

Compilation

0368-3133 (Semester A, 2013/14)

Lecture 3: Syntax Analysis (Top-Down Parsing)

Modern Compiler Design: Chapter 2.2

Noam Rinetzky

Admin

- Slides: <http://www.cs.tau.ac.il/~maon/...>
 - All info: [Moodle](#)
- Next week: Dan David 001
 - Vote ...
- Mobiles ...

What is a Compiler?

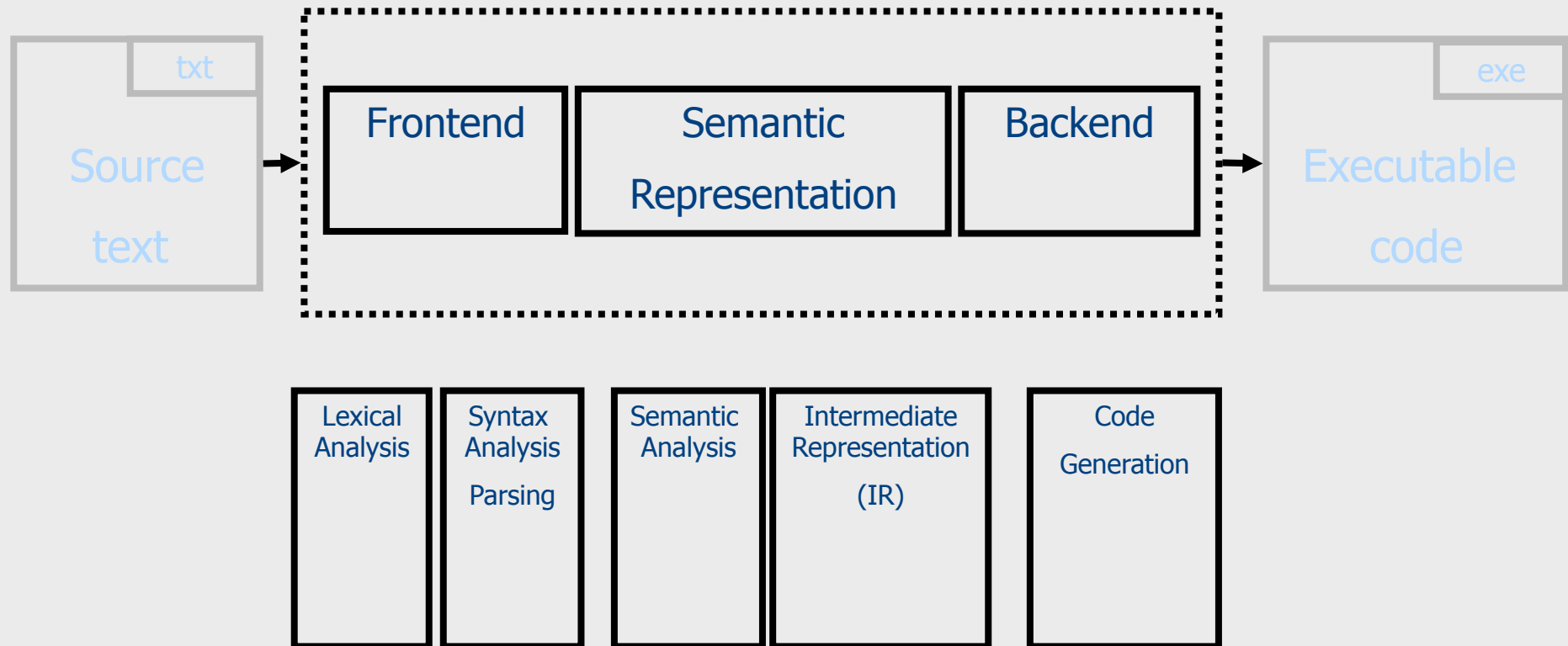
“A compiler is a **computer program** that **transforms** source code written in a programming language (**source language**) into another language (**target language**).

The most common reason for wanting to transform source code is to create an **executable program**.”

--Wikipedia

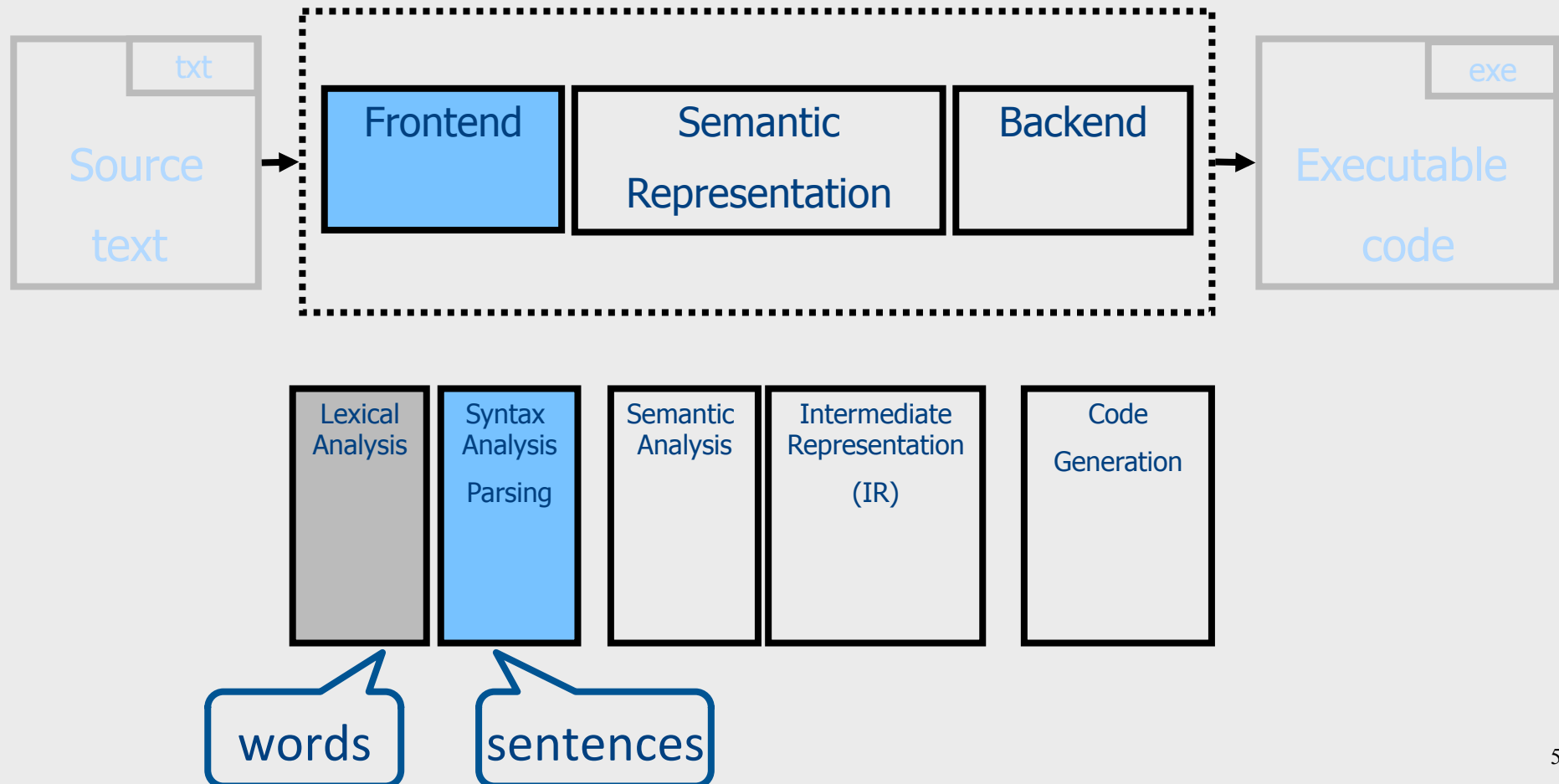
Conceptual Structure of a Compiler

Compiler



Conceptual Structure of a Compiler

Compiler



Lexical Analysis Refresher

- Example: Lang. of fully parenthesized exprs

Expr \rightarrow Num | LP Expr Op Expr RP

Num \rightarrow Dig | Dig Num

Dig \rightarrow '0' | '1' | '2' | '3' | '4' | '5' | '6' | '7' | '8' | '9'

LP \rightarrow '('

RP \rightarrow ')'

Op \rightarrow '+' | '*'

((23 + 7) * 19)

What does Lexical Analysis do?

- Language: fully parenthesized expressions

Context free
language

$\text{Expr} \rightarrow \text{Num} \mid \text{LP Expr Op Expr RP}$

Regular
languages

$\text{Num} \rightarrow \text{Dig} \mid \text{Dig Num}$

$\text{Dig} \rightarrow '0' \mid '1' \mid '2' \mid '3' \mid '4' \mid '5' \mid '6' \mid '7' \mid '8' \mid '9'$

$\text{LP} \rightarrow '('$

$\text{RP} \rightarrow ')'$

$\text{Op} \rightarrow '+' \mid '*'$

((23 + 7) * 19)

What does Lexical Analysis do?

- Language: fully parenthesized expressions

Context free language

$\text{Expr} \rightarrow \text{Num} \mid \text{LP Expr Op Expr RP}$

Regular languages

$\text{Num} \rightarrow \text{Dig} \mid \text{Dig Num}$

$\text{Dig} \rightarrow '0' \mid '1' \mid '2' \mid '3' \mid '4' \mid '5' \mid '6' \mid '7' \mid '8' \mid '9'$

$\text{LP} \rightarrow '('$

$\text{RP} \rightarrow ')'$

$\text{Op} \rightarrow '+' \mid '*'$

	Token	Token	...						
Value	((23	+	7)	*	19)
Kind	LP	LP	Num	Op	Num	RP	Op	Num	RP

Scanners: Spec. & Impl.

if	{ return IF; }
[a-z][a-z0-9]*	{ return ID; }
[0-9]+	{ return NUM; }
[0-9]"."[0-9]* [0-9]*"."[0-9]+	{ return REAL; }
(\\-\\-[a-z]*\\n) (“ ” \\n \\t)	{ ; }
.	{ error(); }

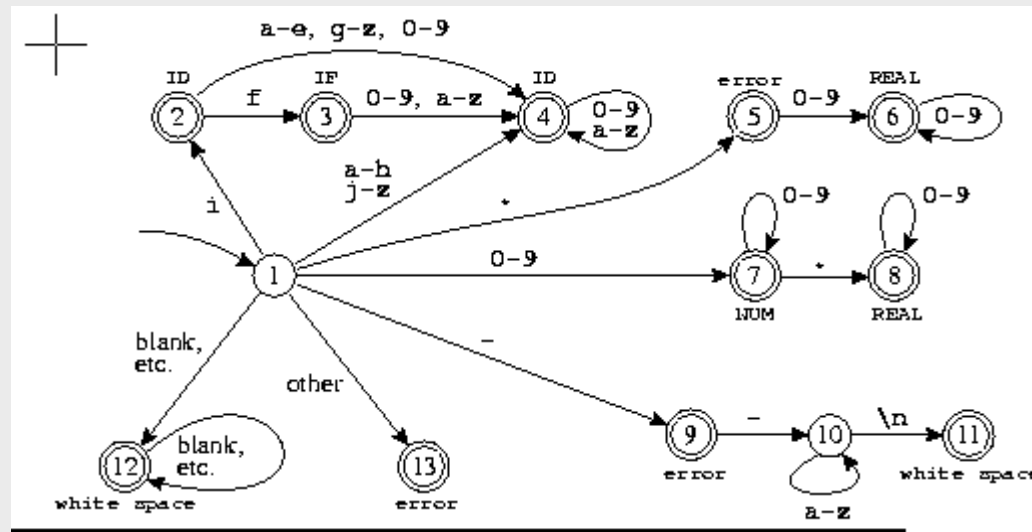


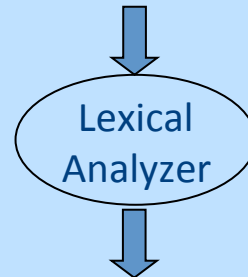
FIGURE 2.4. Combined finite automaton.
 From *Modern Compiler Implementation in ML*,
 Cambridge University Press, ©1998 Andrew W. Appel

Scanner usage of DFSM goes beyond language recognition

Frontend: Scanning & Parsing

program text

((23 + 7) * x)



token stream

((23	+	7)	*	x)
LP	LP	Num	OP	Num	RP	OP	Id	RP

Grammar:

$E \rightarrow \dots \mid \text{Id}$

$\text{Id} \rightarrow \text{'a'} \mid \dots \mid \text{'z'}$



syntax error

valid

Op(*)

Abstract Syntax Tree

Op(+)

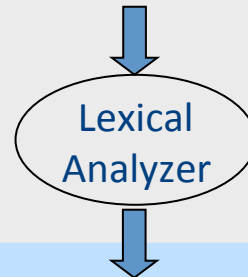
Id(b)

Num(23) Num(7)

From scanning to parsing

program text

((23 + 7) * x)



token stream

((23	+	7)	*	x)
LP	LP	Num	OP	Num	RP	OP	Id	RP



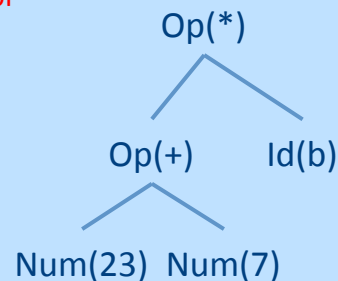
Grammar:

$E \rightarrow \dots \mid \text{Id}$

$\text{Id} \rightarrow \text{'a'} \mid \dots \mid \text{'z'}$

syntax error

valid



Abstract Syntax Tree


Parsing

- Goals
 - Is a sequence of tokens a valid program in the **language**?
 - Construct a **structured representation** of the input text
 - Error detection and reporting
- Challenges
 - How do you describe the programming language?
 - How do you check validity of an input?
 - How do you construct the structured representation?
 - Where do you report an error?

Lecture Outline

- ✓ Role & place of syntax analysis
- Context free languages
 - Grammars
 - Push down automata (PDA) ★
- Predictive parsing
- Error handling
- Semantics actions
- Earley parser

Context free languages

- Example: $\{ 0^n 1^n \mid n > 0 \}$
- Context Free Grammars (CFG)
- 
- Push Down Automata (PDA)
 - Non Deterministic Pushdown Automata

Context free grammars (CFG)

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

- **V** – non terminals (syntactic variables)
- **T** – terminals (tokens)
- **P** – derivation rules
 - Each rule of the form $V \rightarrow (T \cup V)^*$
- **S** – start symbol

What can CFGs do?

- Recognize CFLs
- $S \rightarrow 0T1$
- $T \rightarrow 0T1 \mid \epsilon$

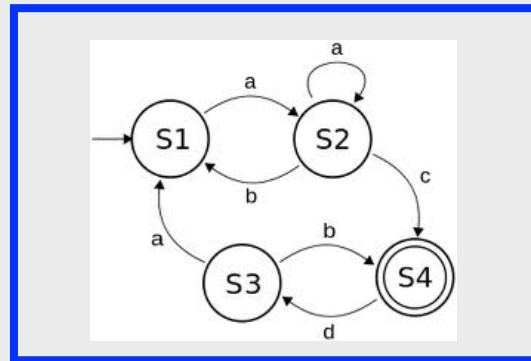
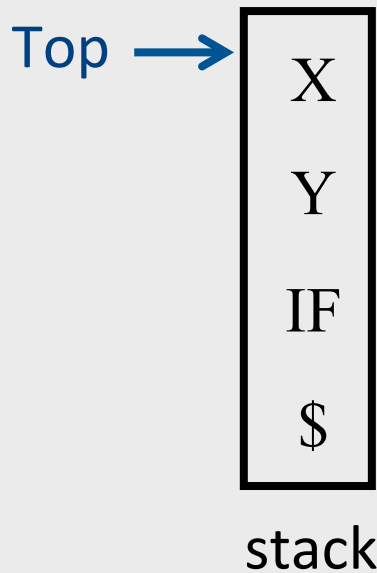
Pushdown Automata

- Nondeterministic PDA is an automaton that defines a CFL language
 - Nondeterministic PDAs define all CFLs
 - Equivalent to the CFG in language-defining power
- The deterministic version models parsers.
 - Most programming languages have deterministic PDAs
 - Efficient implementation



Intuition: PDA

- An ϵ -NFA with the additional power to manipulate **a** stack



control (ϵ -NFA)



Intuition: PDA

- Think of an ϵ -NFA with the additional power that it can manipulate a stack
- Its moves are determined by:
 - The current state (of its “ ϵ -NFA”)
 - The current input symbol (or ϵ)
 - The current symbol on top of its stack

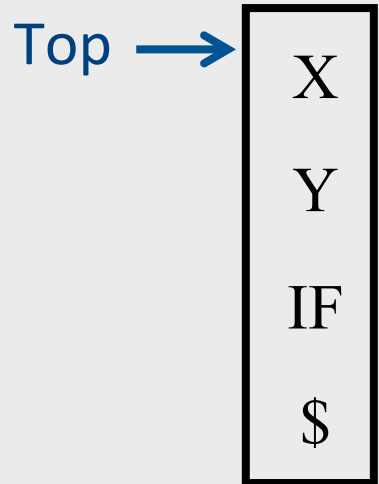


Intuition: PDA

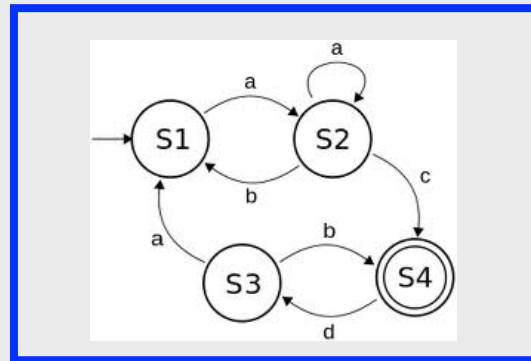
Current



input `if (oops) then stat:= blah else abort`



stack



control (ϵ -NFA)



Intuition: PDA – (2)

- Moves:
 - Change state
 - Replace the top symbol by 0...n symbols
 - 0 symbols = “pop” (“reduce”)
 - $0 < \text{symbols}$ = sequence of “pushes” (“shift”)
- Nondeterministic choice of next move



PDA Formalism

- PDA = $(A, \Sigma, \Gamma, q_0, \$, F)$:
 - A finite set of **states** (Q , typically).
 - An input **alphabet** (Σ , typically).
 - A **stack** alphabet (Γ , typically).
 - A **transition function** (δ , typically).
 - A **start state** (q_0 , in Q , typically).
 - A **start symbol** ($\$,$ in Γ , typically).
 - A set of **final states** ($F \subseteq Q$, typically).



The Transition Function

- $\delta(q, a, X) = \{ (p_1, \sigma_1), \dots, (p_n, \sigma_n) \}$
 - Input: triplet
 - A state $q \in Q$
 - An input symbol $a \in \Sigma$ or ε
 - A stack symbol $X \in \Gamma$
 - Output: set of 0 ... k **actions** of the form (p, σ)
 - A state $q \in Q$
 - σ a sequence $X_1 \dots X_n \in \Gamma^*$ of stack symbols



Actions of the PDA

- Say $(p, \sigma) \in \delta(q, a, X)$
 - If the PDA is in state q and X is the top symbol and a is at the front of the input
 - Then it can
 - Change the state to p .
 - Remove a from the front of the input
 - (but a may be ϵ).
 - Replace X on the top of the stack by σ .



Example: Deterministic PDA

- Design a PDA to accept $\{0^n 1^n \mid n > 1\}$.
- The states:
 - q = start state. We are in state q if we have seen only 0 's so far.
 - p = we've seen at least one 1 and may now proceed only if the inputs are 1 's.
 - f = final state; accept.



Example: Stack Symbols

- $\$$ = start symbol. Also marks the bottom of the stack, so we know when we have counted the same number of 1's as 0's.
- X = “counter”, used to count the number of 0's seen on the input.

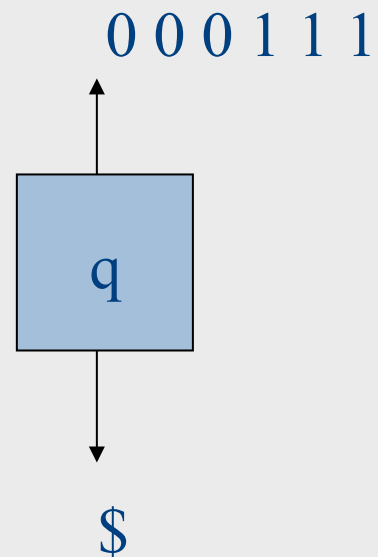


Example: Transitions

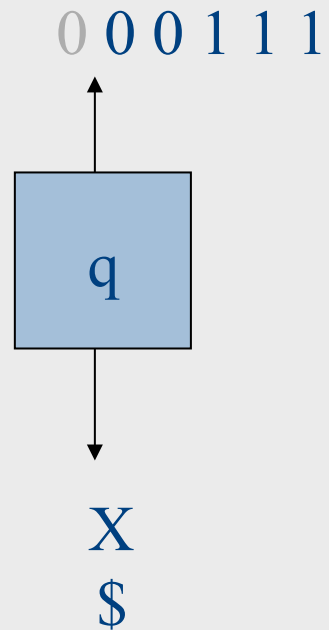
- $\delta(q, 0, \$) = \{(q, X\$)\}$.
- $\delta(q, 0, X) = \{(q, XX)\}$.
 - These two rules cause one X to be pushed onto the stack for each 0 read from the input.
- $\delta(q, 1, X) = \{(p, \varepsilon)\}$.
 - When we see a 1, go to state p and pop one X.
- $\delta(p, 1, X) = \{(p, \varepsilon)\}$.
 - Pop one X per 1.
- $\delta(p, \varepsilon, \$) = \{(f, \$)\}$.
 - Accept at bottom.



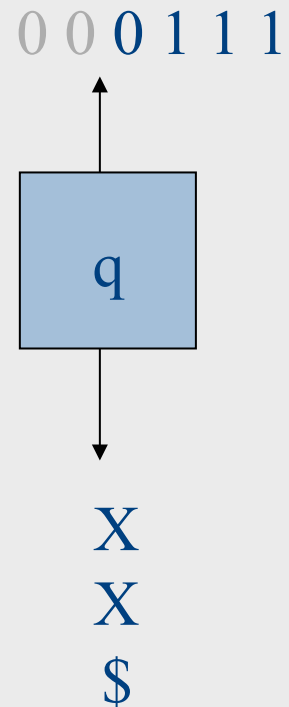
Actions of the Example PDA



Actions of the Example PDA

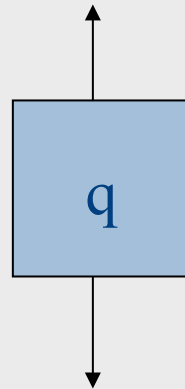


Actions of the Example PDA



Actions of the Example PDA

0 0 0 1 1 1

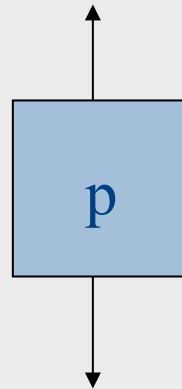


X
X
X
\$



Actions of the Example PDA

0 0 0 1 1 1

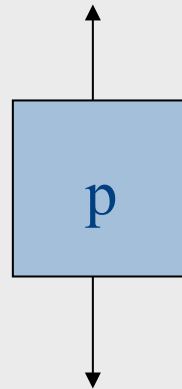


X
X
\$



Actions of the Example PDA

0 0 0 1 1 1

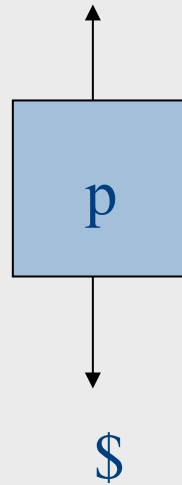


X
\$



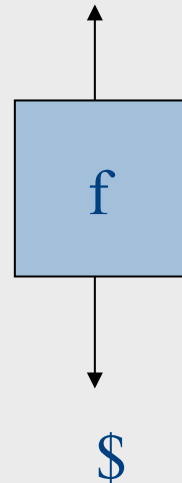
Actions of the Example PDA

0 0 0 1 1 1



Actions of the Example PDA

0 0 0 1 1 1



Example: Non Deterministic PDA

- A PDA that accepts palindromes
 - $L \{pp' \in \Sigma^* \mid p' = \text{reverse}(p)\}$



Instantaneous Descriptions

- We can formalize the pictures just seen with an instantaneous description (ID).
- A ID is a triple (q, w, α) , where:
 - q is the current state.
 - w is the remaining input.
 - α is the stack contents, top at the left.
- Define a transition relation between IDs

Context free grammars

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

- V – non terminals (syntactic variables)
- T – terminals (tokens)
- P – derivation rules
 - Each rule of the form $V \rightarrow (T \cup V)^*$
- S – start symbol

Example grammar

$S \rightarrow S ; S$

$S \rightarrow \text{id} := E$

$E \rightarrow \text{id}$

$E \rightarrow \text{num}$

$E \rightarrow E + E$

S shorthand
for Statement

E shorthand
for Expression

CFG terminology

$S \rightarrow S ; S$

$S \rightarrow \text{id} := E$

$E \rightarrow \text{id}$

$E \rightarrow \text{num}$

$E \rightarrow E + E$

Symbols:

Terminals (tokens): ; := () id num print

Non-terminals: S E L

Start non-terminal: S

Convention: the non-terminal appearing in the first derivation rule

Grammar productions (rules)

$N \rightarrow \mu$

CFG terminology

- **Derivation** - a sequence of replacements of non-terminals using the derivation rules
- **Language** - the set of strings of terminals derivable from the start symbol
- **Sentential form** - the result of a partial derivation in which there may be non-terminals

Derivations

- Show that a sentence ω is in a grammar G
 - Start with the start symbol
 - Repeatedly replace one of the non-terminals by a right-hand side of a production
 - Stop when the sentence contains only terminals
- Given a sentence $\alpha N \beta$ and rule $N \rightarrow \mu$
 $\alpha N \beta \Rightarrow \alpha \mu \beta$
- ω is in $L(G)$ if $S \Rightarrow^* \omega$

Derivation

sentence

```
x := z;  
y := x + z
```

grammar

```
S → S;S  
S → id := E | ...  
E → id | E + E | E * E | ...
```

S	$S \rightarrow S;S$
$S ; S$	$S \rightarrow id := E$
$id := E; S$	$S \rightarrow id := e$
$id := E; id := E$	$E \rightarrow id$
$id := id; id := E$	$E \rightarrow E + E$
$id := id; id := E + E$	$E \rightarrow id$
$id := id; id := E + id$	$E \rightarrow id$
$id := id; id := id + id$	
$\langle id, "x" \rangle \langle ASS \rangle \langle id, "z" \rangle \langle SEMI \rangle \langle id, "y" \rangle \langle ASS \rangle \langle id, "x" \rangle \langle PLUS \rangle \langle id, "z" \rangle$	

Parse trees: Traces of Derivations

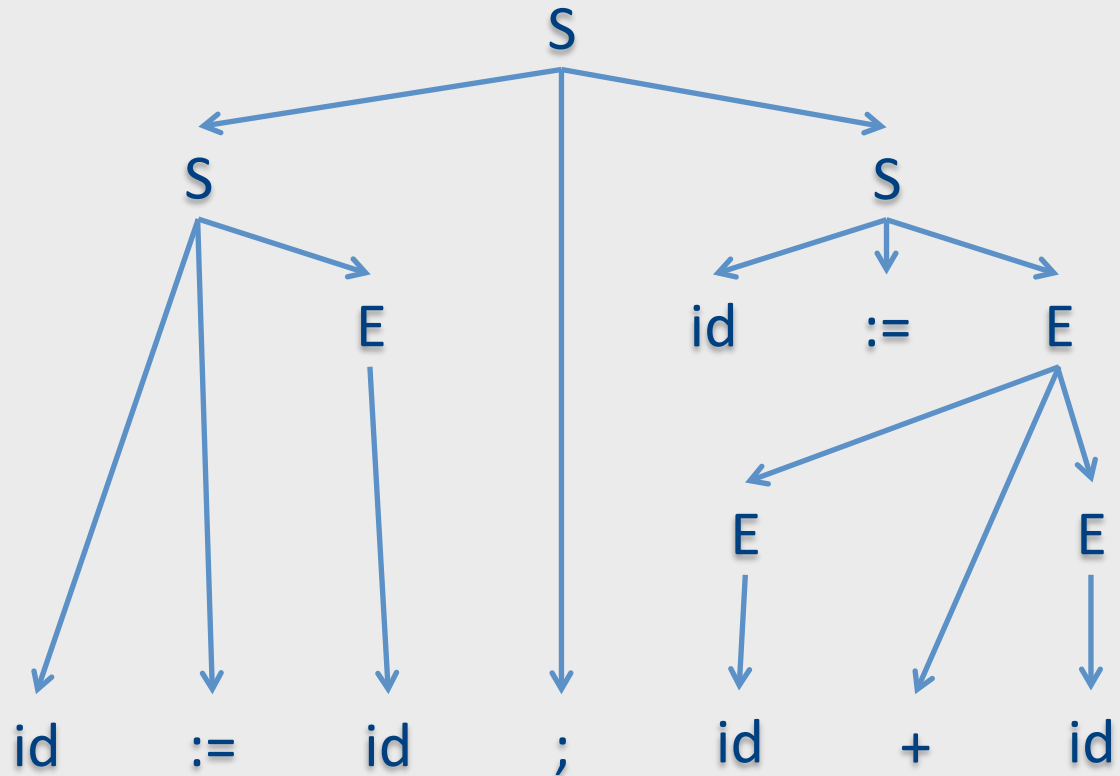
- Tree nodes are symbols, children ordered left-to-right
- Each internal node is non-terminal and its children correspond to one of its productions



- Root is start non-terminal
- Leaves are tokens
- *Yield* of parse tree: left-to-right walk over leaves

Parse Tree

```
S  
S ; S  
id := E; S  
id := id; S  
id := id; id := E  
id := id; id := E + E  
id := id; id := E + id  
id := id; id := id + id  
x := z ; y := x + z
```



Language of a CFG

- A sentence ω is in $L(G)$ (valid program) if
 - There exists a corresponding derivation
 - There exists a corresponding parse tree

Questions

- How did we know which rule to apply on every step?
- Would we always get the same result?
- Does it matter?

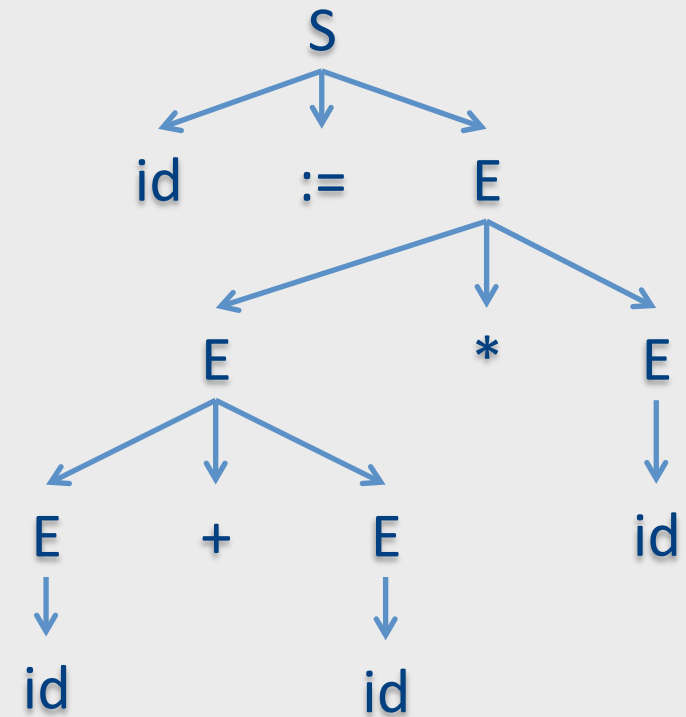
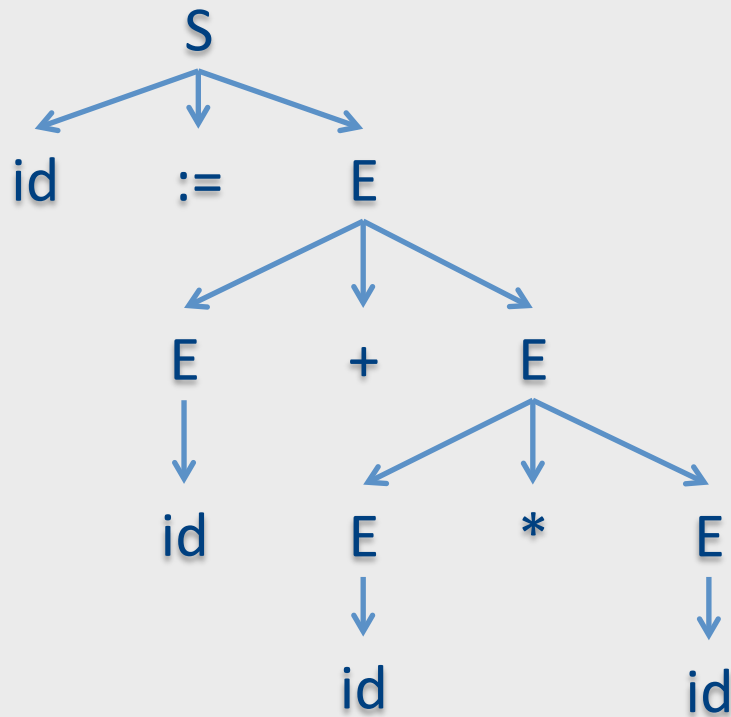
Ambiguity

$x := y+z*w$

$S \rightarrow S ; S$

$S \rightarrow id := E \mid \dots$

$E \rightarrow id \mid E + E \mid E * E \mid \dots$



Leftmost/rightmost Derivation

- Leftmost derivation
 - always expand leftmost non-terminal
- Rightmost derivation
 - Always expand rightmost non-terminal
- Allows us to describe derivation by listing the sequence of rules
 - always know what a rule is applied to

Leftmost Derivation

x := z;
y := x + z

$S \rightarrow S;S$

$S \rightarrow id := E$

$E \rightarrow id \mid E + E \mid E * E \mid (E)$

S

S ; S

id := **E** ; S

id := id; **S**

id := id; id := **E**

id := id; id := **E + E**

id := id; id := id + **E**

id := id; id := id + id

x := z ; y := x + z

$S \rightarrow S;S$

$S \rightarrow id := E$

$E \rightarrow id$

$S \rightarrow id := E$

$E \rightarrow E + E$

$E \rightarrow id$

$E \rightarrow id$

Rightmost Derivation

$\langle \text{id}, "x" \rangle$ ASS $\langle \text{id}, "z" \rangle$;
 $\langle \text{id}, "y" \rangle$ ASS $\langle \text{id}, "x" \rangle$ PLUS
 $\langle \text{id}, "z" \rangle$

$S \rightarrow S;S$
 $S \rightarrow \text{id} := E \mid \dots$
 $E \rightarrow \text{id} \mid E + E \mid E * E \mid \dots$

S		$S \rightarrow S;S$
$S ; S$		$S \rightarrow \text{id} := E$
$S ; \text{id} := E$		$E \rightarrow E + E$
$S ; \text{id} := E + E$		$E \rightarrow \text{id}$
$S ; \text{id} := E + \text{id}$		$E \rightarrow \text{id}$
$S ; \text{id} := \text{id} + \text{id}$		$S \rightarrow \text{id} := E$
$\text{id} := E ; \text{id} := \text{id} + \text{id}$		$E \rightarrow \text{id}$
$\text{id} := \text{id} ; \text{id} := \text{id} + \text{id}$		
$\langle \text{id}, "x" \rangle$ ASS $\langle \text{id}, "z" \rangle$; $\langle \text{id}, "y" \rangle$ ASS $\langle \text{id}, "x" \rangle$ PLUS $\langle \text{id}, "z" \rangle$		

Sometimes there are two parse trees

Arithmetic expressions:

$E \rightarrow \text{id}$

$E \rightarrow \text{num}$

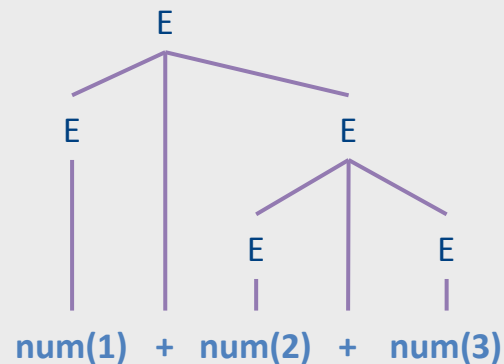
$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow (E)$

1 + 2 + 3

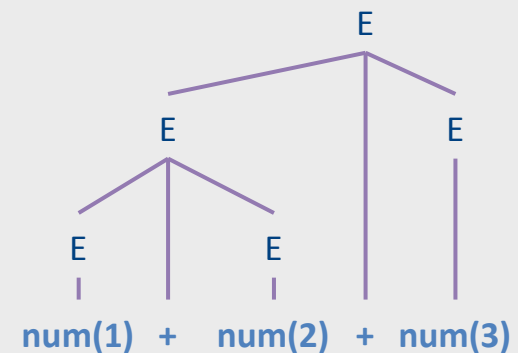
"1 + (2 + 3)"



Leftmost derivation

E
E + E
num + E
num + E + E
num + num + E
num + num + num

"(1 + 2) + 3"



Rightmost derivation

E
E + E
E + num
E + E + num
E + num + num
num + num + num

Is ambiguity a problem?

Arithmetic expressions:

$E \rightarrow \text{id}$

$E \rightarrow \text{num}$

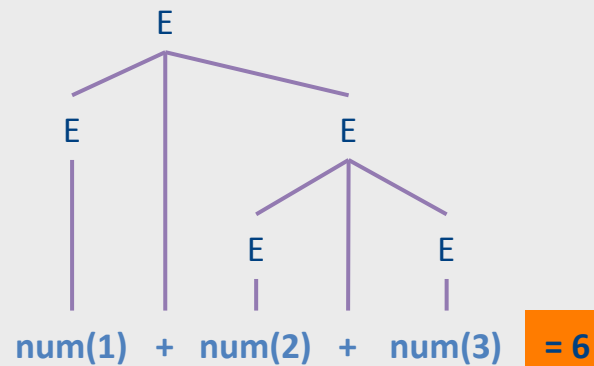
$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow (E)$

$1 + 2 + 3$

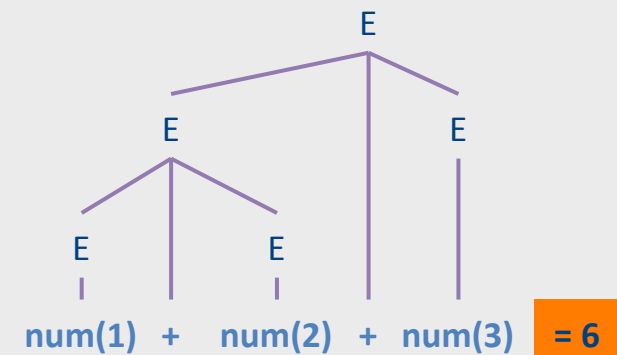
$"1 + (2 + 3)"$



Leftmost derivation

E
E + E
num + E
num + E + E
num + num + E
num + num + num

$"(1 + 2) + 3"$



Rightmost derivation

E
E + E
E + num
E + E + num
E + num + num
num + num + num

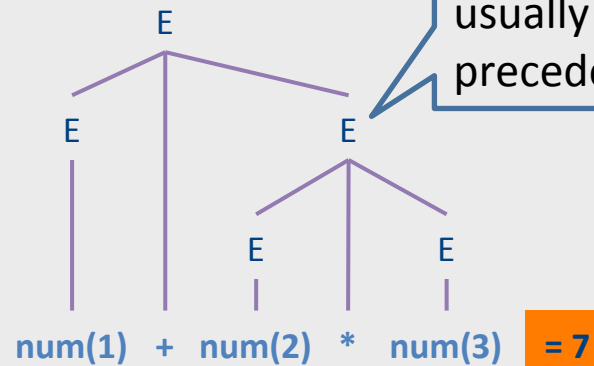
Problematic ambiguity example

Arithmetic expressions:

- $E \rightarrow \text{id}$
- $E \rightarrow \text{num}$
- $E \rightarrow E + E$
- $E \rightarrow E * E$
- $E \rightarrow (E)$

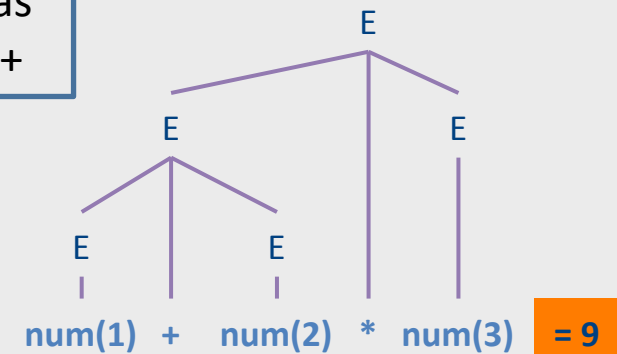
$1 + 2 * 3$

"1 + (2 * 3)"



This is what we usually want: * has precedence over +

"(1 + 2) * 3"



Leftmost derivation

E
 $E + E$
 $\text{num} + E$
 $\text{num} + E * E$
 $\text{num} + \text{num} * E$
 $\text{num} + \text{num} * \text{num}$

Rightmost derivation

E
 $E * E$
 $E * \text{num}$
 $E + E * \text{num}$
 $E + \text{num} * \text{num}$
 $\text{num} + \text{num} * \text{num}$

Ambiguous grammars

- A *grammar is ambiguous* if there exists a sentence for which there are
 - Two different leftmost derivations
 - Two different rightmost derivations
 - Two different parse trees
- Property of *grammars*, not *languages*
- Some languages are inherently ambiguous – no unambiguous grammars exist
- No algorithm to detect whether arbitrary grammar is ambiguous

Drawbacks of ambiguous grammars

- Ambiguous semantics
 - $1 + 2 * 3 = 7$ or 9
- Parsing complexity
- May affect other phases
- Solutions
 - Transform grammar into non-ambiguous
 - Handle as part of parsing method
 - Using special form of “precedence”

Transforming ambiguous grammars to non-ambiguous by layering

Ambiguous grammar

$E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow \text{id}$
 $E \rightarrow \text{num}$
 $E \rightarrow (E)$

Unambiguous grammar

$E \rightarrow E + T$
 $E \rightarrow T$

 $T \rightarrow T * F$
 $T \rightarrow F$

 $F \rightarrow \text{id}$
 $F \rightarrow \text{num}$
 $F \rightarrow (E)$

Layer 1

Layer 2

Layer 3

Let's derive $1 + 2 * 3$

Each layer takes care of one way of composing substrings to form a string:
1: by +
2: by *
3: atoms

Transformed grammar: * precedes +

Ambiguous grammar

$E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow \text{id}$
 $E \rightarrow \text{num}$
 $E \rightarrow (E)$

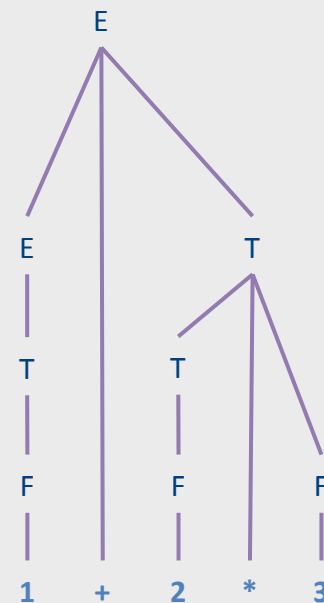
Unambiguous grammar

$E \rightarrow E + T$
 $E \rightarrow T$
 $T \rightarrow T * F$
 $T \rightarrow F$
 $F \rightarrow \text{id}$
 $F \rightarrow \text{num}$
 $F \rightarrow (E)$

Derivation

E
 $\Rightarrow E + T$
 $\Rightarrow T + T$
 $\Rightarrow F + T$
 $\Rightarrow 1 + T$
 $\Rightarrow 1 + T * F$
 $\Rightarrow 1 + F * F$
 $\Rightarrow 1 + 2 * F$
 $\Rightarrow 1 + 2 * 3$

Parse tree



Transformed grammar: + precedes *

Ambiguous grammar

$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow \text{id}$

$E \rightarrow \text{num}$

$E \rightarrow (E)$

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

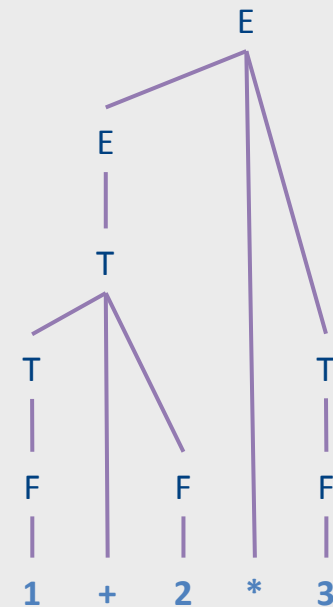
$F \rightarrow \text{num}$

$F \rightarrow (E)$

Derivation

E
 $\Rightarrow E * T$
 $\Rightarrow T * T$
 $\Rightarrow T + F * T$
 $\Rightarrow F + F * T$
 $\Rightarrow 1 + F * T$
 $\Rightarrow 1 + 2 * T$
 $\Rightarrow 1 + 2 * F$
 $\Rightarrow 1 + 2 * 3$

Parse tree



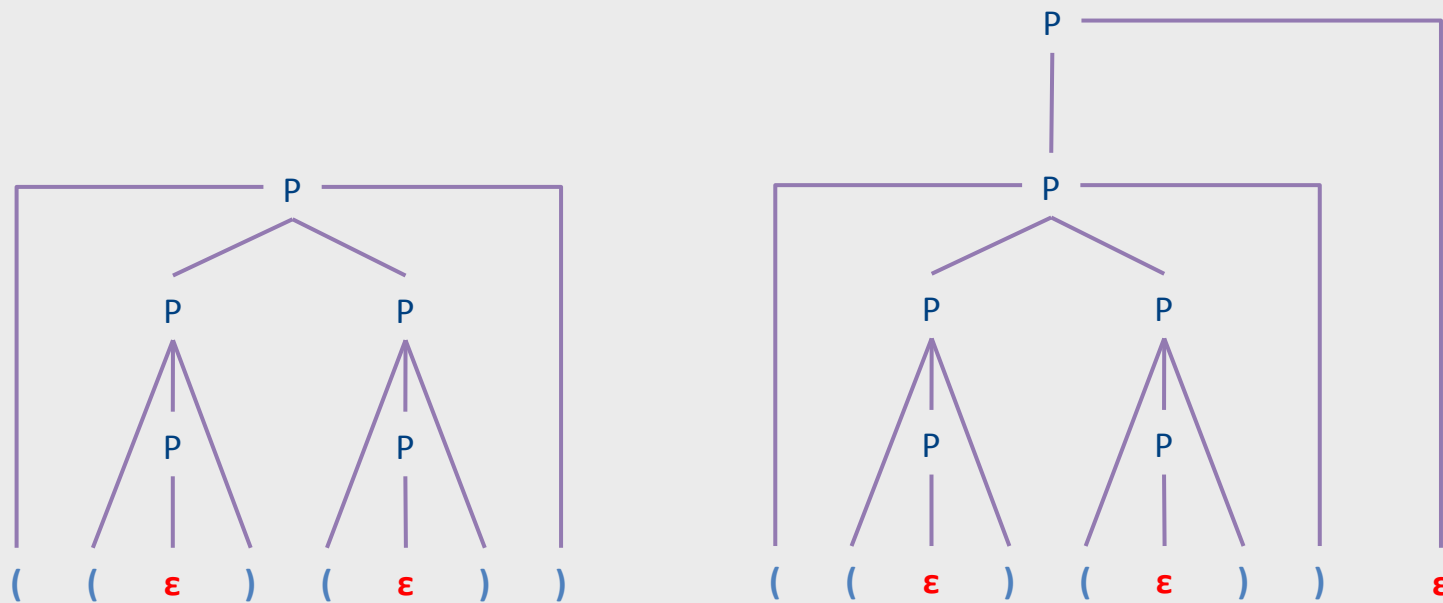
Another example for layering

Ambiguous grammar

$P \rightarrow \epsilon$

$| P P$

$| (P)$



Another example for layering

Ambiguous grammar

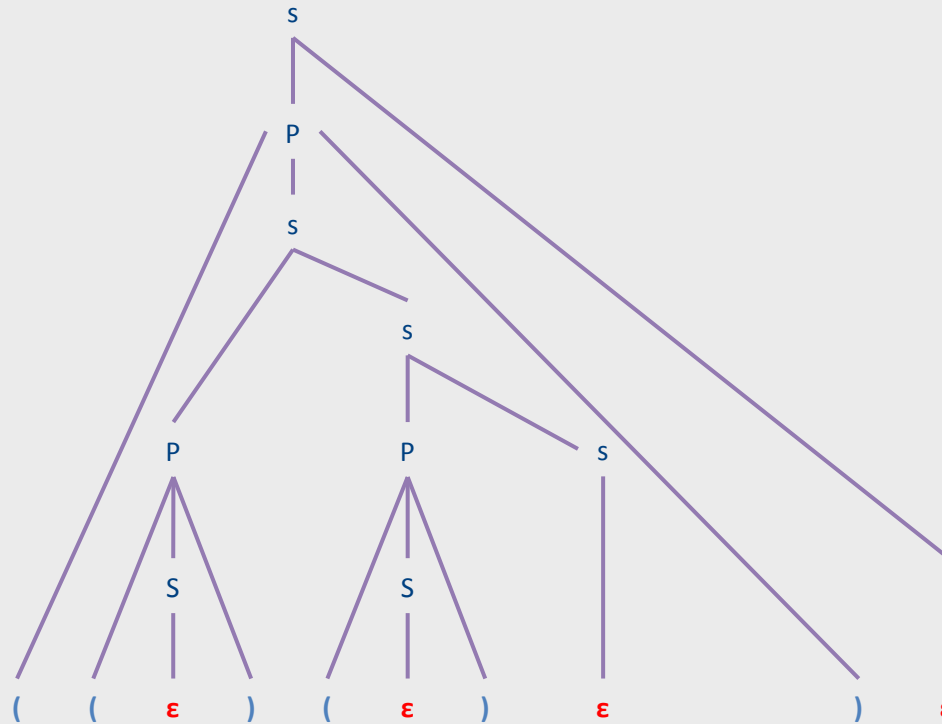
$P \rightarrow \epsilon$
| PP
| (P)

Unambiguous grammar

$S \rightarrow PS$
| ϵ
 $P \rightarrow (S)$

Takes care of
"concatenation"

Takes care of nesting



“dangling-else” example

Ambiguous grammar

$S \rightarrow \text{if } E \text{ then } S$
 $S \mid \text{if } E \text{ then } S \text{ else } S$
 $\mid \text{other}$

This is what we usually want: match **else** to closest unmatched **then**

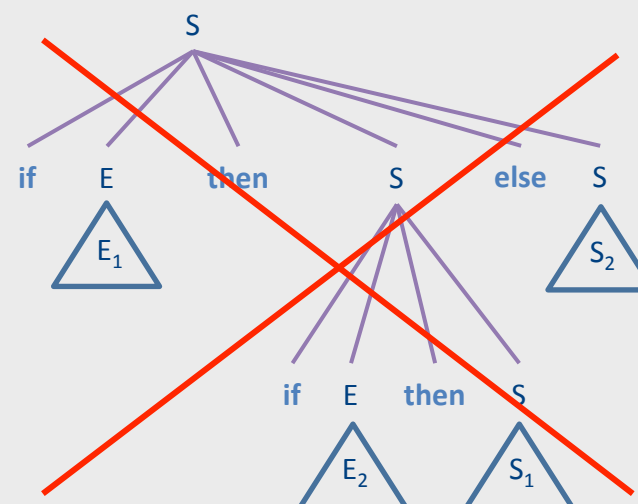
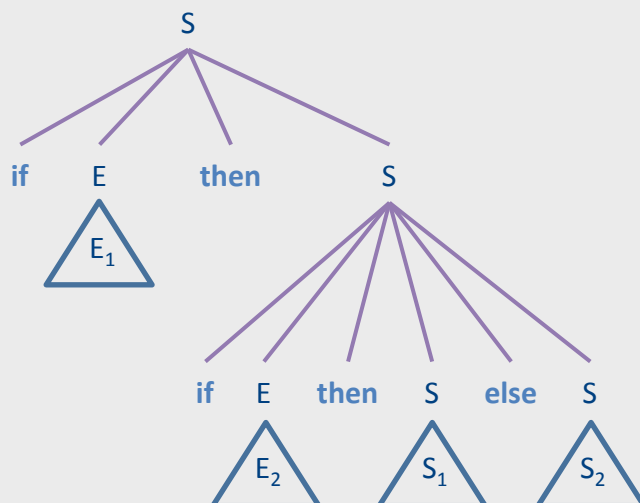
Unambiguous grammar

?

if E_1 then if E_2 then S_1 else S_2

if E_1 then (if E_2 then S_1 else S_2)

if E_1 then (if E_2 then S_1) else S_2



Broad kinds of parsers

- **Parsers for arbitrary grammars**
 - Earley's method, CYK method $O(n^3)$
 - Usually, not used in practice (though might change)
- **Top-Down parsers**
 - Construct parse tree in a top-down matter
 - Find the leftmost derivation
 - Predictive: for every non-terminal and k-tokens predict the next production LL(k)
 - Preorder tree traversal
- **Bottom-Up parsers**
 - Construct parse tree in a bottom-up manner
 - Find the rightmost derivation in a reverse order
 - For every potential right hand side and k-tokens decide when a production is found LR(k)
 - Postorder tree traversal

Broad kinds of parsers

- Parsers for arbitrary grammars
 - Earley's method, CYK method $O(n^3)$
 - Usually, not used in practice (though might change)
- **Top-Down** parsers
 - Construct parse tree in a top-down manner
 - Find the leftmost derivation
- **Bottom-Up** parsers
 - Construct parse tree in a bottom-up manner
 - Find the rightmost derivation in a reverse order

Intuition: Top-down vs. bottom-up

- Top-down parsing
 - Begin with start symbol
 - Guess the productions
 - Check if parse tree yields user's program
- Bottom-up parsing
 - Begin with the user's program
 - Guess parse subtrees
 - Check if root is the start symbol

Intuition: Top-Down Parsing

- Begin with start symbol
- Guess the productions
- Check if parse tree yields user's program

Intuition: Top-Down parsing

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

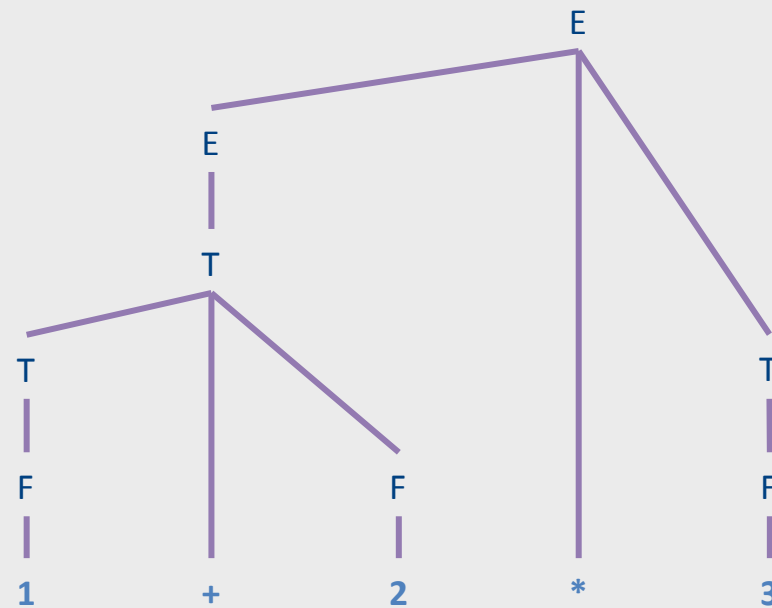
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$



Intuition: Top-Down parsing

**Unambiguous
grammar**

$E \rightarrow E * T$

$E \rightarrow T$

$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$

We need this
rule to get the *

E

1

+

2

*

3



Intuition: Top-down parsing

**Unambiguous
grammar**

$E \rightarrow E * T$

$E \rightarrow T$

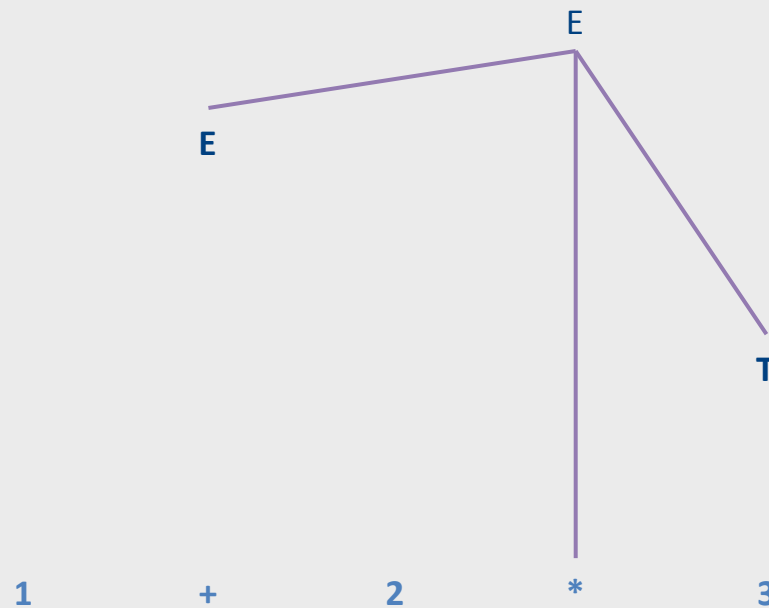
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$



Intuition: Top-down parsing

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

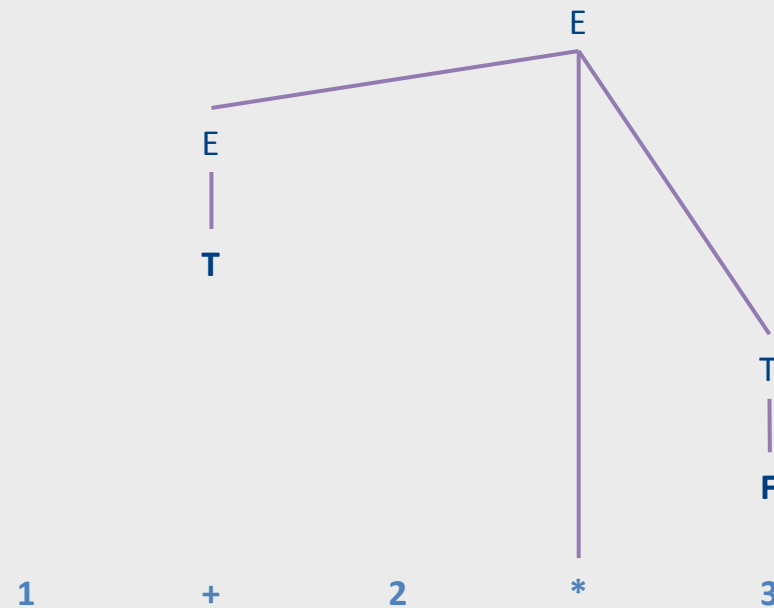
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$



Intuition: Top-down parsing

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

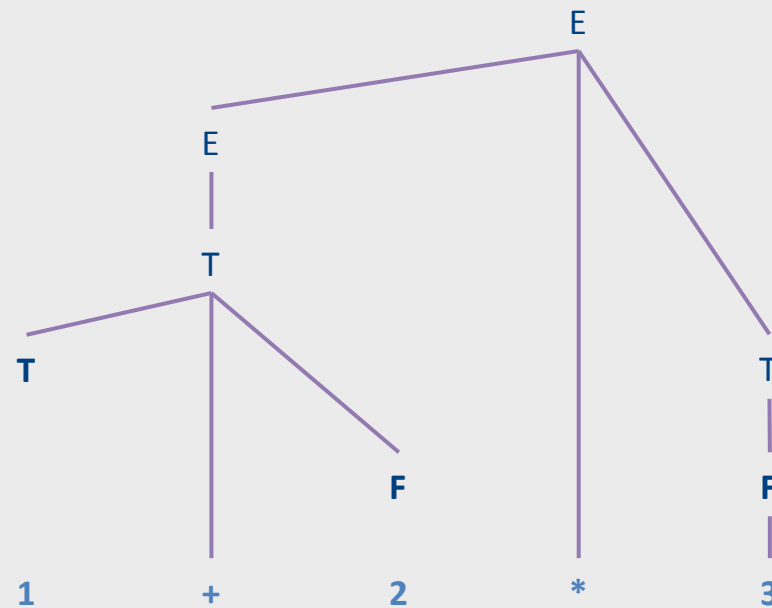
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$



Intuition: Top-Down parsing

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

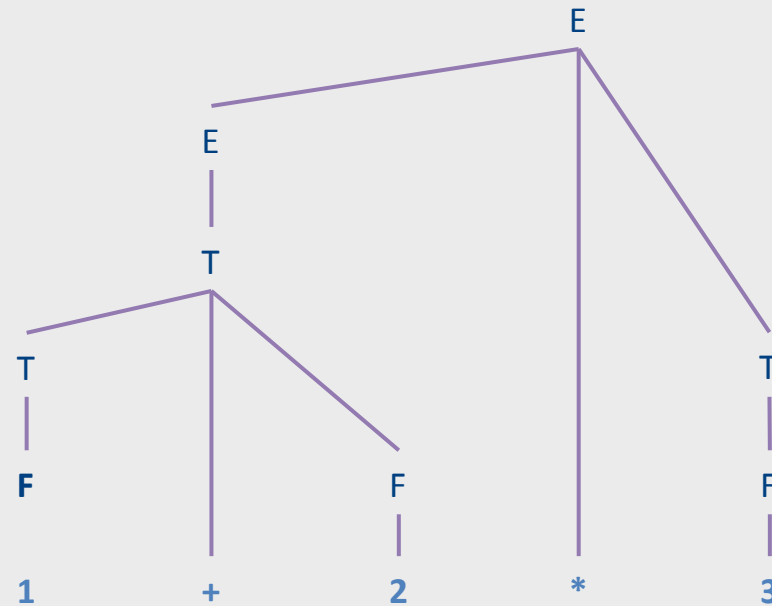
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$



Intuition: Top-Down parsing

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

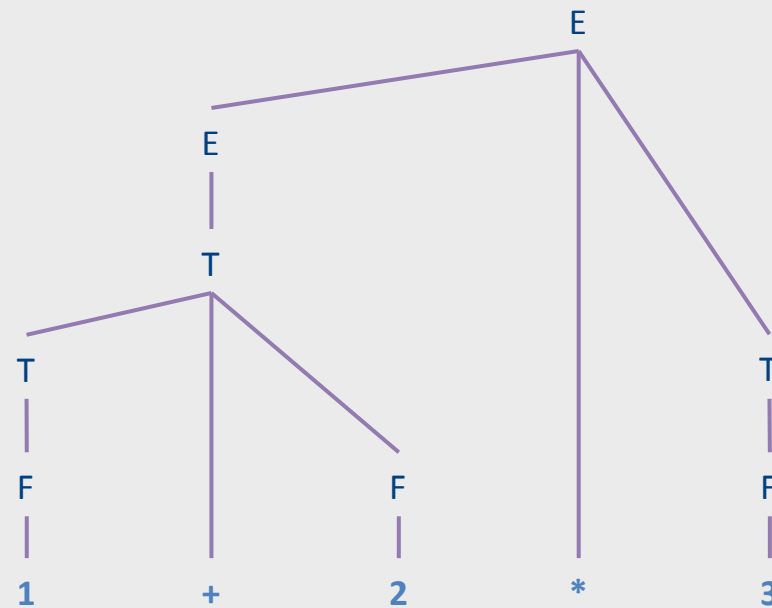
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$



Intuition: Bottom-Up Parsing

- Begin with the user's program
- Guess parse (sub)trees
- Check if root is the start symbol

Bottom-up parsing

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$

1

+

2

*

3

Bottom-up parsing

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$

1

+

2

*

3

F
|

Top-down parsing

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$

1 + F
 |
 2 * F
 |
 3

Top-down parsing

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

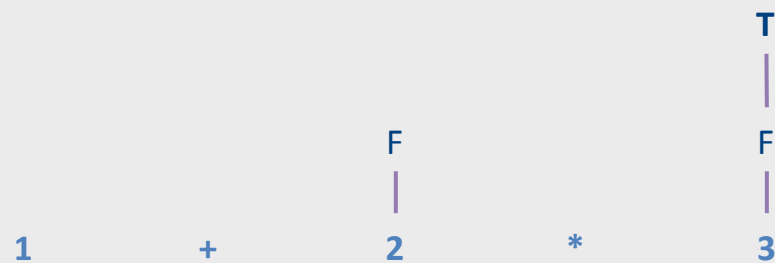
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$



Top-down parsing

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

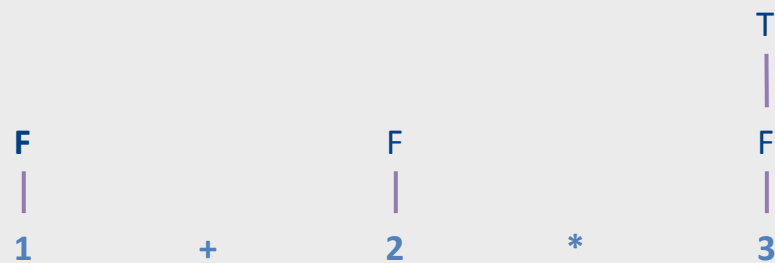
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$



Bottom-up parsing

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

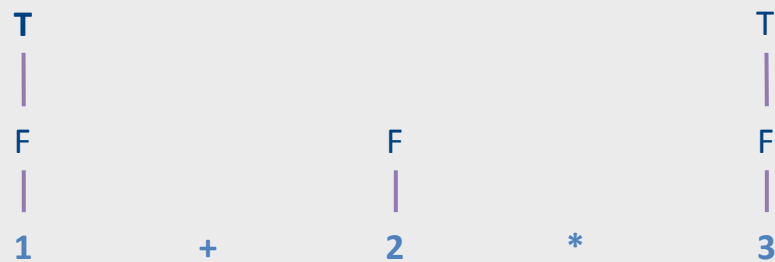
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$



Top-down parsing

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

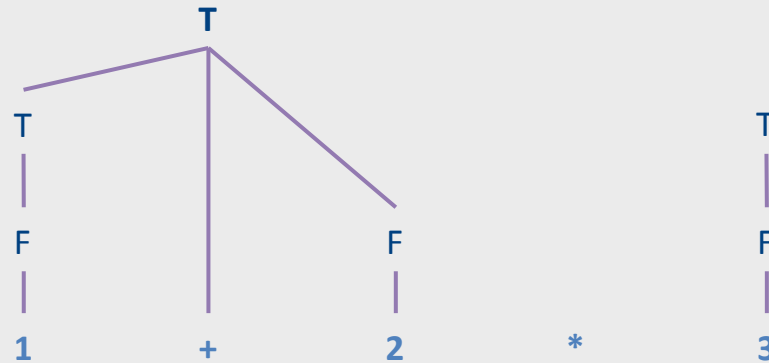
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$



Top-down parsing

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

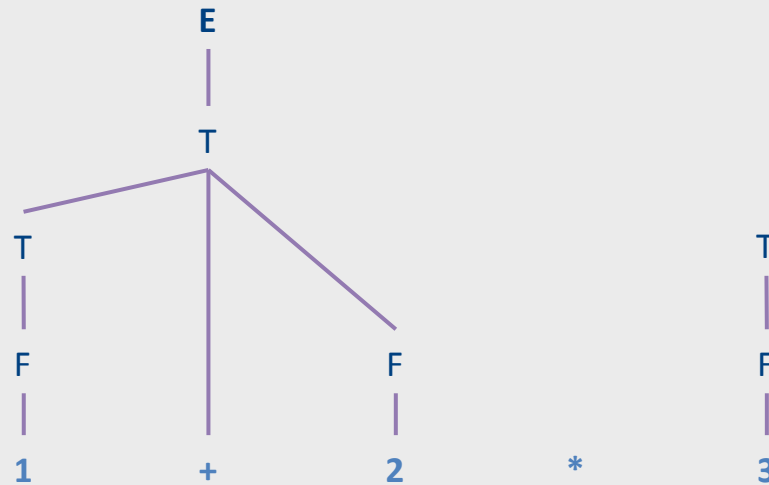
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$



Top-down parsing

Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

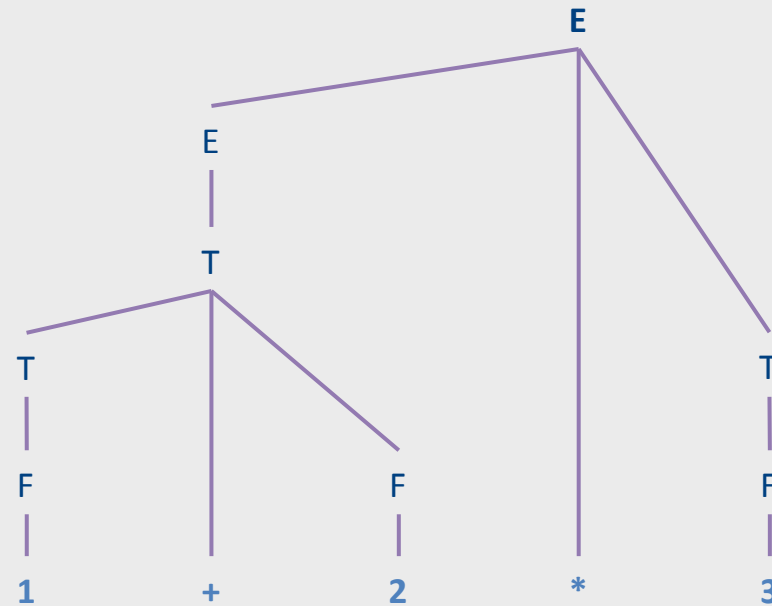
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow (E)$



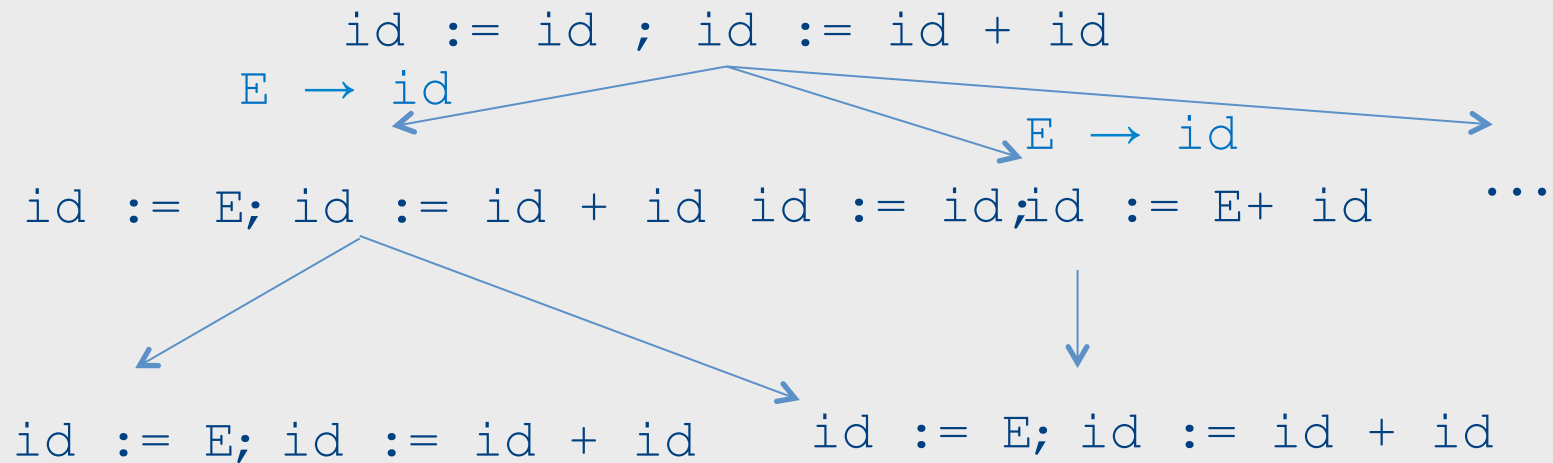
Challenges in top-down parsing

- Top-down parsing begins with virtually no
- information
 - Begins with just the start symbol, which matches *every program*
- How can we know which productions to apply?

“Brute-force” Parsing

`x := z;`
`y := x + z`

`S → S;S`
`S → id := E`
`E → id | E + E | ...`



(not a parse tree... a search for the parse tree by exhaustively applying all rules)

Challenges in top-down parsing

- Top-down parsing begins with virtually no information
 - Begins with just the start symbol, which matches *every program*
- How can we know which productions to apply?
- In general, we can't
 - There are some grammars for which the best we can do is guess and backtrack if we're wrong
- If we have to guess, how do we do it?
 - Parsing as a search algorithm
 - Too expensive in theory (exponential worst-case time) and practice

Predictive parsing

- Given a grammar G and a word w attempt to derive w using G
- Idea
 - Apply production to leftmost nonterminal
 - Pick production rule based on next input token
- General grammar
 - More than one option for choosing the next production based on a token
- Restricted grammars (LL)
 - Know exactly which single rule to apply
 - May require some lookahead to decide

Boolean expressions example

$E \rightarrow \text{LIT} \mid (E \text{ OP } E) \mid \text{not } E$

$\text{LIT} \rightarrow \text{true} \mid \text{false}$

$\text{OP} \rightarrow \text{and} \mid \text{or} \mid \text{xor}$

not (not true or false)

Boolean expressions example

$E \rightarrow LIT \mid (E \text{ OP } E) \mid \text{not } E$

$LIT \rightarrow \text{true} \mid \text{false}$

$OP \rightarrow \text{and} \mid \text{or} \mid \text{xor}$

production to
apply known from
next token

not (not true or false)

$E \Rightarrow$

not $E \Rightarrow$

$\text{not} (E \text{ OP } E) \Rightarrow$

$\text{not} (\text{not } E \text{ OP } E) \Rightarrow$

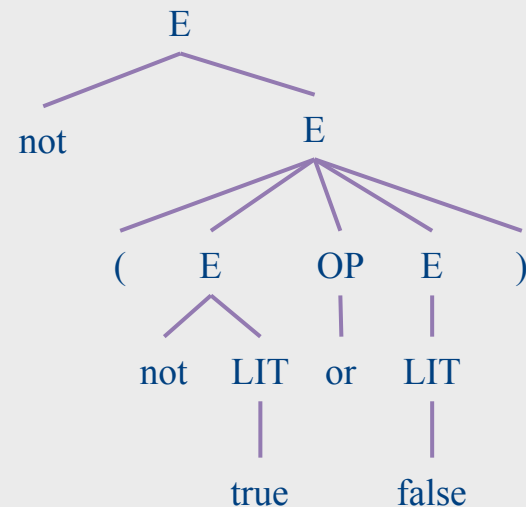
$\text{not} (\text{not } LIT \text{ OP } E) \Rightarrow$

$\text{not} (\text{not } \text{true} \text{ OP } E) \Rightarrow$

$\text{not} (\text{not } \text{true} \text{ or } E) \Rightarrow$

$\text{not} (\text{not } \text{true} \text{ or } LIT) \Rightarrow$

$\text{not} (\text{not } \text{true} \text{ or } \text{false})$



Recursive descent parsing

- Define a function for every nonterminal
- Every function work as follows
 - Find applicable production rule
 - Terminal function checks match with next input token
 - Nonterminal function calls (recursively) other functions
- If there are several applicable productions for a nonterminal, use lookahead

Matching tokens

$E \rightarrow \text{LIT} \mid (\text{E OP E}) \mid \text{not E}$

$\text{LIT} \rightarrow \text{true} \mid \text{false}$

$\text{OP} \rightarrow \text{and} \mid \text{or} \mid \text{xor}$

```
match(token t) {  
    if (current == t)  
        current = next_token()  
    else  
        error  
}
```

- Variable **current** holds the current input token

Functions for nonterminals

$E \rightarrow \text{LIT} \mid (E \text{ OP } E) \mid \text{not } E$

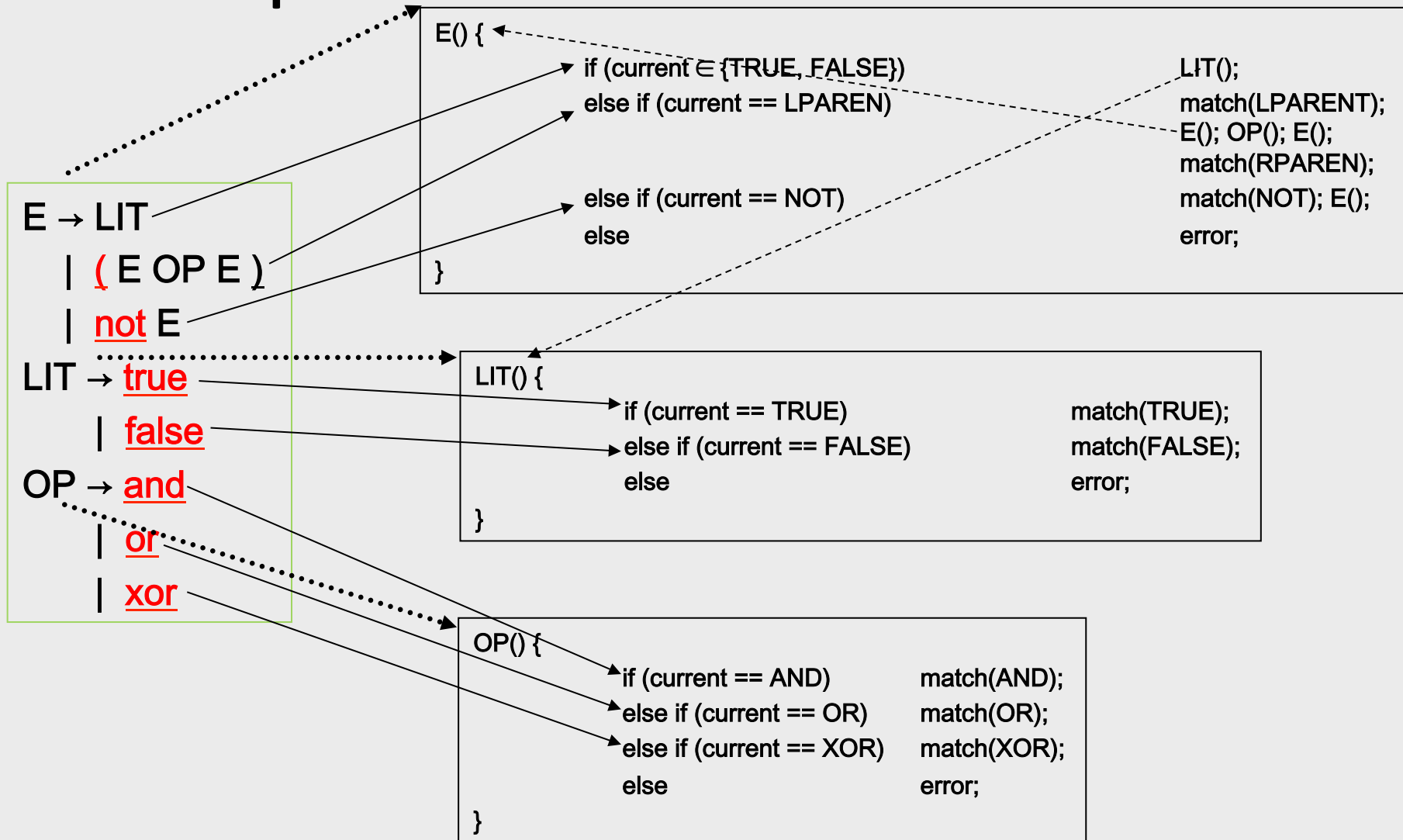
$\text{LIT} \rightarrow \text{true} \mid \text{false}$

$\text{OP} \rightarrow \text{and} \mid \text{or} \mid \text{xor}$

```
E() {
  if (current ∈ {TRUE, FALSE}) // E → LIT
    LIT();
  else if (current == LPAREN) // E → ( E OP E )
    match(LPAREN); E(); OP(); E(); match(RPAREN);
  else if (current == NOT) // E → not E
    match(NOT); E();
  else
    error;
}
```

```
LIT() {
  if (current == TRUE) match(TRUE);
  else if (current == FALSE) match(FALSE);
  else error;
}
```

Implementation via recursion



Recursive descent

```
void A() {
  choose an A-production,  $A \rightarrow X_1X_2\dots X_k$ ;
  for (i=1; i ≤ k; i++) {
    if ( $X_i$  is a nonterminal)
      call procedure  $X_i()$ ;
    elseif ( $X_i == \text{current}$ )
      advance input;
    else
      report error;
  }
}
```

- How do you pick the right A-production?
- Generally – try them all and use backtracking
- In our case – use lookahead

Problem 1: productions with common prefix

term \rightarrow ID | indexed_elem

indexed_elem \rightarrow ID [expr]

- The function for indexed_elem will never be tried...
 - What happens for input of the form ID [expr]

Problem 2: null productions

$S \rightarrow A a b$
 $A \rightarrow a \mid \varepsilon$

```
int S() {  
    return A() && match(token('a')) && match(token('b'));  
}  
int A() {  
    return match(token('a')) || 1;  
}
```

- What happens for input “ab”?
- What happens if you flip order of alternatives and try “aab”?

Problem 3: left recursion

$E \rightarrow E - \text{term} \mid \text{term}$

```
int E() {  
    return E() && match(token('-')) && term();  
}
```

- What happens with this procedure?
- **Recursive descent parsers cannot handle left-recursive grammars**

FIRST sets

- For every production rule $A \rightarrow \alpha$
 - $\text{FIRST}(\alpha)$ = all terminals that α can start with
 - Every token that can appear as first in α under some derivation for α
- In our Boolean expressions example
 - $\text{FIRST}(\text{LIT}) = \{ \text{true}, \text{false} \}$
 - $\text{FIRST}((E \text{ OP } E)) = \{ '(' \}$
 - $\text{FIRST}(\text{not } E) = \{ \text{not} \}$
- No intersection between FIRST sets \Rightarrow can always pick a single rule
- If the FIRST sets intersect, may need longer lookahead
 - $\text{LL}(k)$ = class of grammars in which production rule can be determined using a lookahead of k tokens
 - $\text{LL}(1)$ is an important and useful class

Computing FIRST sets

- Assume no null productions $A \rightarrow \varepsilon$
 1. Initially, for all nonterminals A , set $\text{FIRST}(A) = \{ t \mid A \rightarrow t\omega \text{ for some } \omega \}$
 2. Repeat the following until no changes occur:
for each nonterminal A
for each production $A \rightarrow B\omega$
set $\text{FIRST}(A) = \text{FIRST}(A) \cup \text{FIRST}(B)$
- This is known as fixed-point computation

FIRST sets computation example

STMT \rightarrow if EXPR then STMT
| while EXPR do STMT
| EXPR ;

EXPR \rightarrow TERM \rightarrow id
| zero? TERM
| not EXPR
| ++ id
| -- id

TERM \rightarrow id
| constant

STMT	EXPR	TERM

1. Initialization

STMT \rightarrow if EXPR then STMT
| while EXPR do STMT
| EXPR ;
EXPR \rightarrow TERM \rightarrow id
| zero? TERM
| not EXPR
| ++ id
| -- id
TERM \rightarrow id
| constant

STMT	EXPR	TERM
if while	zero? Not ++ --	id constant

2. Iterate 1

STMT → if EXPR then STMT
| while EXPR do STMT
| **EXPR** ;
EXPR → **TERM** -> id
| zero? **TERM**
| not **EXPR**
| ++ id
| -- id
TERM → id
| constant

STMT	EXPR	TERM
if while	zero? Not ++ --	id constant
zero? Not ++ --		

2. Iterate 2

STMT → if EXPR then STMT
| while EXPR do STMT
| **EXPR** ;

EXPR → **TERM** -> id
| zero? **TERM**
| not **EXPR**
| ++ id
| -- id

TERM → id
| constant

STMT	EXPR	TERM
if while	zero? Not ++ --	id constant
zero? Not ++ --	id constant	

2. Iterate 3 – fixed-point

STMT → if EXPR then STMT
 | while EXPR do STMT
 | **EXPR** ;
EXPR → TERM -> id
 | zero? TERM
 | not EXPR
 | ++ id
 | -- id
TERM → id
 | constant

STMT	EXPR	TERM
if while	zero? Not ++ --	id constant
zero? Not ++ --	id constant	
id constant		

FOLLOW sets

- What do we do with nullable (ϵ) productions?
 - $A \rightarrow B C D \quad B \rightarrow \epsilon \quad C \rightarrow \epsilon$
 - Use what comes afterwards to predict the right production
- For every production rule $A \rightarrow \alpha$
 - $\text{FOLLOW}(A)$ = set of tokens that can immediately follow A
- Can predict the alternative A_k for a non-terminal N when the lookahead token is in the set
 - $\text{FIRST}(A_k) \cup (\text{if } A_k \text{ is nullable then } \text{FOLLOW}(N))$

LL(k) grammars

- A grammar is in the class LL(K) when it can be derived via:
 - Top-down derivation
 - Scanning the input from left to right (L)
 - Producing the leftmost derivation (L)
 - With lookahead of k tokens (k)
 - For every two productions $A \rightarrow \alpha$ and $A \rightarrow \beta$ we have $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \{\}$
and $\text{FIRST}(A) \cap \text{FOLLOW}(A) = \{\}$
- A language is said to be LL(k) when it has an LL(k) grammar

Back to problem 1

term \rightarrow ID | indexed_elem
indexed_elem \rightarrow ID [expr]

- FIRST(term) = { ID }
- FIRST(indexed_elem) = { ID }
- FIRST/FIRST conflict

Solution: left factoring

- Rewrite the grammar to be in LL(1)

term \rightarrow ID | indexed_elem
indexed_elem \rightarrow ID [expr]



term \rightarrow ID after_ID
After_ID \rightarrow [expr] | ϵ

Intuition: just like factoring $x*y + x*z$ into $x*(y+z)$

Left factoring – another example

$S \rightarrow$ if E then S else S
| if E then S
| T



$S \rightarrow$ if E then S S'
| T
 $S' \rightarrow$ else S | ϵ

Back to problem 2

$S \rightarrow A a b$

$A \rightarrow a \mid \varepsilon$

- $\text{FIRST}(S) = \{ a \}$ $\text{FOLLOW}(S) = \{ \}$
- $\text{FIRST}(A) = \{ a \varepsilon \}$ $\text{FOLLOW}(A) = \{ a \}$
- **FIRST/FOLLOW conflict**

Solution: substitution

$$S \rightarrow A a b$$
$$A \rightarrow a \mid \varepsilon$$


Substitute A in S

$$S \rightarrow a a b \mid a b$$


Left factoring

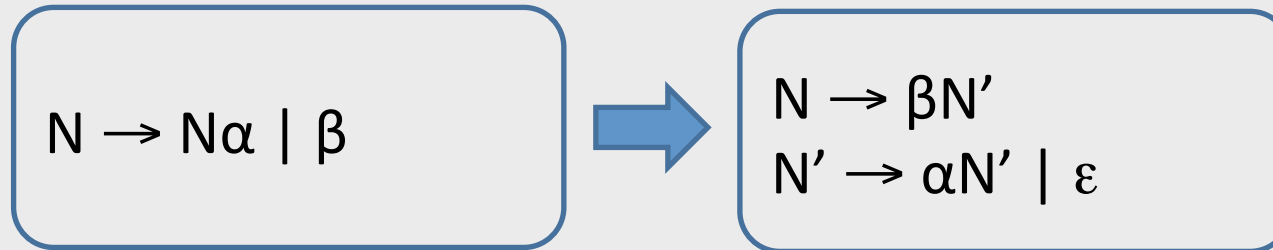
$$S \rightarrow a \text{ after_}A$$
$$\text{after_}A \rightarrow a b \mid b$$

Back to problem 3

$E \rightarrow E - \text{term} \mid \text{term}$

- Left recursion cannot be handled with a bounded lookahead
- What can we do?

Left recursion removal



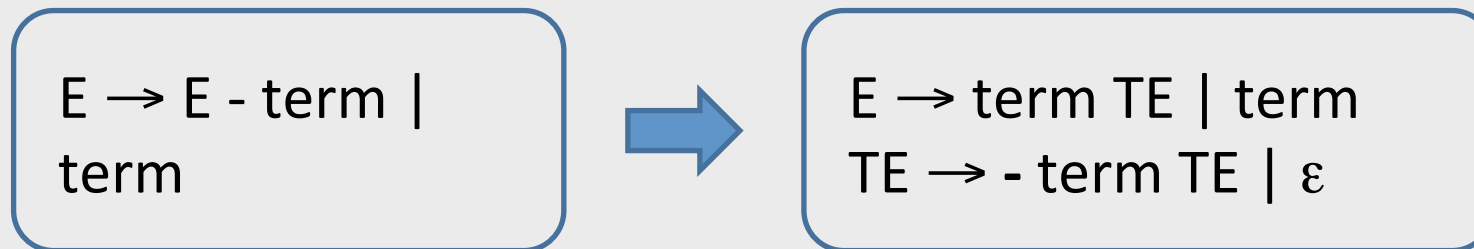
G_1

G_2

- $L(G_1) = \beta, \beta\alpha, \beta\alpha\alpha, \beta\alpha\alpha\alpha, \dots$
- $L(G_2) = \text{same}$

Can be done algorithmically.
Problem: grammar becomes mangled beyond recognition

- For our 3rd example:



LL(k) Parsers

- Recursive Descent
 - Manual construction
 - Uses recursion
- Wanted
 - A parser that can be generated automatically
 - Does not use recursion

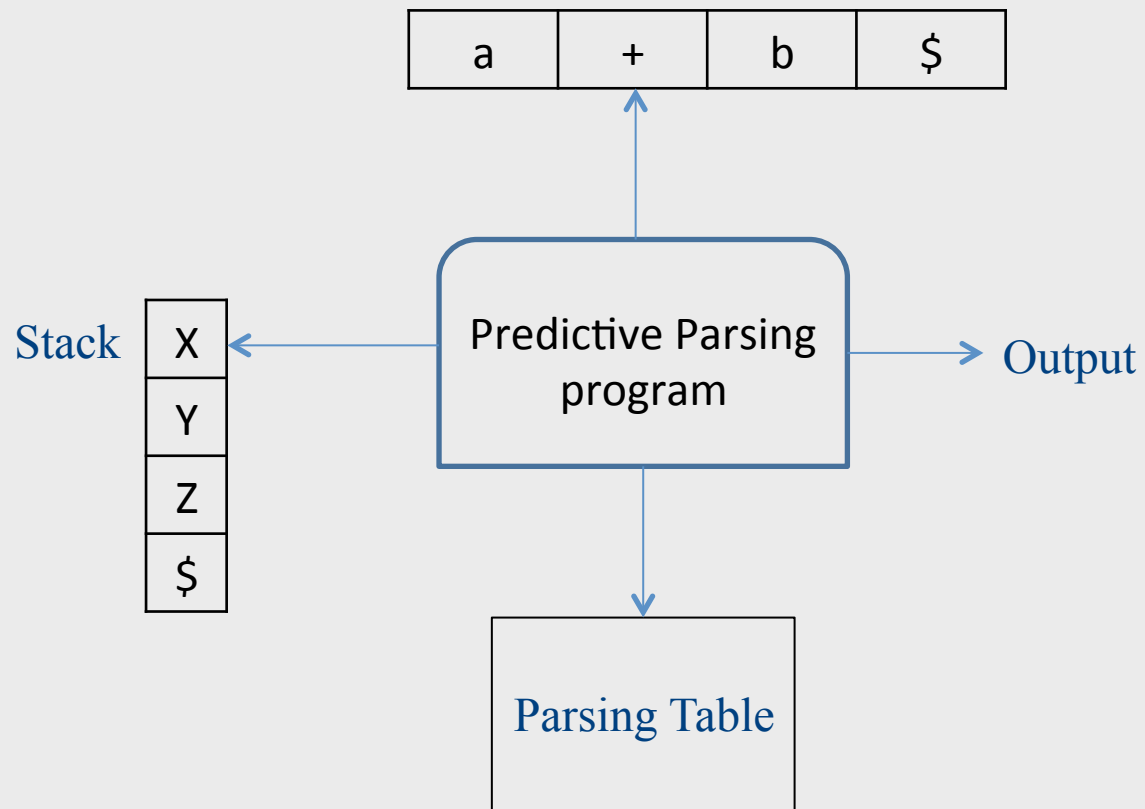
LL(k) parsing via pushdown automata

- Pushdown automaton uses
 - Prediction stack
 - Input stream
 - Transition table
 - nonterminals x tokens \rightarrow production alternative
 - Entry indexed by nonterminal N and token t contains the alternative of N that must be predicated when current input starts with t

LL(k) parsing via pushdown automata

- Two possible moves
 - **Prediction**
 - When top of stack is nonterminal N , pop N , lookup $\text{table}[N,t]$. If $\text{table}[N,t]$ is not empty, push $\text{table}[N,t]$ on prediction stack, otherwise – syntax error
 - **Match**
 - When top of prediction stack is a terminal T , must be equal to next input token t . If $(t == T)$, pop T and consume t . If $(t \neq T)$ syntax error
- Parsing terminates when prediction stack is empty
 - If input is empty at that point, success. Otherwise, syntax error

Model of non-recursive predictive parser



Example transition table

- (1) $E \rightarrow LIT$
- (2) $E \rightarrow (E OP E)$
- (3) $E \rightarrow \text{not } E$
- (4) $LIT \rightarrow \text{true}$
- (5) $LIT \rightarrow \text{false}$
- (6) $OP \rightarrow \text{and}$
- (7) $OP \rightarrow \text{or}$
- (8) $OP \rightarrow \text{xor}$

Which rule should be used

Nonterminals

Input tokens

	()	not	true	false	and	or	xor	\$
E	2		3	1	1				
LIT				4	5				
OP						6	7	8	

Running parser example

aacbb\$

$A \rightarrow aAb \mid c$

Input suffix	Stack content	Move
aacbb\$	A\$	predict(A,a) = $A \rightarrow aAb$
aacbb\$	aAb\$	match(a,a)
acbb\$	Ab\$	predict(A,a) = $A \rightarrow aAb$
acbb\$	aAbb\$	match(a,a)
cbb\$	Abb\$	predict(A,c) = $A \rightarrow c$
cbb\$	cbb\$	match(c,c)
bb\$	bb\$	match(b,b)
b\$	b\$	match(b,b)
\$	\$	match(\$,\$) – success

	a	b	c
A	$A \rightarrow aAb$		$A \rightarrow c$

Erors

Handling Syntax Errors

- Report and locate the error
- Diagnose the error
- Correct the error
- Recover from the error in order to discover more errors
 - without reporting too many “strange” errors

Error Diagnosis

- Line number
 - may be far from the actual error
- The current token
- The expected tokens
- Parser configuration

Error Recovery

- Becomes less important in interactive environments
- Example heuristics:
 - Search for a semi-column and ignore the statement
 - Try to “replace” tokens for common errors
 - Refrain from reporting 3 subsequent errors
- Globally optimal solutions
 - For every input w , find a valid program w' with a “minimal-distance” from w

Illegal input example

abcbb\$

$A \rightarrow aAb \mid c$

Input suffix	Stack content	Move
abcbb\$	A\$	predict(A,a) = $A \rightarrow aAb$
abcbb\$	aAb\$	match(a,a)
bcbb\$	Ab\$	predict(A,b) = ERROR

	a	b	c
A	$A \rightarrow aAb$		$A \rightarrow c$

Error handling in LL parsers

c\$

$S \rightarrow a c \mid b S$

Input suffix	Stack content	Move
c\$	S\$	predict(S,c) = ERROR

- Now what?
 - Predict $b S$ anyway “missing token b inserted in line XXX”

	a	b	c
S	$S \rightarrow a c$	$S \rightarrow b S$	

Error handling in LL parsers

c\$

$S \rightarrow a c \mid b S$

Input suffix	Stack content	Move
bc\$	S\$	predict(b,c) = $S \rightarrow bS$
bc\$	bS\$	match(b,b)
c\$	S\$	Looks familiar?

- Result: infinite loop

	a	b	c
S	$S \rightarrow a c$	$S \rightarrow b S$	

Error handling and recovery

- $x = a * (p+q * (-b * (r-s)));$
- Where should we report the error?
- The valid prefix property

The Valid Prefix Property

- For every prefix tokens
 - t_1, t_2, \dots, t_i that the parser identifies as legal:
 - there exists tokens $t_{i+1}, t_{i+2}, \dots, t_n$ such that t_1, t_2, \dots, t_n is a syntactically valid program
- If every token is considered as single character:
 - For every prefix word u that the parser identifies as legal there exists w such that $u.w$ is a valid program

Bring any memories?

Recovery is tricky

- Heuristics for dropping tokens, skipping to semicolon, etc.

Building the Parse Tree

Adding semantic actions

- Can add an action to perform on each production rule
- Can build the parse tree
 - Every function returns an object of type Node
 - Every Node maintains a list of children
 - Function calls can add new children

Building the parse tree

```
Node E() {
    result = new Node();
    result.name = "E";
    if (current ∈ {TRUE, FALSE}) // E → LIT
        result.addChild(LIT());
    else if (current == LPAREN) // E → ( E OP E )
        result.addChild(match(LPAREN));
        result.addChild(E());
        result.addChild(OP());
        result.addChild(E());
        result.addChild(match(RPAREN));
    else if (current == NOT) // E → not E
        result.addChild(match(NOT));
        result.addChild(E());
    else error;
    return result;
}
```

Parser for Fully Parenthesized Exprs

```
static int Parse_Expression(Expression **expr_p) {
    Expression *expr = *expr_p = new_expression() ;
    /* try to parse a digit */
    if (Token.class == DIGIT) {
        expr->type='D';  expr->value=Token.repr -'0';
        get_next_token();
        return 1;      }

    /* try parse parenthesized expression */
    if (Token.class == '(') {
        expr->type='P';  get_next_token();
        if (!Parse_Expression(&expr->left))  Error("missing expression");
        if (!Parse_Operator(&expr->oper))  Error("missing operator");
        if (Token.class != ')') Error("missing )");

        get_next_token();
        return 1; }

    return 0;
}
```

Famous last words

Syntax Analysis (Parsing)

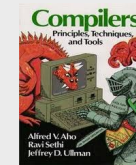
- input
 - Sequence of tokens
- output
 - Abstract Syntax Tree
- Report syntax errors
 - unbalanced parentheses
- Create “symbol-table”
- Create pretty-printed version of the program

Why use context free grammars for defining PL syntax?

- Captures program structure (hierarchy)
- Employ formal theory results
- Automatically create “efficient” parsers

Capabilities and limitations of CFGs

- CFGs naturally express
 - Hierarchical structure
 - A program is a list of classes,
A Class is a list of definition,
A definition is either...
 - Beginning-end type of constraints
 - Balanced parentheses $S \rightarrow (S)S \mid \epsilon$
- Cannot express
 - Correlations between unbounded strings (identifiers)
 - Variables are declared before use: $\omega S \omega$
 - Handled by semantic analysis



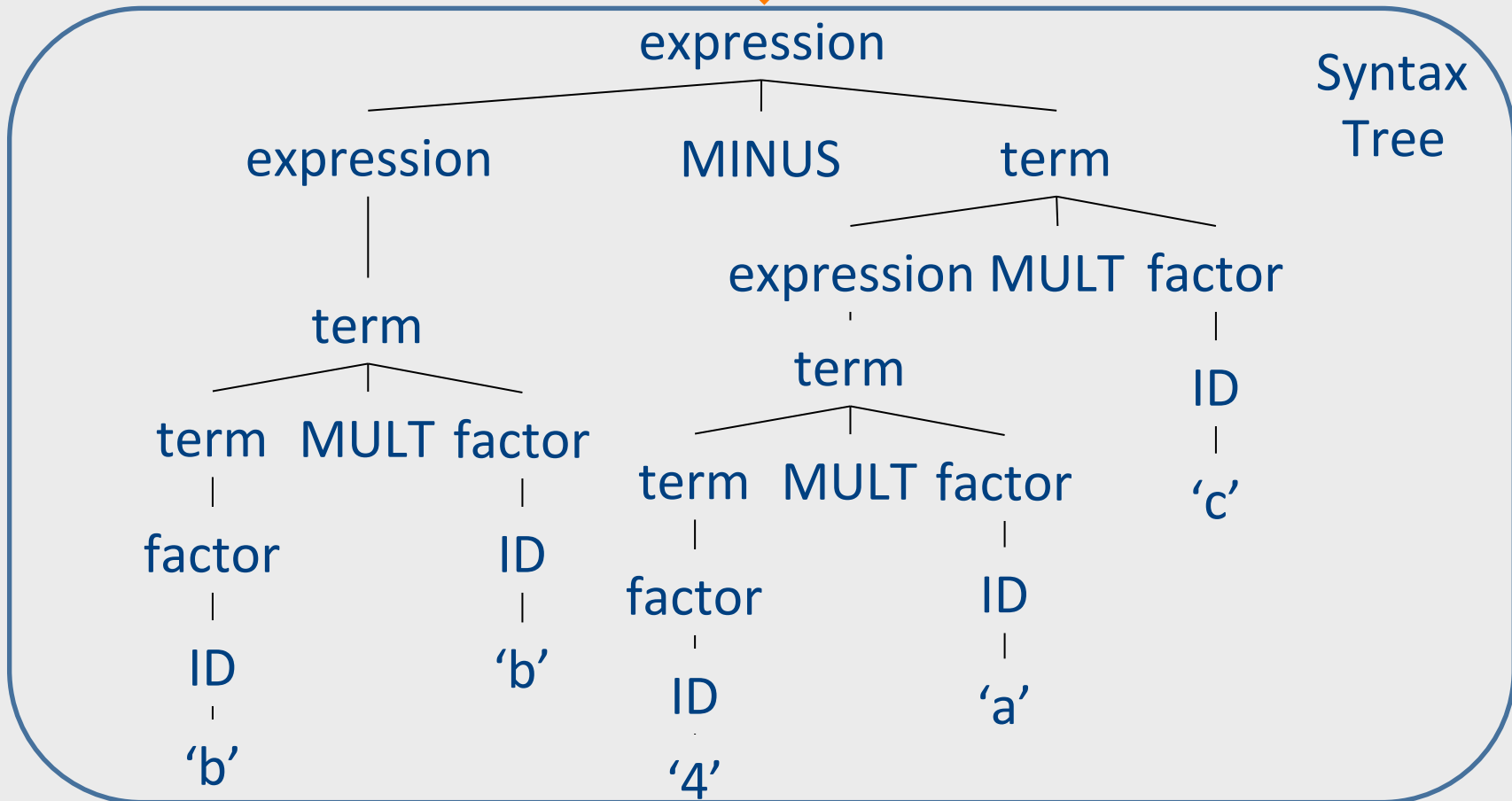
p. 173

See you next week

- Dan David 001

Parsers: from tokens to AST

<ID,"b"> <MULT> <ID,"b"> <MINUS> <INT,4> <MULT> <ID,"a"> <MULT> <ID,"c">



Lexical
Analysis

Syntax
Analysis

Sem.
Analysis

Inter.
Rep.

Code
Gen.

Conventions

- a, b, \dots are input symbols.
 - But sometimes we allow ϵ as a possible value.
- \dots, X, Y, Z are stack symbols.
- \dots, w, x, y, z are strings of input symbols.
- α, β, \dots are strings of stack symbols.



“dangling-else” example

Ambiguous grammar

$S \rightarrow \text{if } E \text{ then } S$
 $S \mid \text{if } E \text{ then } S \text{ else } S$
 $\mid \text{other}$

Unambiguous grammar

$S \rightarrow M \mid U$
 $M \rightarrow \text{if } E \text{ then } M \text{ else } M$
 $\mid \text{other}$
 $U \rightarrow \text{if } E \text{ then } S$
 $\mid \text{if } E \text{ then } M \text{ else } U$

Matched statements

Unmatched statements

$\text{if } E_1 \text{ then if } E_2 \text{ then } S_1 \text{ else } S_2$

$\text{if } E_1 \text{ then (if } E_2 \text{ then } S_1 \text{ else } S_2)$

$\text{if } E_1 \text{ then (if } E_2 \text{ then } S_1) \text{ else } S_2$

