

**Compiler Construction**  
**Lent Term 2020**  
**Parsing Part I**  
**Lecture 14**

**Timothy G. Griffin**  
**tgg22@cam.ac.uk**  
**Computer Laboratory**  
**University of Cambridge**

# Lecture 14 -- 16

## Parsing (some) context-free grammars

- 1. Context-free grammars**
- 2. The ambiguity problem**
- 3. Top-down parsing (LL(1))**
  - 1. Recursive descent parsing (hand coded)**
  - 2. Predictive parsing (table driven)**
- 4. Bottom-up parsing (Lectures 15,16)**
  - 1. SLR(1) (Simple LR)**
  - 2. LR(1)**

# Programming Language Syntax

## 6.7 Declarations

### Syntax

*declaration:*

*declaration-specifiers init-declarator-list<sub>opt</sub> ;*  
*static\_assert-declaration*

*declaration-specifiers:*

*storage-class-specifier declaration-specifiers<sub>opt</sub>*  
*type-specifier declaration-specifiers<sub>opt</sub>*  
*type-qualifier declaration-specifiers<sub>opt</sub>*  
*function-specifier declaration-specifiers<sub>opt</sub>*  
*alignment-specifier declaration-specifiers<sub>opt</sub>*

*init-declarator-list:*

*init-declarator*  
*init-declarator-list , init-declarator*

*init-declarator:*

*declarator*  
*declarator = initializer*

A small fragment of the C standard. How can we turn this specification into a parser that reads a text file and produces a syntax tree?

# Context-Free Grammars (CFGs)

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

$N$  : set of nonterminals

$T$  : set of terminals

$P \subseteq N \times (N \cup T)^*$  : a set of productions

$S \in N$  : start symbol

Each  $(A, \alpha) \in P$  is written as  $A \rightarrow \alpha$

## Example CFG

$$G_1 = (N_1, T_1, P_1, E)$$

$$N_1 = \{E\} \quad T_1 = \{+, *, (, ), \text{id}\}$$

$P_1$ :

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

This is shorthand for

$$P_1 = \{(E, E + E), (E, E * E), (E, (E)), (E, \text{id})\}$$

# Derivations

Notation convention s :

$$\alpha, \beta, \gamma, \dots \in (N \cup T)^*$$

$$A, B, C, \dots \in N$$

Given :  $\alpha A \beta$  and a production  $A \rightarrow \gamma$

a derivation step is written as

$$\alpha A \beta \Rightarrow \alpha \gamma \beta$$

$\Rightarrow^+$  means one or more derivation steps and

$\Rightarrow^*$  means zero or more derivation steps <sup>6</sup>

# Example derivations

$$E \Rightarrow E * E$$

$$\Rightarrow (E) * E$$

$$\Rightarrow (E + E) * E$$

$$\Rightarrow (x + E) * E$$

$$\Rightarrow (x + y) * E$$

$$\Rightarrow (x + y) * (E)$$

$$\Rightarrow (x + y) * (E + E)$$

$$\Rightarrow (x + y) * (z + E)$$

$$\Rightarrow (x + y) * (z + x)$$

$$E \Rightarrow E * E$$

$$\Rightarrow E * (E)$$

$$\Rightarrow E * (E + E)$$

$$\Rightarrow E * (E + x)$$

$$\Rightarrow E * (z + x)$$

$$\Rightarrow (E) * (z + x)$$

$$\Rightarrow (E + E) * (z + x)$$

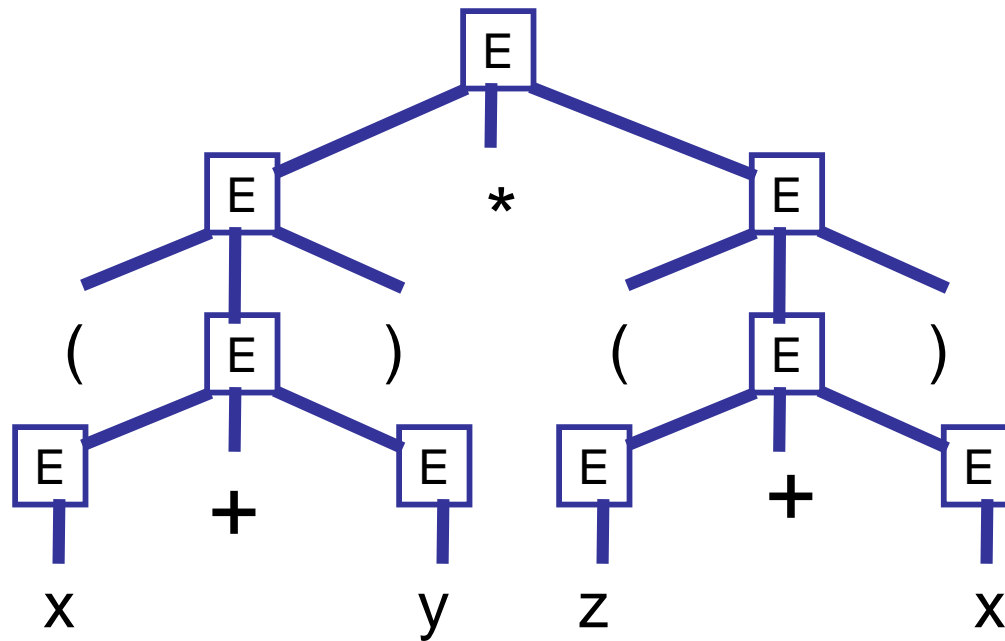
$$\Rightarrow (E + y) * (z + x)$$

$$\Rightarrow (x + y) * (z + x)$$

**A leftmost derivation**

**A rightmost derivation**

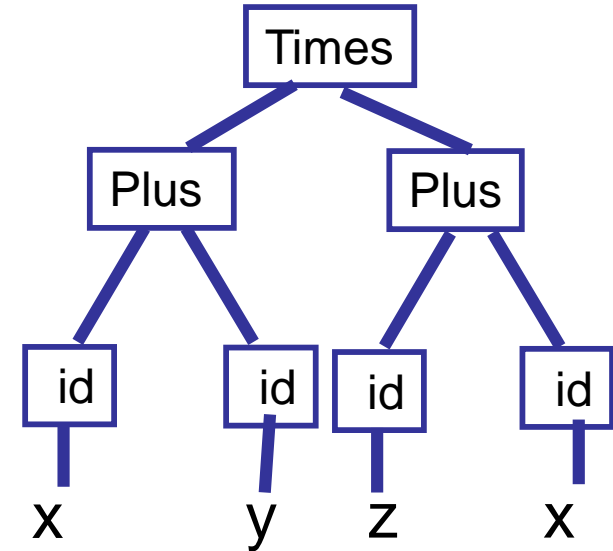
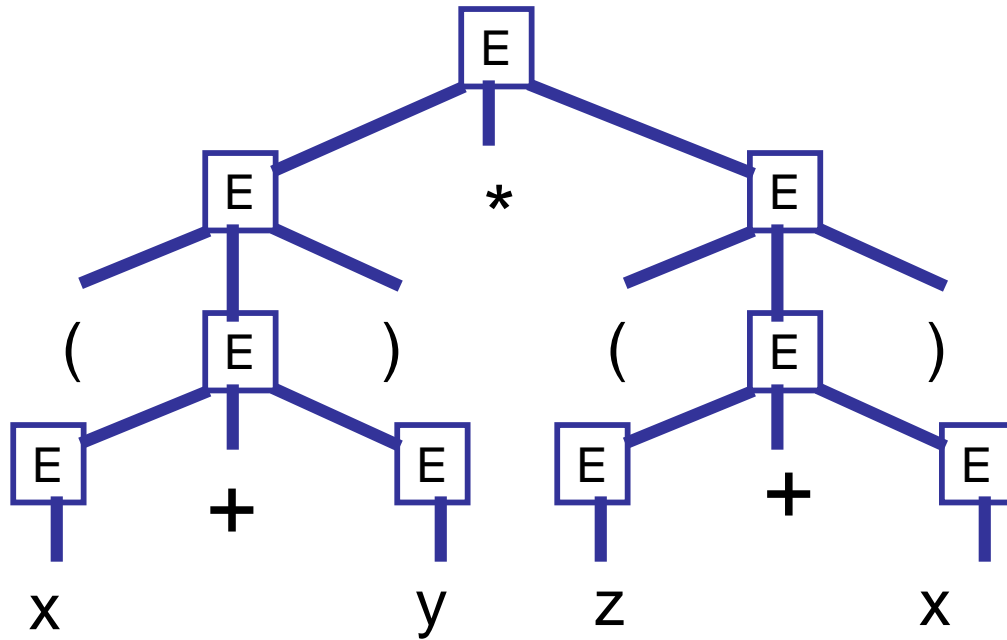
# Derivation Trees



The derivation tree for  $(x + y) * (z + x)$ .  
All derivations of this expression will produce the same derivation tree.



# Concrete vs. Abstract Syntax Trees



**parse tree =  
derivation tree =  
concrete syntax tree**

**An AST contains only the  
information needed to  
generate an intermediate  
representation**

## $L(G)$ = The Language Generated by Grammar $G$

$$L(G) = \{w \in T^* / S \Rightarrow^+ w\}$$

For example, if  $G$  has production  $s$

$$S \rightarrow aSb \mid \varepsilon$$

then

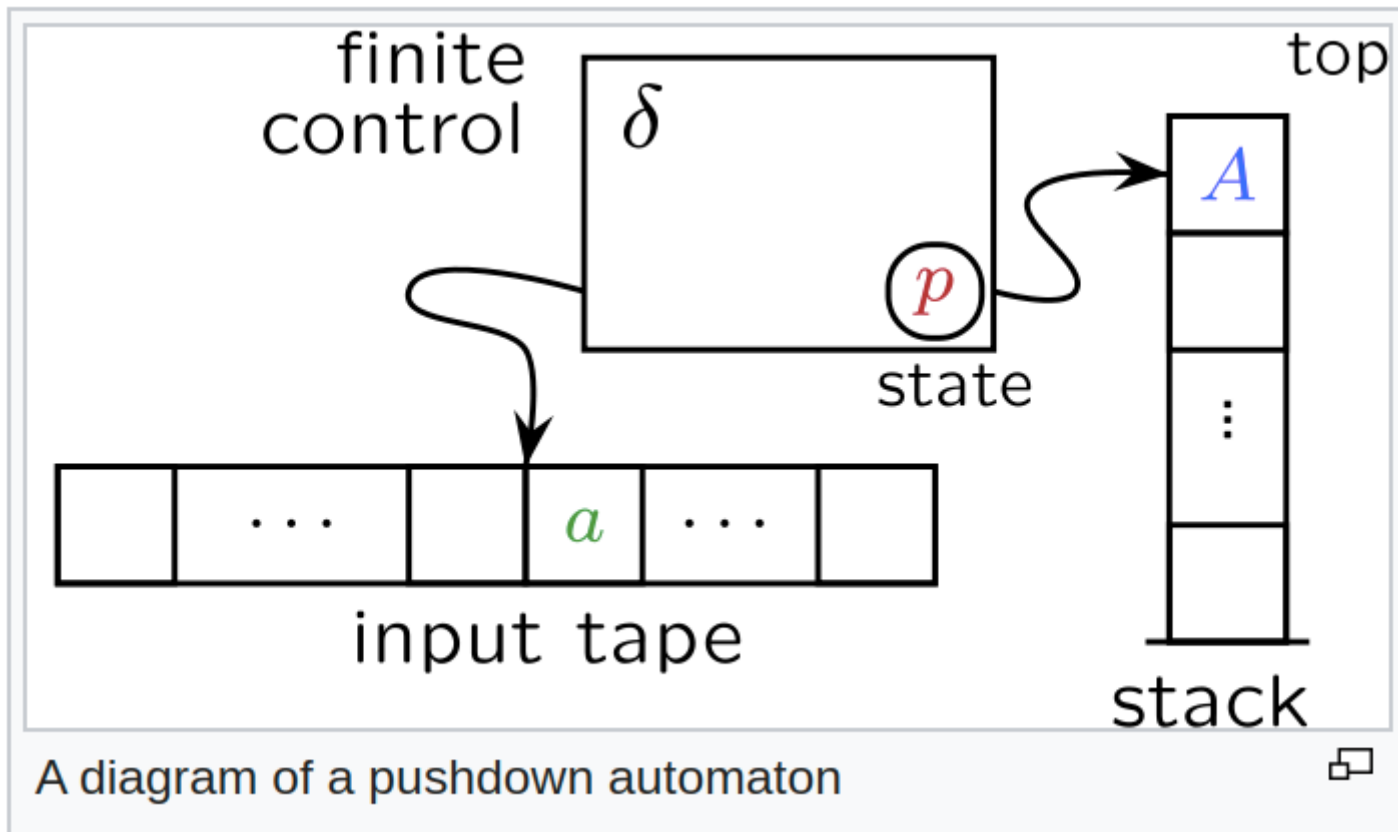
$$L(G) = \{a^n b^n / n \geq 0\}.$$

So CFGs can capture more than regular languages!

# Pushdown Automata (PDAs)

**Regular languages are accepted by Finite Automata. Context-free languages are accepted by Pushdown Automata, a finite automata augmented with a stack.**

**Illustration from [https://en.wikipedia.org/wiki/Pushdown\\_automaton](https://en.wikipedia.org/wiki/Pushdown_automaton)**



# Pushdown Automata (PDAs)

$$M = (Q, \Sigma, \Gamma, \delta, q_0, Z)$$

$Q$  : states     $\Sigma$  : alphabet     $\Gamma$  : stack symbols

$q_0 \in Q$  : start state

$Z \in \Gamma$  : initial stack symbol

$$\delta : \forall q \in Q, a \in (\Sigma \cup \{\varepsilon\}), X \in \Gamma,$$

$$\delta(q, a, X) \subseteq Q \times \Gamma^*$$

# Pushdown Automata (PDAs)

$(q', \beta) \in \delta(q, a, X)$  means that when the machine is in state  $q$  reading  $a$  with  $X$  on top of the stack, it can move to state  $q'$  and replace  $X$  with  $\beta$ . That is, it "pops"  $X$  and "pushes"  $\beta$  (leftmost symbol is top of stack).

# Pushdown Automata (PDAs)

For  $q \in Q, w \in \Sigma^*, \alpha \in \Gamma^*$

$(q, w, \alpha)$

is called an instantaneous

description (ID). It denotes the PDA

in state  $q$  looking at the first symbol

of  $w$ , with  $\alpha$  on the stack (top at left).

## Language accepted by a PDA

For  $(q, \beta) \in \delta(q, a, X)$ ,  $a \in \Sigma$  define the relation  $\rightarrow$  on IDs as

$$(q, aw, X\alpha) \rightarrow (q', w, \beta\alpha)$$

and for  $(q, \beta) \in \delta(q, \varepsilon, X)$  as

$$(q, w, X\alpha) \rightarrow (q', w, \beta\alpha)$$

$$L(M) =$$

$$\{w \in \Sigma^* \mid \exists q \in Q, (q_0, w, Z) \rightarrow^+ (q, \varepsilon, \varepsilon)\}$$

## Exercise : work out the details of this PDA

$(q_0, aaabbb, Z)$

$\rightarrow (q_a, aabbb, A)$

$\rightarrow (q_a, abbb, AA)$

$\rightarrow (q_a, bbb, AAA)$

$\rightarrow (q_b, bb, AA)$

$\rightarrow (q_b, b, A)$

$\rightarrow (q_b, \varepsilon, \varepsilon)$

$L(M) =$

$\{ a^n b^n \mid n \geq 0 \}$



# PDAs and CFGs Facts

(we will not prove them)

1) For every CFG  $G$  there is a PDA  $M$  such that  $L(G) = L(M)$ .

2) For every PDA  $M$  there is a CFG  $G$  such that  $L(G) = L(M)$ .

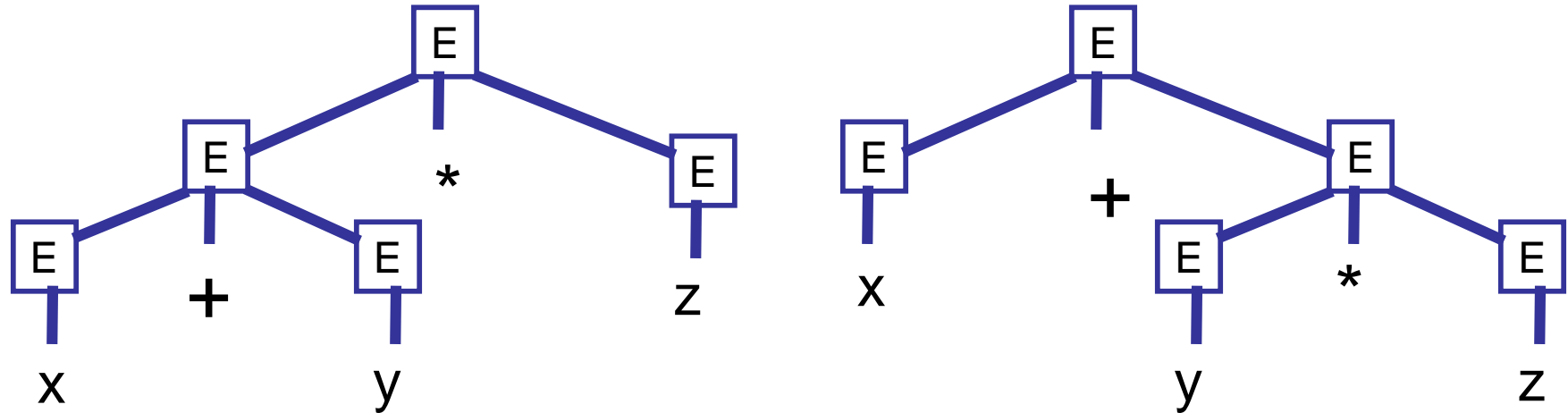
---

Parsing problem solved? Given a CFG  $G$  just construct the PDA  $M$ ? Not so fast!

For programming languages we want

$M$  to be deterministic!

# Origins of nondeterminism? Ambiguity!



Both derivation trees correspond "x + y \* z".  
But  $(x+y) * z$  is not the same as  $x + (y * z)$ .

This type of ambiguity will cause problems  
when we try to go from program texts to  
derivation trees! Semantic ambiguity!

$$G_2 = (N_2, T_1, P_2, E)$$

$$N_2 = \{E, T, F\} \quad T_1 = \{+, *, (, ), \text{id}\}$$

$P_2$  :

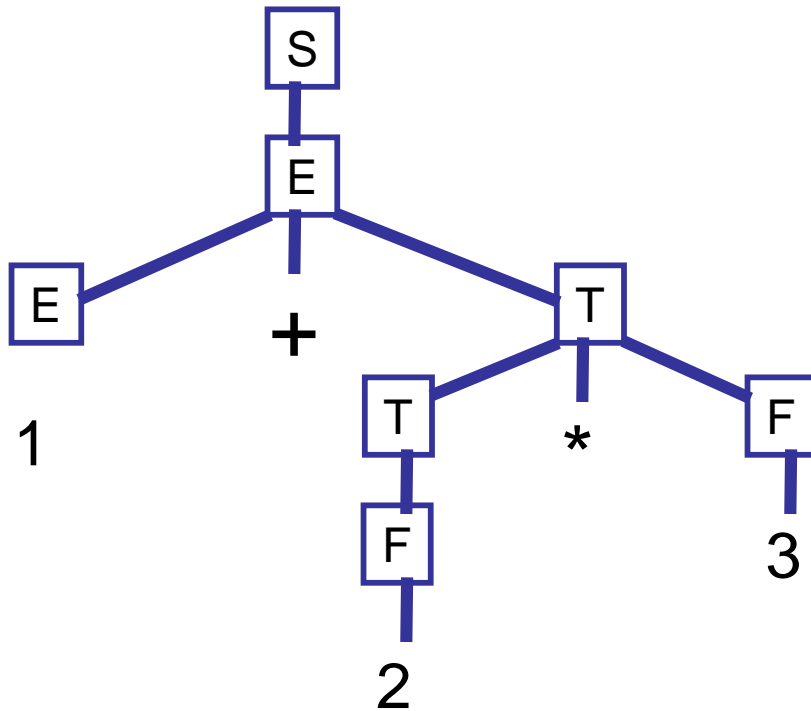
$$E \rightarrow E + T \mid T \quad (\text{expressions})$$

$$T \rightarrow T * F \mid F \quad (\text{terms})$$

$$F \rightarrow (E) \mid \text{id} \quad (\text{factors})$$

Can you prove that  $L(G_1) = L(G_2)$ ?

# The modified grammar eliminates ambiguity



This is now  
the unique  
derivation  
tree for  
 $x + y * z$

# Fun Fun Facts

- (1) Some context-free languages are *inherently ambiguous* --- every context-free grammar for them will be ambiguous. For example:

$$L = \left\{ a^n b^n c^m d^m \mid m \geq 1, n \geq 1 \right\} \\ \cup \left\{ a^n b^m c^m d^n \mid m \geq 1, n \geq 1 \right\}$$

- (2) Checking for ambiguity in an arbitrary context-free grammar is not decidable! Ouch!
- (3) Given two grammars  $G_1$  and  $G_2$ , checking  $L(G_1) = L(G_2)$  is not decidable! Ouch!

# Two approaches to building stack-based parsing machines: top-down and bottom-up

- **Top Down** : attempts a left-most derivation. We will look at two techniques:
  - Recursive decent (hand coded)
  - Predictive parsing (table driven)
- **Bottom-up** : attempts a right-most derivation backwards. We will look at two techniques:
  - SL(1) : Simple LR(1)
  - LR(1)

Bottom-up techniques are strictly more powerful. That is, they can parse more grammars.

# Recursive Descent Parsing

(G5)

```
S ::= if E then S else S
    | begin S L
    | print E
```

```
E ::= NUM = NUM
```

```
L ::= end
    | ; S L
```

Parse corresponds to  
a left-most derivation  
constructed in  
a “top-down” manner

```
int tok = getToken();

void advance() {tok = getToken();}
void eat (int t) {if (tok == t) advance(); else
error();}

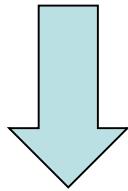
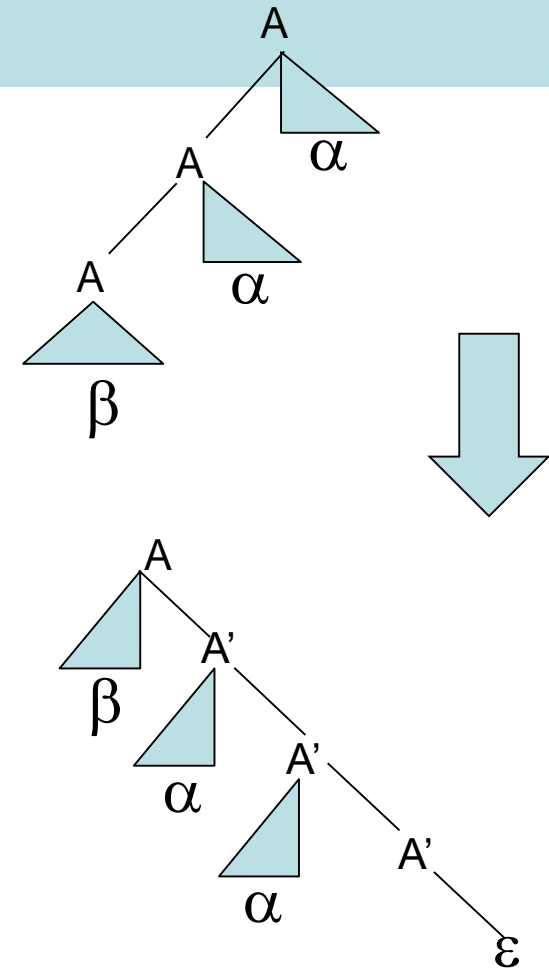
void S() {switch(tok) {
    case IF:    eat(IF); E(); eat(THEN);
                S(); eat(ELSE); S(); break;
    case BEGIN: eat(BEGIN); S(); L(); break;
    case PRINT: eat(PRINT); E(); break;
    default: error();
}}

void L() {switch(tok) {
    case END:   eat(END); break;
    case SEMI: eat(SEMI); S(); L(); break;
    default: error();
}}

void E() {eat(NUM) ; eat(EQ); eat(NUM); }
```

But the production  $E \rightarrow T$  in  $G_2$  will lead to an infinite loop!

Eliminate left recursion!

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_k \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$$
$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_k A' \mid \varepsilon$$


For eliminating left-recursion in general, see Aho and Ullman.<sup>24</sup>



# Eliminate left recursion

$$G_3 = (N_3, T_1, P_3, E)$$

$$N_2 = \{E, E', T, T', F\} \quad T_1 = \{+, *, (, ), \text{id}\}$$

$P_2$ :

$$E \rightarrow T E'$$

$$E' \rightarrow +T E' \mid \varepsilon$$

$$T \rightarrow F T'$$

$$T' \rightarrow * F T' \mid \varepsilon$$

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

Can you prove that  $L(G_2) = L(G_3)$ ?

# Recursive descent pseudocode

$\text{getE}() = \text{getT}(); \text{getE}'()$

$\text{getE}'() = \text{if token}() = "+" \text{ then eat}("+"); \text{getT}(); \text{getE}'()$

$\text{getT}() = \text{getF}(); \text{getT}'()$

$\text{getT}'() = \text{if token}() = "*" \text{ then eat}("*"); \text{getF}(); \text{getT}'()$

$\text{getF}() = \text{if token}() = id$

$\text{then eat}(id)$

$\text{else eat}("("); \text{getE}(); \text{eat}(")")$

# Where's the stack machine? It's implicit in the call stack!

Parsing  $(x+y)^*(z+x)$  using a call to `getE()`

				<code>eat("(")</code>	<code>getE()</code>	
			<code>getF()</code>	<code>getF()</code>	<code>getF()</code>	
	<code>getT()</code>	<code>getT()</code>	<code>getT()</code>	<code>getT()</code>		...
<code>getE()</code>	<code>getE()</code>	<code>getE()</code>	<code>getE()</code>	<code>getE()</code>	<code>getE()</code>	



call stack over time ...