstLanguage.g4 -o . -listener -no-visitor -encoding UTF-8
```

This is the output directory option

- However, IntelliJ does *not* need the **package** declaration in **@header!**
 - There is no way to get compatible behaviour; strip the .g4-file!

ANTLR REGULAR EXPRESSION SYNTAX

Lexer rules always start with uppercase

- Single quotes for literal characters or character sequences

```
KEYWORD : 'key' ;  
LBRACE  : '{'  ;  
NEWLINE : '\n' ;
```

- Choice specified in multiple ways:

```
LETDIG : '0'..'9' | 'a'..'z' ;  
HEX    : [0123456789ABCDEF] ;  
NO_QUOTES : ~('\'' | '\"') ;
```

Operator or range (char ordering!)
No quotes! Rule order is important
Anything but the single or double quote

- Sequencing, repetition (zero-or-more, at least once), optional

```
NUMBER : '0' | '1'..'9' ('0'..'9')* ;  
WORD   : ('a'..'z' | 'A'..'Z')+ ;  
OPT    : 'end' ' please'? ;
```

- Special features

```
fragment LETTER : ('a'..'z' | 'A'..'Z') ;  
WORD : LETTER+ ;  
WS : [ \t\r\n ]+ -> skip ;  
STRING : '\"' .*? '\"' ;
```

Auxiliary rule, *not* a token type

Skip: not emitted as token
Turn non-greedy at *? (= any)

JAVA BLOCK COMMENTS

- Correct solution:

```
/* (^* | *+ ^)* *+ /
```

- In Antlr:

```
COMMENT : '/'* ( ~'*' | '*'+ ~'/' )* '*'+ '/' ;
```

- Or, using “non-greedy” special feature of Antlr:

```
COMMENT : '/'* .*? '*/' ;
```

NEXT STEP: ANTLR PARSER GENERATION

parse
rules:
start
with
lower
case

```
parser grammar MyLanguageParser1;  
options { tokenVocab = MyLanguage; }  
@header{package pp.slides.cc1;}  
main : BEGIN stat* END ;  
stat : RETURN SEMI ;
```

specify tokens

```
lexer grammar MyLanguage;  
@header{package pp.slides.cc1;}  
RETURN : 'return' ;  
BEGIN : '{' ;  
END : '}' ;  
SEMI : ';' ;  
WS : [ \r\n\t ]+ -> skip ;
```

Alternative: combined grammar

```
grammar MyLanguageParser2;  
import MyLanguage;  
main : BEGIN stat* END ;  
stat : RETURN SEMI ;
```

Or everything in single file:

```
grammar MyLanguageCombined;  
@header{package pp.slides.cc1;}  
main : BEGIN stat* END ;  
stat : 'return' ';' ; unnamed terminal  
BEGIN : '{' ;  
END : '}' ;  
WS : [ \r\n\t ]+ -> skip ;
```