

# Compilation

0368-3133 (Semester A, 2013/14)

## Lecture 4: Syntax Analysis (Top-Down Parsing)

Modern Compiler Design: Chapter 2.2

Noam Rinetzky

# Admin

- Next week: Trubowicz 101 (Law school)
- Mobiles ...

# What is a Compiler?

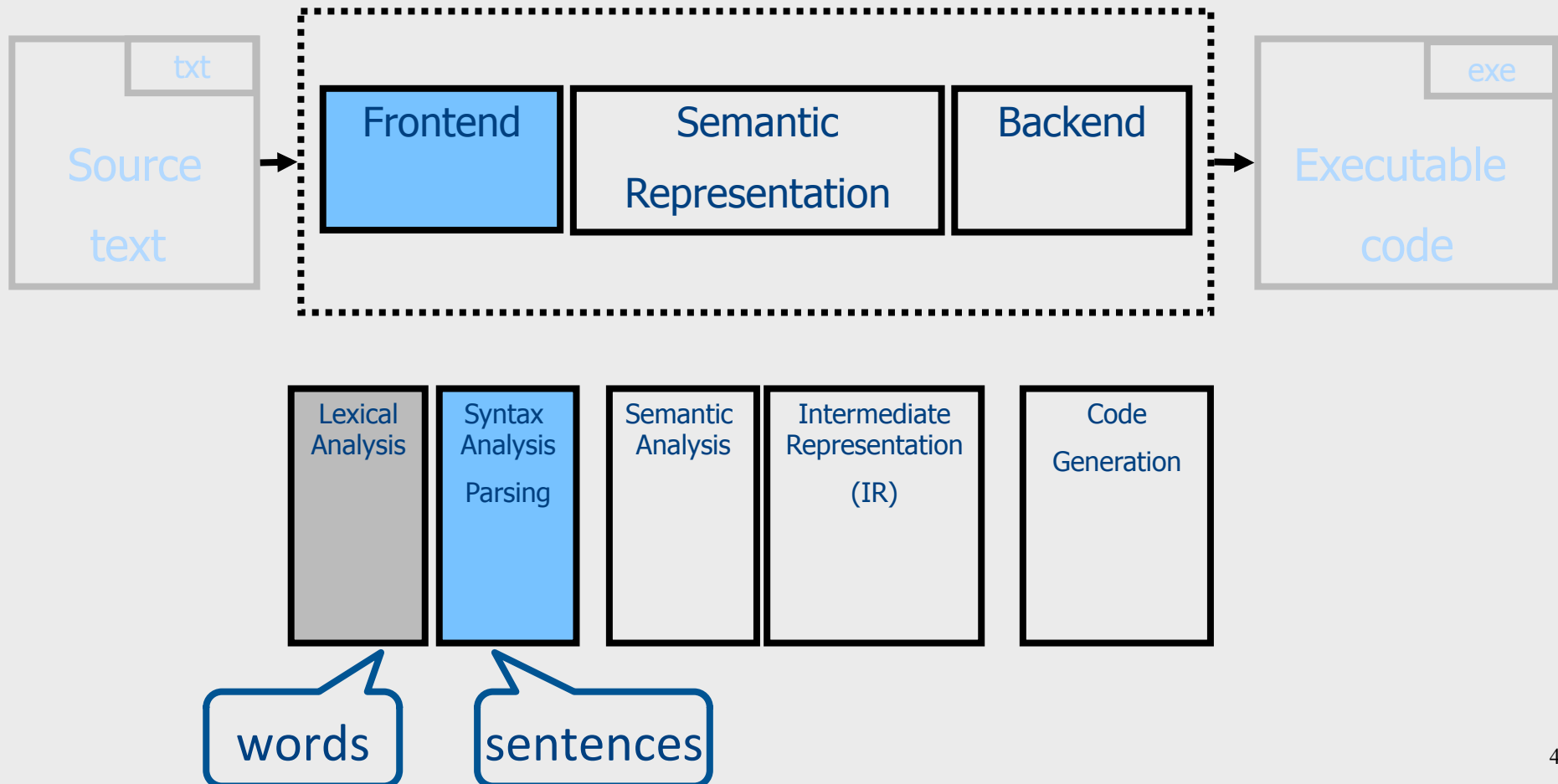
“A compiler is a **computer program** that **transforms** source code written in a programming language (**source language**) into another language (**target language**).

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

*--Wikipedia*

# Conceptual Structure of a Compiler

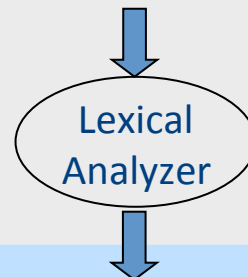
## Compiler



# From scanning to parsing

*program text*

$((23 + 7) * x)$



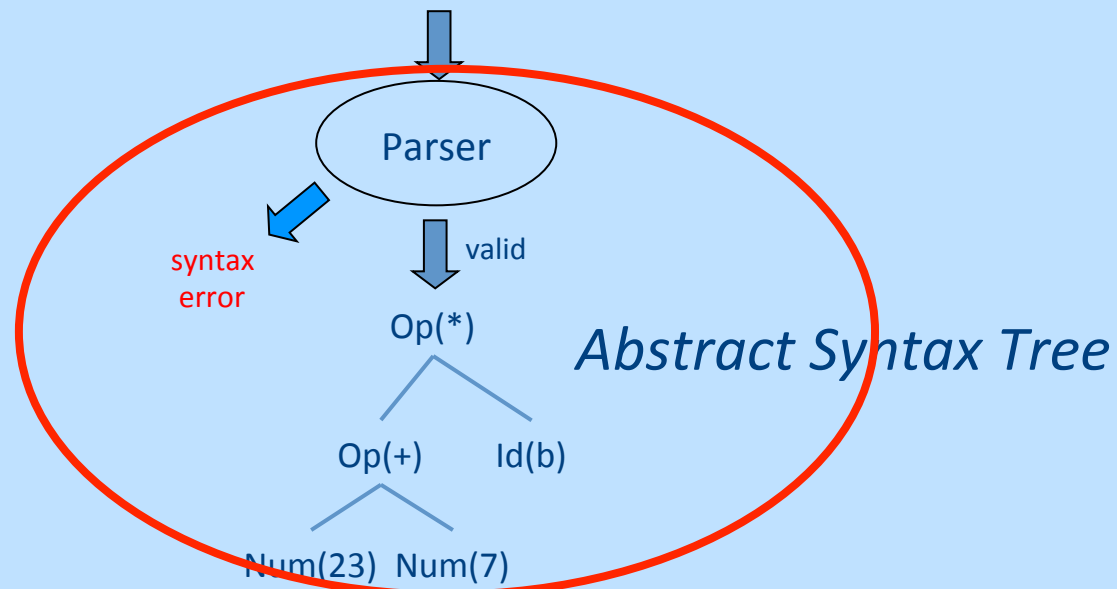
*token stream*

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

Grammar:

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

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



# 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

# 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)*

$N \rightarrow \mu$

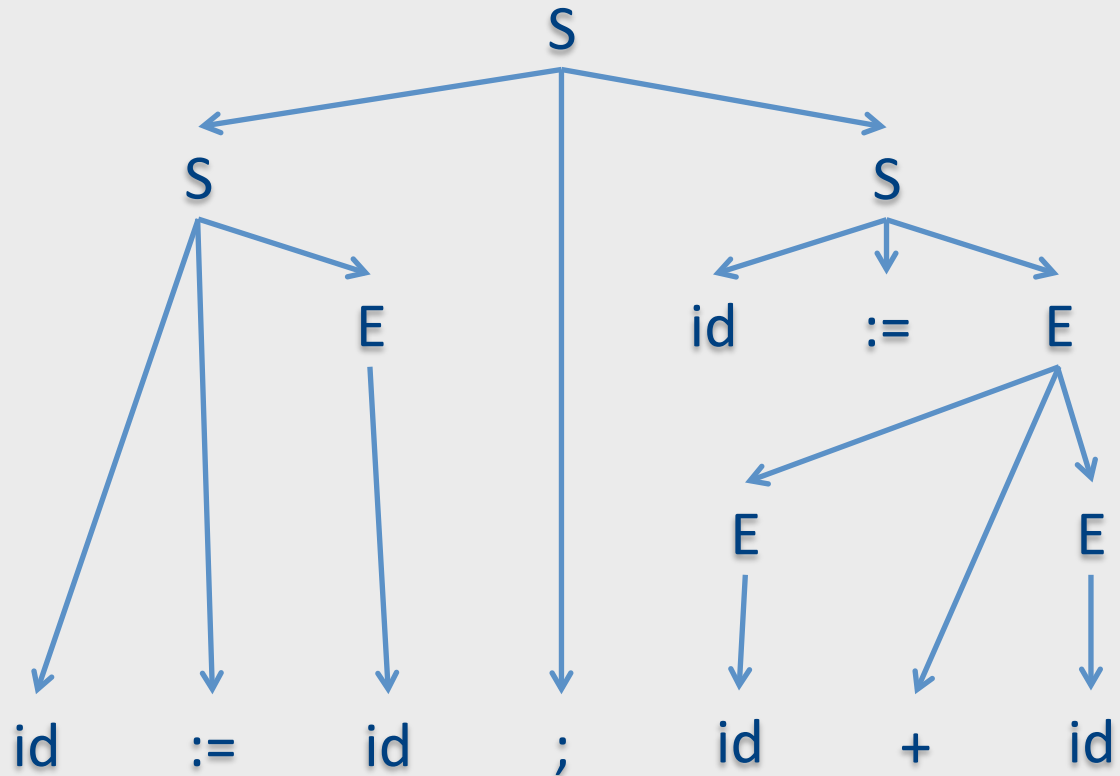
# CFG terminology

- **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



# 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



# 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**

$S \rightarrow S;S$

$S \rightarrow id := E$

$E \rightarrow id$

$S \rightarrow id := E$

$E \rightarrow E + E$

$E \rightarrow id$

$E \rightarrow id$

# Broad kinds of parsers

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

# 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

$E \rightarrow E * T$

$E \rightarrow T$

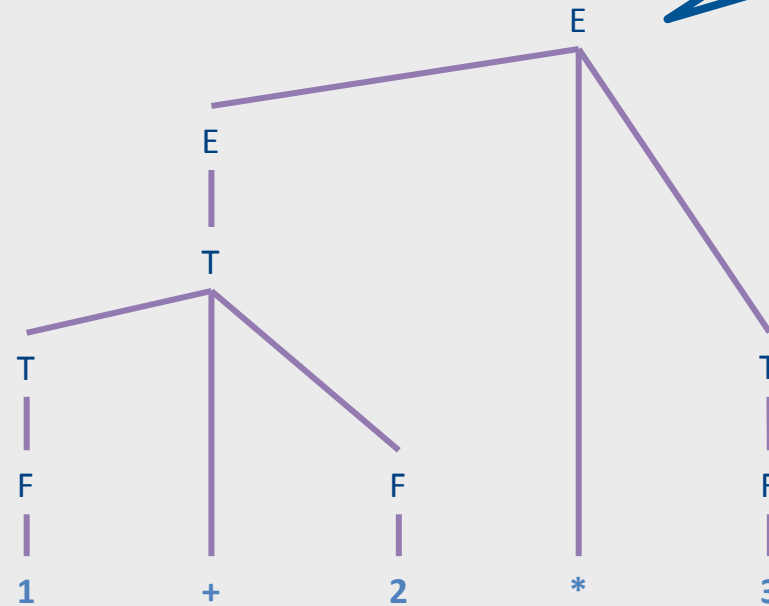
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow ( E )$



*Recall: Non standard  
precedence ...*

# 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**

$E \rightarrow E * T$

$E \rightarrow T$

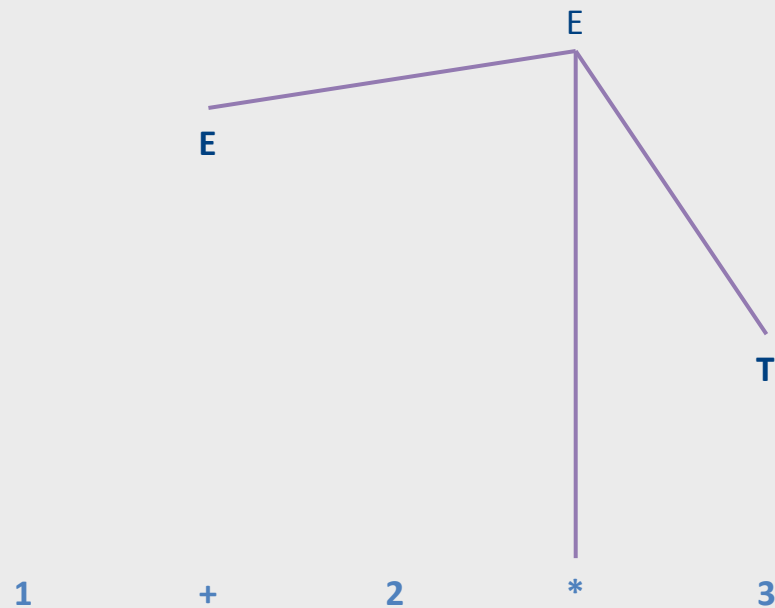
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow ( E )$





# Intuition: Top-Down parsing

## Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

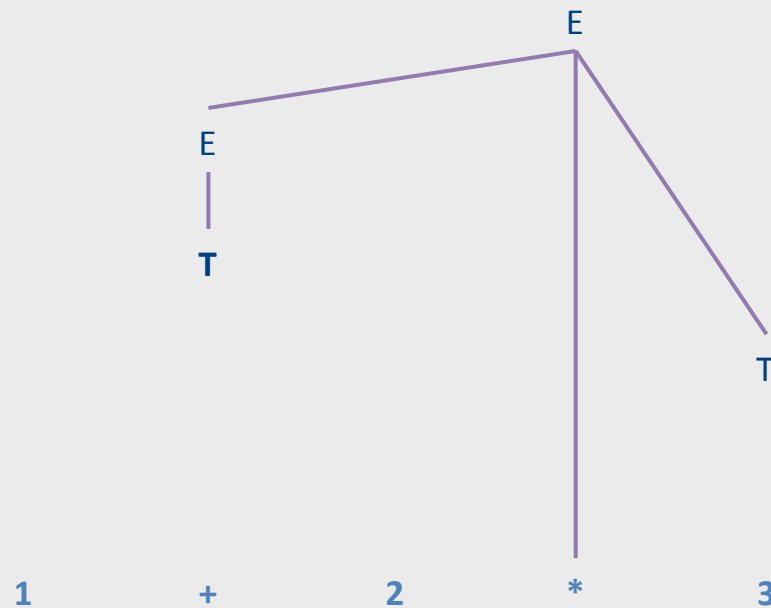
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow ( E )$



# Intuition: Top-Down parsing

## Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

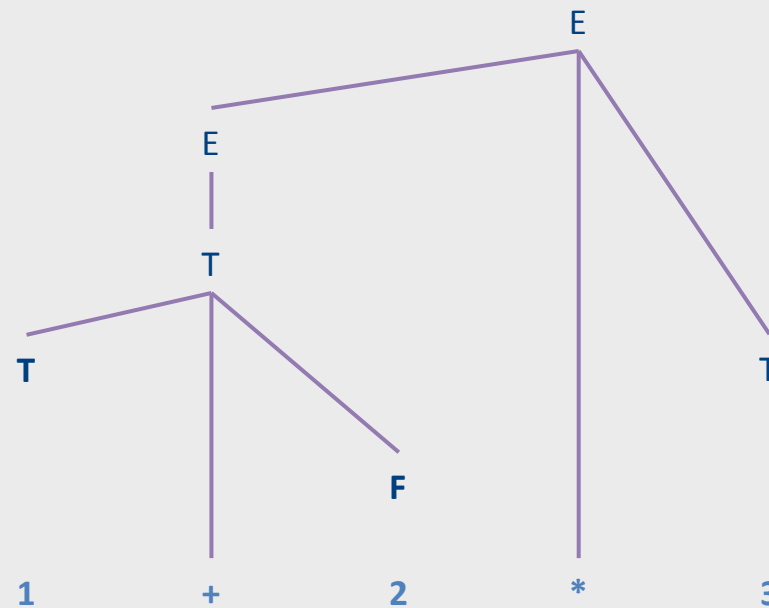
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow ( E )$



# Intuition: Top-Down parsing

## Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

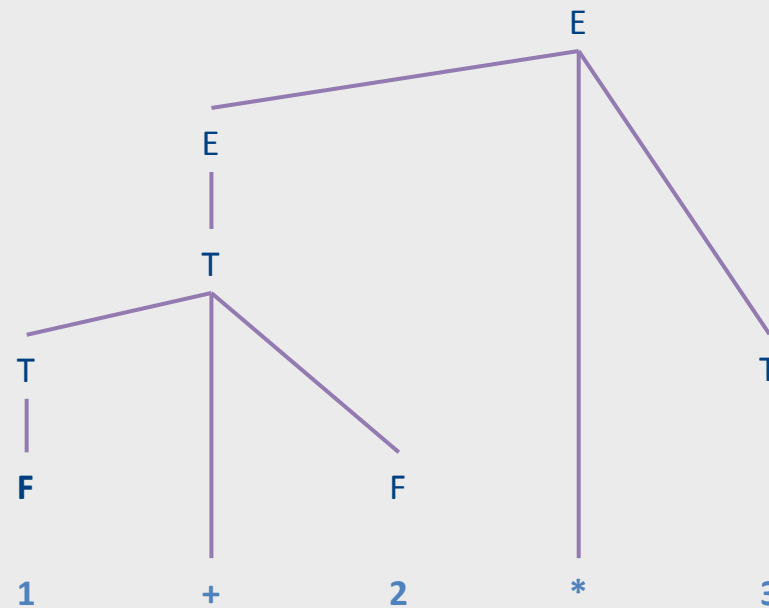
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow ( E )$



# Intuition: Top-Down parsing

## Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

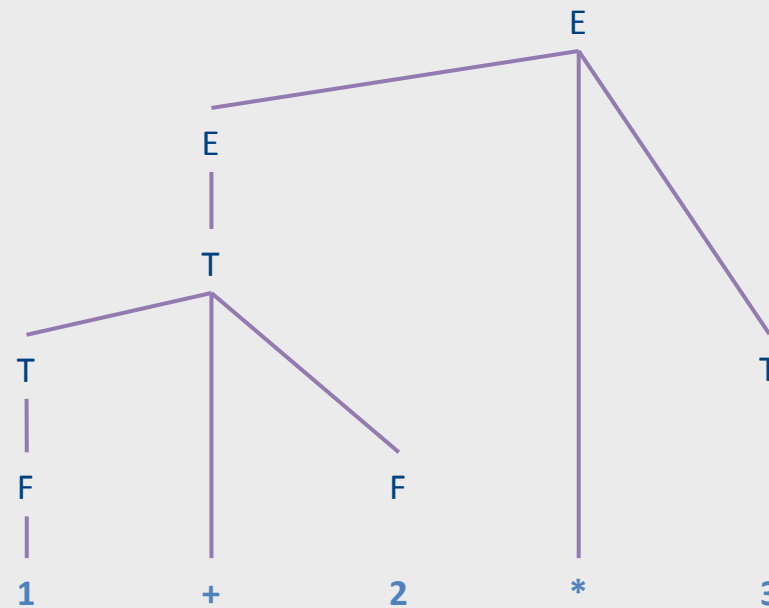
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow ( E )$



# Intuition: Top-Down parsing

## Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

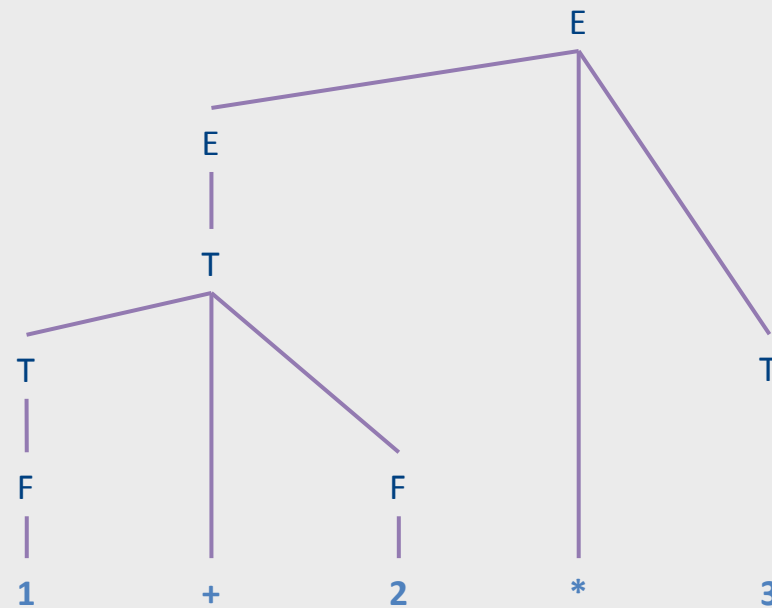
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow ( E )$



# Intuition: Top-Down parsing

## Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

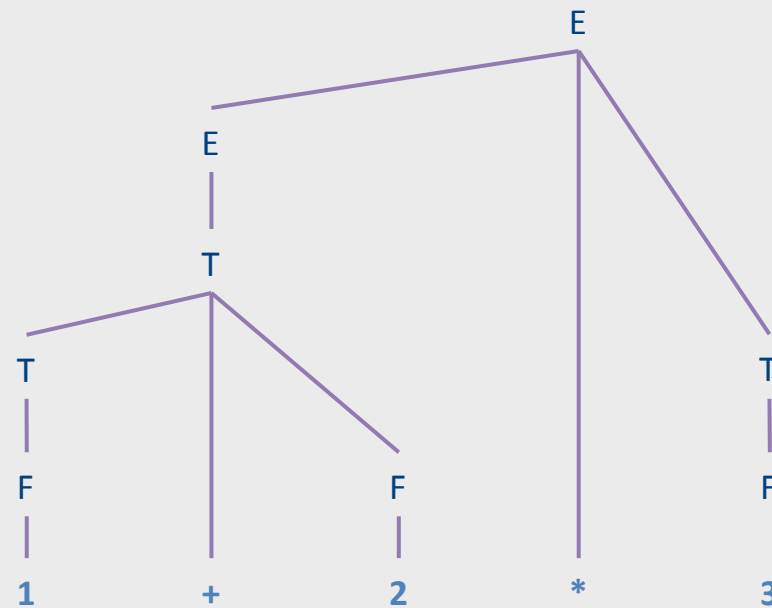
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow ( E )$



# Intuition: Top-Down parsing

## Unambiguous grammar

$E \rightarrow E * T$

$E \rightarrow T$

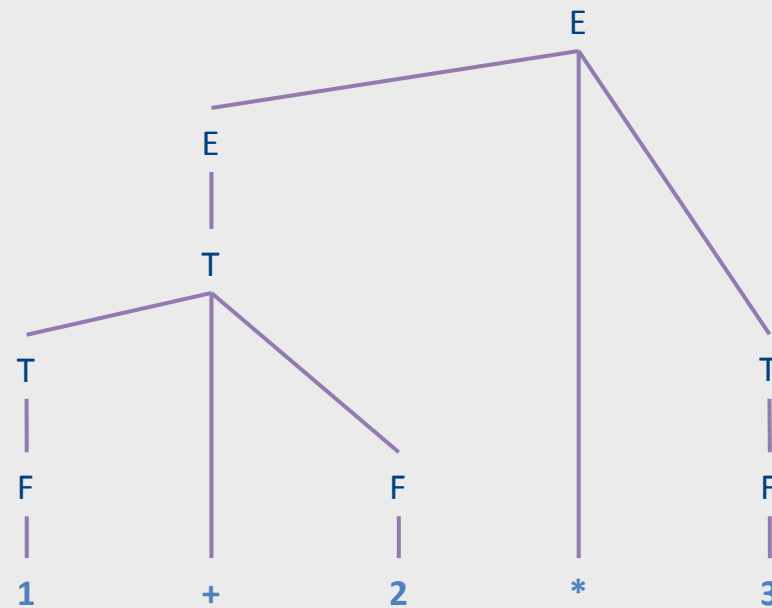
$T \rightarrow T + F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow ( E )$



# 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?



# Recursive Descent

- Blind exhaustive search
  - Goes over all possible production rules
  - Read & parse prefixes of input
  - Backtracks if guesses wrong
- Implementation
  - Uses (possibly recursive) functions for every production rule
  - Backtracks → “rewind” input

# Recursive descent

```
bool A() { // A → A1 | ... | An
    pos = recordCurrentPosition();

    for (i = 1; i ≤ n; i++) {
        if (Ai())
            return true;
        rewindCurrent(pos);
    }
    return false;
}

bool Ai() { // Ai = X1X2...Xk
    for (j=1; j ≤ k; j++)
        if (Xj is a terminal)
            if (Xj == current) match(current);
            else return false;
        else if (! Xj()) return false;
    return true;
}
```

*token stream*

current  
↓

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

# Example

- Grammar
  - $E \rightarrow T \mid T + E$
  - $T \rightarrow \text{int} \mid \text{int} * T \mid ( E )$
- Input: ( 5 )
- Token stream: LP int RP

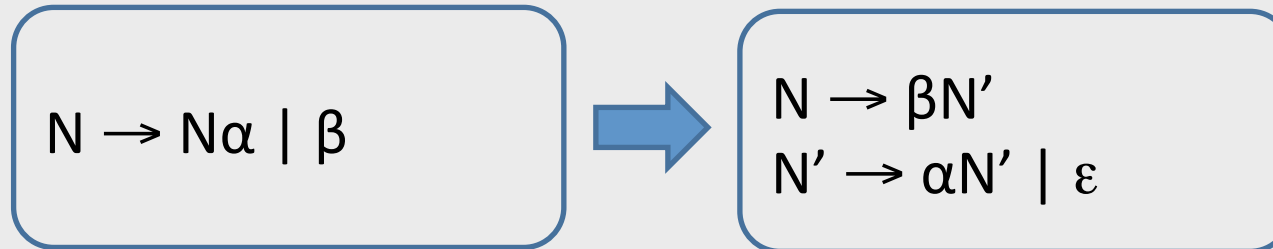
# Problem: Left Recursion

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

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

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

# Left recursion removal



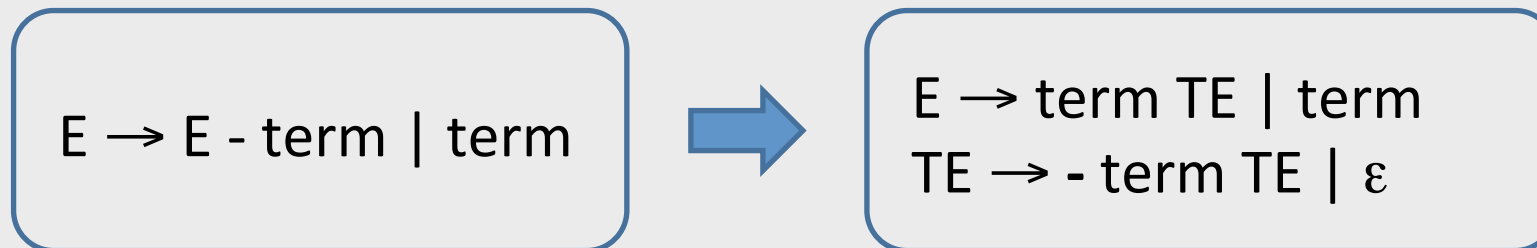
$G_1$

$G_2$

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

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

- For our 3<sup>rd</sup> example:



# 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?
- **Wanted:** Top-Down parsing without backtracking

# Predictive parsing

- Given a grammar  $G$  and a word  $w$  derive  $w$  using  $G$ 
  - 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 based on
    - Non terminal
    - Next (k) tokens (lookahead)

# 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**  $E \Rightarrow$

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

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

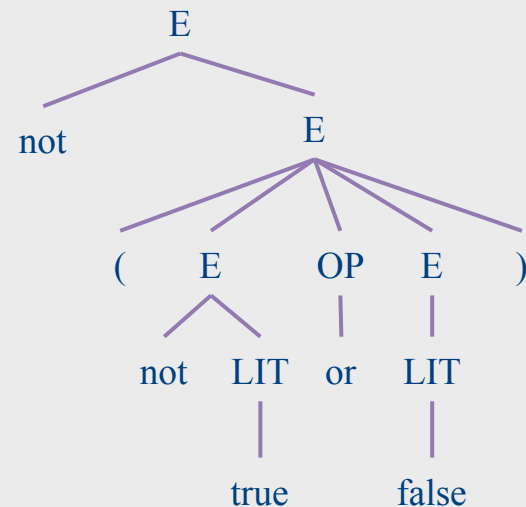
$\text{not} ( \text{not } LIT \text{ OP } E ) \Rightarrow$

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

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

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

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



## Problem: productions with common prefix

term  $\rightarrow$  ID | ID [ expr ]

- Cannot tell which rule to use based on lookahead (ID)

# Solution: left factoring

- Rewrite the grammar to be in LL(1)

term  $\rightarrow$  ID | ID [ expr ]



term  $\rightarrow$  ID after\_ID  
After\_ID  $\rightarrow$  [ expr ] |  $\epsilon$

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 |  $\epsilon$

# LL(k) Parsers

- Predictive parser
  - Can be generated automatically
  - Does not use recursion
  - Efficient
- In contrast, recursive descent
  - Manual construction
  - Recursive
  - Expensive

# LL(k) parsing via pushdown automata and prediction table

- Pushdown automaton uses
  - Prediction stack
  - Input token 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  $t$

# 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

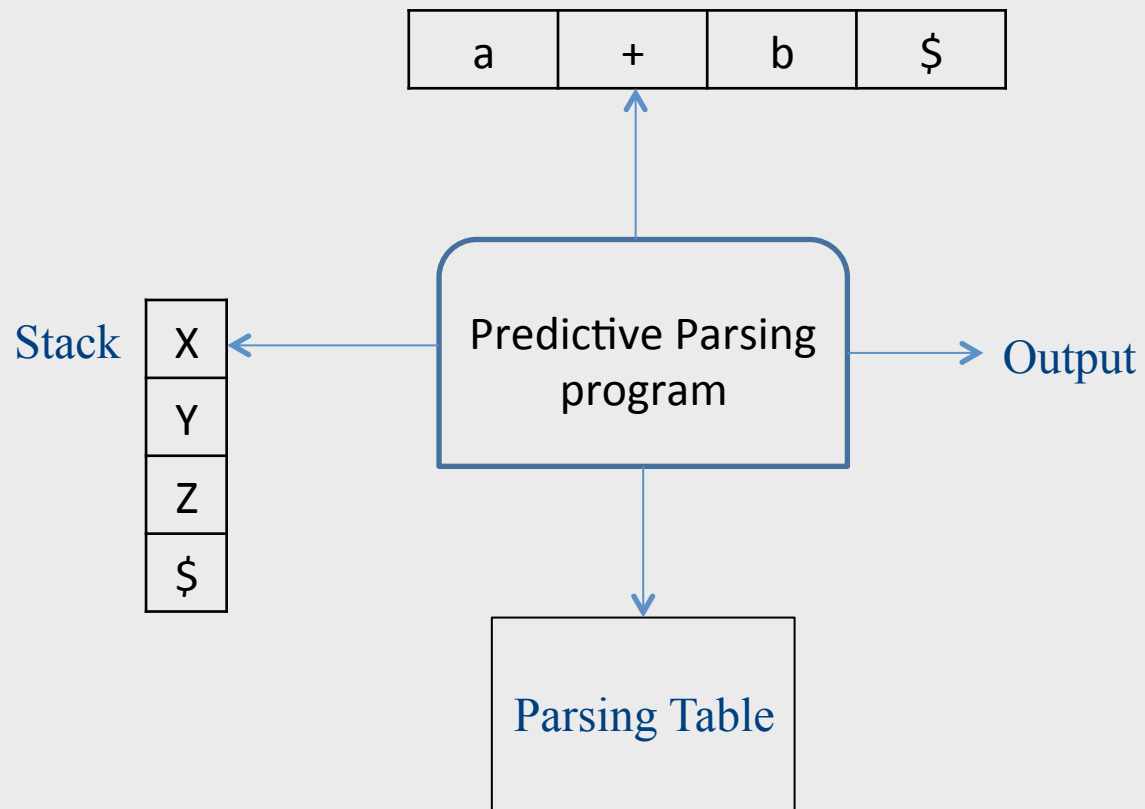
Nonterminals

Input tokens

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

Table

# Non-Recursive Predictive Parser





# LL(k) parsing via pushdown automata and prediction table

- 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

# 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	$A \rightarrow aAb$		$A \rightarrow c$

# FIRST sets

- $\text{FIRST}(\alpha) = \{ t \mid \alpha \rightarrow^* t \beta \} \cup \{ \epsilon \mid X \rightarrow^* \epsilon \}$ 
  - $\text{FIRST}(\alpha)$  = all terminals that  $\alpha$  can appear as first in some derivation for  $\alpha$ 
    - +  $\epsilon$  if can be derived from  $X$
- Example:
  - $\text{FIRST}(\text{LIT}) = \{ \text{true}, \text{false} \}$
  - $\text{FIRST}(( E \text{ OP } E )) = \{ '(' \}$
  - $\text{FIRST}(\text{not } E) = \{ \text{not} \}$

# Computing FIRST sets

- $\text{FIRST}(t) = \{ t \}$  // “t” non terminal
- $\epsilon \in \text{FIRST}(X)$  if
  - $X \rightarrow \epsilon$  or
  - $X \rightarrow A_1 .. A_k$  and  $\epsilon \in \text{FIRST}(A_i) \ i=1\dots k$
- $\text{FIRST}(\alpha) \subseteq \text{FIRST}(X)$  if
  - $X \rightarrow A_1 .. A_k \ \alpha$  and  $\epsilon \in \text{FIRST}(A_i) \ i=1\dots k$

# FIRST sets computation example

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

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

## 2. Iterate 1

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

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

## 2. Iterate 2

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
zero? Not ++ --	id constant	



## 2. Iterate 3 – fixed-point

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

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

# FOLLOW sets

- $\text{FOLLOW}(X) = \{ t \mid S \rightarrow^* \alpha X t \beta \}$ 
  - $\text{FOLLOW}(X)$  = set of tokens that can immediately follow  $X$  in some sentential form
  - NOT related to what can be derived from  $X$
- Intuition:  $X \rightarrow A B$ 
  - $\text{FIRST}(B) \subseteq \text{FOLLOW}(A)$
  - $\text{FOLLOW}(B) \subseteq \text{FOLLOW}(X)$
  - If  $B \rightarrow^* \varepsilon$  then  $\text{FOLLOW}(A) \subseteq \text{FOLLOW}(X)$

# FOLLOW sets: Constraints

- $\$ \in \text{FOLLOW}(S)$
- $\text{FIRST}(\beta) - \{\epsilon\} \subseteq \text{FOLLOW}(X)$ 
  - For each  $A \rightarrow \alpha X \beta$
- $\text{FOLLOW}(A) \subseteq \text{FOLLOW}(X)$ 
  - For each  $A \rightarrow \alpha X \beta$  and  $\epsilon \in \text{FIRST}(\beta)$

# Example: FOLLOW sets

- $E \rightarrow TX$                        $X \rightarrow + E \mid \epsilon$
- $T \rightarrow (E) \mid \text{int } Y$                $Y \rightarrow * T \mid \epsilon$

Terminal	+	(	*	)	int
FOLLOW	int, (	int, (	int, (	_, ), \$	*, ), +, \$

Non. Term.	E	T	X	Y
FOLLOW	), \$	+, ), \$	\$, )	_, ), \$

# Prediction Table

- $A \rightarrow \alpha$
- $T[A,t] = \alpha$  if  $t \in \text{FIRST}(\alpha)$
- $T[A,t] = \alpha$  if  $\epsilon \in \text{FIRST}(\alpha)$  and  $t \in \text{FOLLOW}(A)$ 
  - $t$  can also be  $\$$
- $T$  is not well defined  $\Rightarrow$  the grammar is not LL(1)

# Problem: Non LL Grammars

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

```
bool S() {  
    return A() && match(token('a')) && match(token('b'));  
}
```

```
bool A() {  
    return match(token('a')) || true;  
}
```

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

# Problem: Non LL Grammars

$S \rightarrow A a b$

$A \rightarrow a \mid \varepsilon$

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

$S \rightarrow a \text{ after\_}A$   
 $\text{after\_}A \rightarrow a b \mid b$



# LL(k) grammars

- A grammar is LL(k) if 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)
  - T is well defined
- A language is said to be LL(k) when it has an LL(k) grammar

# LL(k) grammars

- A grammar is not LL(k) if it is
  - Ambiguous
  - Left recursive
  - Not left factored
  - ...

# Earley Parsing

- Invented by Earley [PhD. 1968]
- Handles arbitrary CFG
- Can handle ambiguous grammars
- Complexity  $O(N^3)$  when  $N = |\text{input}|$
- Uses dynamic programming
  - Compactly encodes ambiguity

# Dynamic programming

- Break a problem  $P$  into subproblems  $P_1 \dots P_k$ 
  - Solve  $P$  by combining solutions for  $P_1 \dots P_k$
  - Memo-ize (store) solutions to subproblems instead of re-computation
- Bellman-Ford shortest path algorithm
  - $Sol(x,y,i) = \text{minimum of}$ 
    - $Sol(x,y,i-1)$
    - $Sol(t,y,i-1) + \text{weight}(x,t)$  for edges  $(x,t)$

# Earley Parsing

- Dynamic programming implementation of a recursive descent parser
  - $S[N+1]$  Sequence of sets of “Earley states”
    - $N = |\text{INPUT}|$
    - Earley state (item)  $s$  is a sentential form + aux info
  - $S[i]$  All parse tree that can be produced (by a RDP) after reading the first  $i$  tokens
    - $S[i+1]$  built using  $S[0] \dots S[i]$

# Earley States

- $s = \langle \text{constituent}, \text{back} \rangle$ 
  - **constituent** (dotted rule) for  $A \rightarrow \alpha\beta$ 
    - $A \rightarrow \bullet \alpha\beta$  **predicated** constituents
    - $A \rightarrow \alpha \bullet \beta$  **in-progress** constituents
    - $A \rightarrow \alpha\beta \bullet$  **completed** constituents
  - **back** previous Early state in derivation

# Earley Parser

Input =  $x[1...N]$

$S[0] = \langle E' \rightarrow \bullet E, 0 \rangle$ ;  $S[1] = \dots$   $S[N] = \{\}$

for  $i = 0 \dots N$  do

  until  $S[i]$  does not change do

    foreach  $s \in S[i]$

      if  $s = \langle A \rightarrow \dots \bullet a \dots, b \rangle$  and  $a = x[i+1]$  then

$S[i+1] = S[i+1] \cup \{ \langle A \rightarrow \dots a \bullet \dots, b \rangle \}$  // scan

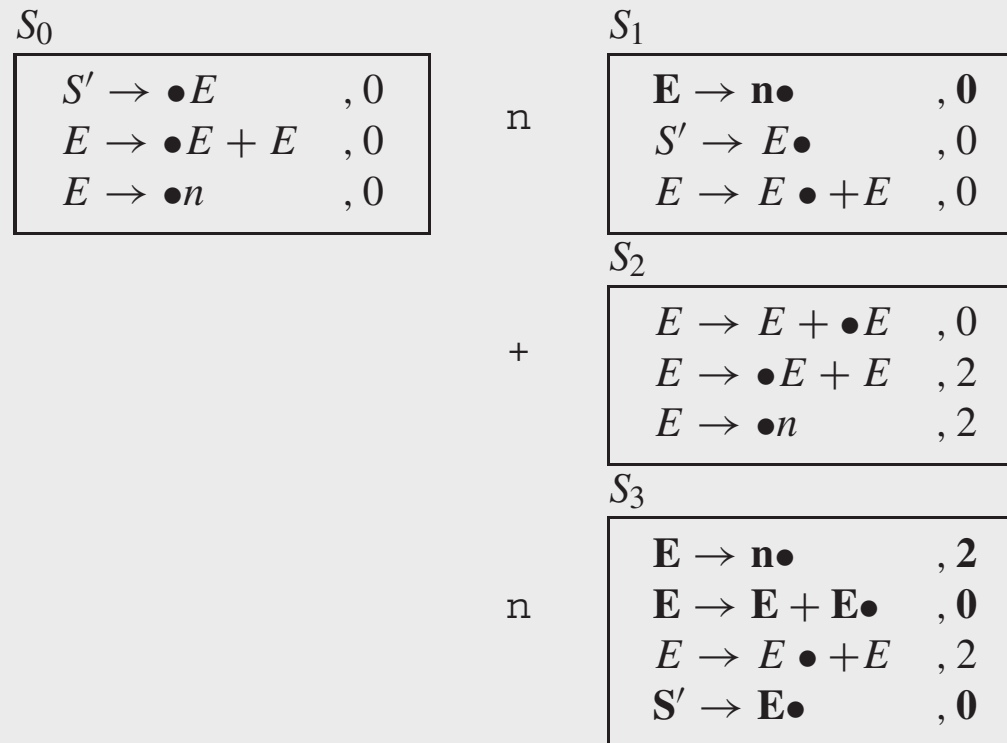
      if  $s = \langle A \rightarrow \dots \bullet X \dots, b \rangle$  and  $X \rightarrow \alpha$  then

$S[i] = S[i] \cup \{ \langle X \rightarrow \bullet \alpha, i \rangle \}$  // predict

      if  $s = \langle A \rightarrow \dots \bullet, b \rangle$  and  $\langle X \rightarrow \dots \bullet A \dots, k \rangle \in S[b]$  then

$S[i] = S[i] \cup \{ \langle X \rightarrow \dots A \bullet \dots, k \rangle \}$  // complete

# Example



**FIGURE 1.** Earley sets for the grammar  $E \rightarrow E + E \mid n$  and the input  $n + n$ . Items in bold are ones which correspond to the input's derivation.



