

Lecture 8: Recursive-descent parsing

- **Recursive-descent formalized**
 - **FIRST sets**
 - **LL(1) condition**
 - **Transformations to LL(1) form**
- **Grammars for expressions - a difficult case**

(Next week: LR(1) parsing, ocamllyacc)

Top-down parsing

- For each non-terminal with productions:

$$A \rightarrow \vec{X} \mid \vec{Y} \mid \dots \mid \vec{Z}$$

$$\downarrow X_1 X_2 \dots X_n, \quad n \geq 0, X_i \in G$$

define parseA:

parseA toklis = choose production based on hd toklis:

if $A \rightarrow \vec{X}$: handle \vec{X} ; produce A

else if $A \rightarrow \vec{Y}$: handle \vec{Y} , etc

handle $X_1 X_2 \dots X_n$: handle X_1 ; handle X_2 ; ...; handle X_n

where handle t : if hd toklis = t

then remove t and continue

else error

handle B : parseB toklis

“choose production based on hd toklis”

derives in one step

- Need to formalize some things...
- Define “ \Rightarrow ”: $X_1 \dots X_n \Rightarrow w_1 \dots w_n$, where $X_i \in G$ and $w_i \in G^*$, if there exists j such that $X_j \rightarrow w_j$ is a production in G , and for all $i \neq j$, $X_i = w_i$.
- \Rightarrow^+ and \Rightarrow^* are the transitive and reflexive-transitive closures of \Rightarrow . (Say \vec{X} derives α if $\vec{X} \Rightarrow^* \alpha$.)

(E.g. α is a sentence of G if the start symbol of G derives α and α consists solely of tokens.)

$$\begin{array}{l|l}
 A \rightarrow (B) & A \Rightarrow (B) \Rightarrow (A) \Rightarrow ((B)) \\
 B \rightarrow id \mid A & \Rightarrow ((id)) \\
 B \Rightarrow id &
 \end{array}$$

$$B \Rightarrow A \Rightarrow (A)$$

“choose production based on hd toklis” (cont.) $A \rightarrow A + id$

- \vec{X} is *nullable* if it can derive ϵ .
- Define: $\text{FIRST}(\vec{X}) = \{t \in T \mid \vec{X} \Rightarrow^* t\alpha \text{ for some } \alpha\} \cup \{\bullet \mid \vec{X} \text{ nullable}\}$.
- Define: G is *left-recursive* if $\exists A : A \Rightarrow^+ A\alpha$ for some α .
- Define: G is *LL(1)* if
 1. G is not left-recursive, and
 2. For all non-terminals A , if the productions of A are $A \rightarrow \vec{X} \mid \dots \mid \vec{Y}$, the sets $\text{FIRST}(\vec{X}), \dots, \text{FIRST}(\vec{Y})$ are pairwise disjoint.

↑
sequence of
grammar symbols

$$\text{FIRST}(A) = \{(\}\ , \text{FIRST}(B) = \{id, (\}$$

Top-down parsing revisited

If G is LL(1), then for each non-terminal A with productions

$$A \rightarrow \vec{X} \mid \vec{Y} \mid \dots \mid \vec{Z}$$

construct `parseA`:

```
parseA toklis = let t = hd toklis in
  if t ∈ FIRST( $\vec{X}$ ) then handle  $\vec{X}$ 
  else if t ∈ FIRST( $\vec{Y}$ ) then handle  $\vec{Y}$ 
  ...else if t ∈ FIRST( $\vec{Z}$ ) or • ∈ FIRST( $\vec{Z}$ ) then handle  $\vec{Z}$ 
    ( $\vec{Z}$  is the unique nullable right-hand side of  $A$ , if any)
  else error

handle  $X_1, X_2, \dots, X_n$  : handle  $X_1$ ; handle  $X_2$ ; ...; handle  $X_n$ 

handle t : if hd toklis = t
  then remove t and continue
  else error

handle  $B$  : parseB toklis
```

Transformation to LL(1)

- **Left refactoring:**

$$A \rightarrow \alpha\beta \mid \alpha\gamma$$

$$\Rightarrow A \rightarrow \alpha B$$

$$B \rightarrow \beta \mid \gamma$$

$$\begin{aligned} A &\Rightarrow \alpha B \Rightarrow \alpha\beta \\ &\Rightarrow \alpha B \Rightarrow \alpha\gamma \end{aligned}$$

- **Left-recursion removal:**

$$A \rightarrow A\alpha \mid \beta$$

$$\Rightarrow A \rightarrow \beta B$$

$$B \rightarrow \epsilon \mid \alpha B$$

A has form $\beta\alpha^$*

$$\begin{aligned} A &\rightarrow \beta B \Rightarrow \beta \\ &\Rightarrow \beta\alpha B \Rightarrow \beta\alpha \\ &\Rightarrow \beta\alpha\alpha B \Rightarrow \beta\alpha\alpha \\ &\vdots \end{aligned}$$

Example

- Consider non-LL(1) grammar 3 from the previous class:

$$A \rightarrow \text{id} \mid '(' B ')'$$
$$B \rightarrow A \mid A '+' B$$

- Grammar 3 transformed to LL(1) form:

$$A \rightarrow \text{id} \mid '(' B ')'$$
$$B \rightarrow A C$$
$$C \rightarrow '+' A C \mid \epsilon$$

} left-factoring

Ambiguity

- **No test for ambiguity**
- **Recursive descent and LR(1) parsing not applicable to ambiguous grammar (possible to “cheat” with LR parser - will see how next week)**

Expression grammars

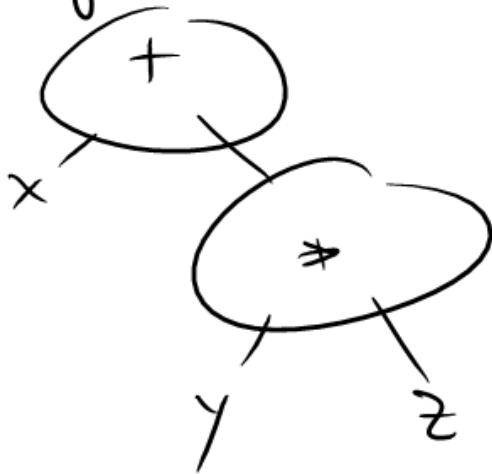
- Expressions are challenging for several reasons:
 - *Should be LL(1) and LR(1)*
 - Grammar should enforce precedence, if possible
 - Grammar should enforce associativity, if possible
 - Grammar shouldn't be ambiguous
 - Should be easy to construct *abstract* syntax tree
- Especially hard to write LL(1) parser for expressions. Not so hard for LR(1).

Enforcing precedence

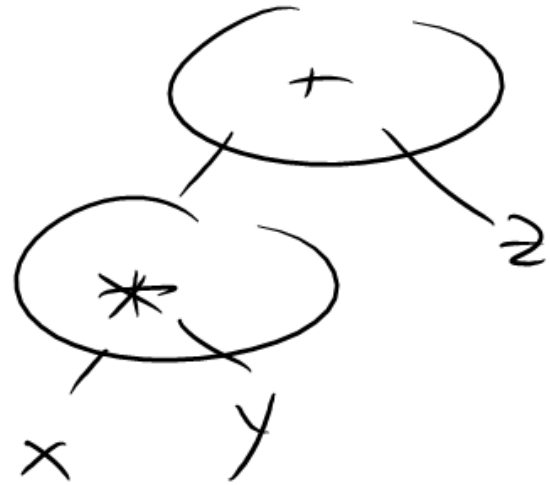
Consider

$$x + y * z$$

Should parse

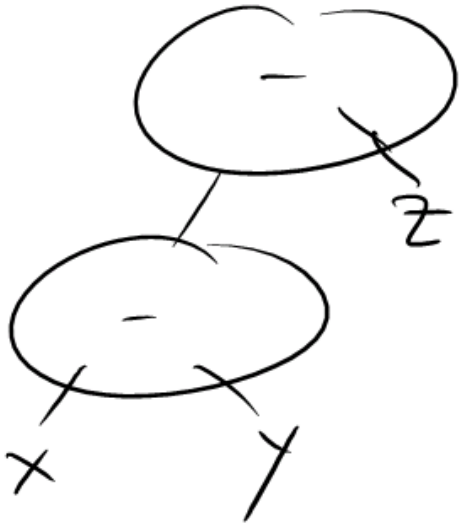


$$x * y + z$$

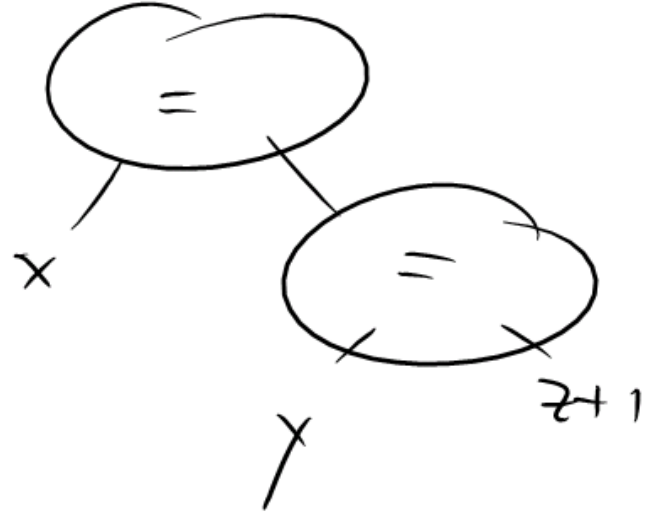


Enforcing associativity

$$x - y - z$$



$$x = y = z + 1$$



Some expression grammars

$$G_A: E \rightarrow \text{id} \mid E - E \mid E * E$$

$$G_B: E \rightarrow \text{id} \mid \text{id} - E \mid \text{id} * E$$

$$G_C: E \rightarrow \text{id} \mid E - \text{id} \mid E * \text{id}$$

$$G_D: E \rightarrow T - E \mid T \\ T \rightarrow \text{id} \mid \text{id} * T$$

$$G_E: E \rightarrow E - T \mid T \\ T \rightarrow \text{id} \mid T * \text{id}$$

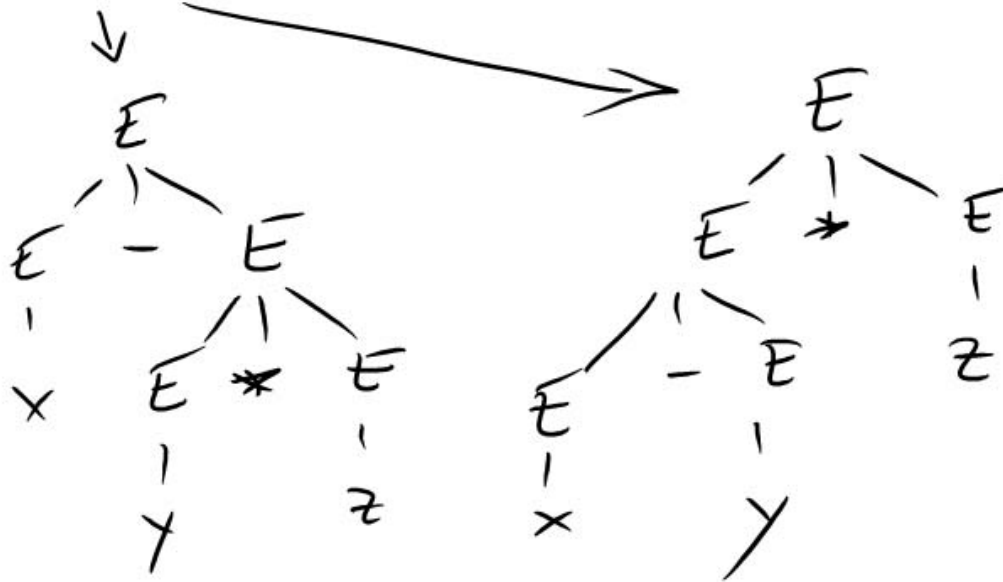
$$G_F: E \rightarrow T E' \\ E' \rightarrow \epsilon \mid - E \\ T \rightarrow \text{id} T' \\ T' \rightarrow \epsilon \mid * T$$

$\text{id} \left(\frac{-}{+} \right) \text{id} \left(\frac{-}{+} \right) - -$

● $G_A: E \rightarrow id \mid E - E \mid E * E$

- Ambiguous
- Does not enforce precedence or associativity

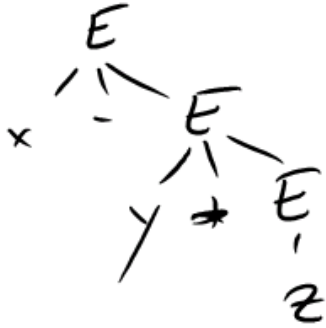
● $x - y * z$



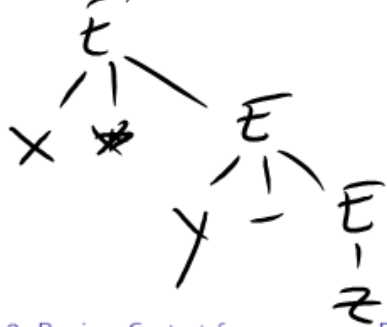
● $G_B: E \rightarrow id \mid id - E \mid id * E$

- Unambiguous
- LR(1)
- Not LL(1), but could be left-factored
- ~~• No precedence~~
- No precedence
- Right associativity

● $x - y * z$



● $x * y - z$



● $x - y - z$



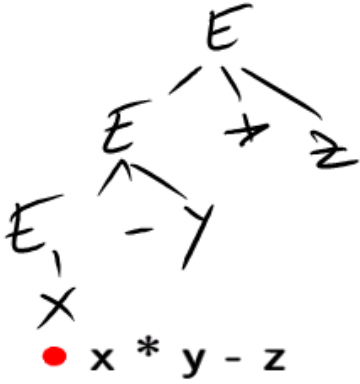
● $G_C: E \rightarrow id \mid E - id \mid E * id$

- Unambiguous
- LR(1), not LL(1) because left-recursive

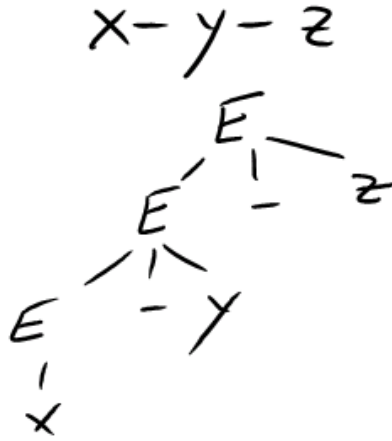
• No precedence

~~Enforces left-associativity~~
• Enforces left-associativity

• $x - y * z$



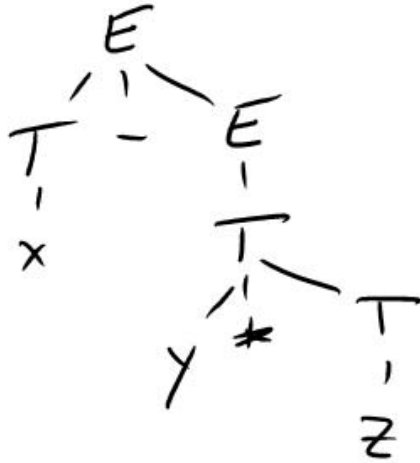
• $x * y - z$



● $G_D: E \rightarrow T - E \mid T$
 $T \rightarrow \text{id} \mid \text{id} * T$

- Unambiguous
- LR(1), not LH(1) (but left-factorable)
- Enforces precedence
- Enforces right-assoc.

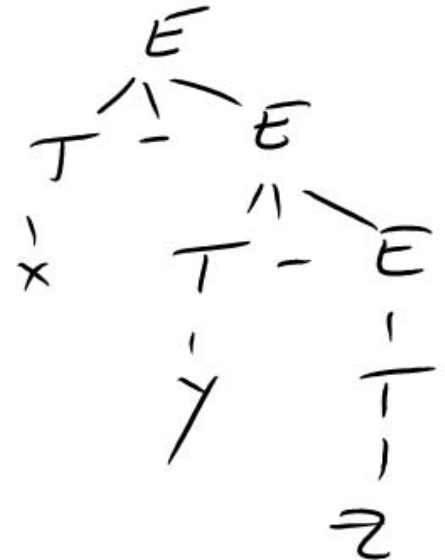
● $x - y * z$



$x * y - z$



$x - y - z$



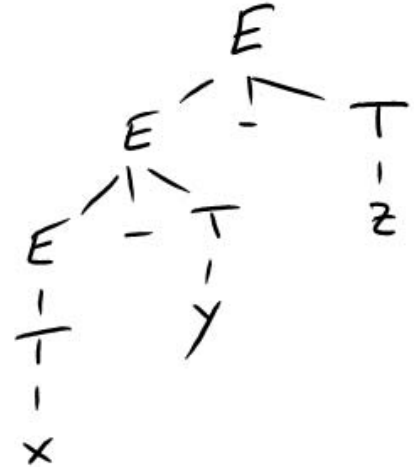
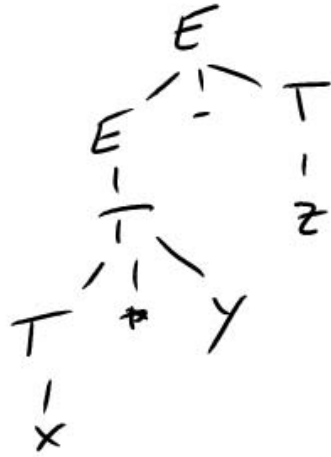
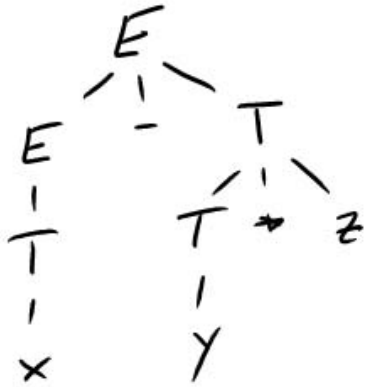
● $G_E: E \rightarrow E - T \mid T$
 $T \rightarrow id \mid T * id$

- Unambiguous
- LR(1), not LL(1) (left-recursive)
- Enforces prec.
- Enforces left-assoc.

● $x - y * z$

$x * y - z$

$x - y - z$

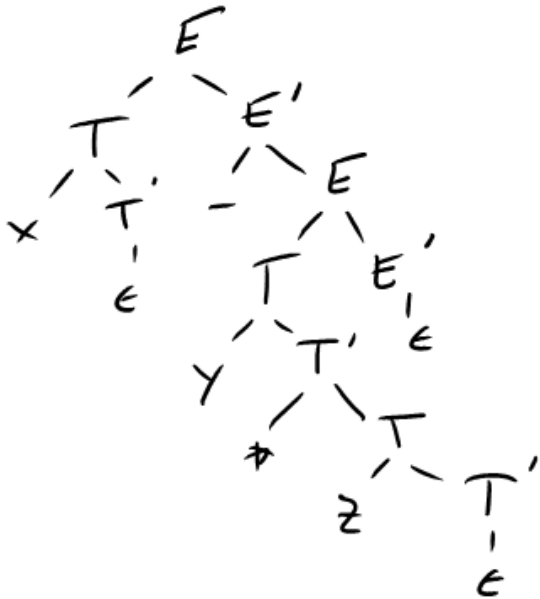


- $G_F: E \rightarrow T E'$
 $E' \rightarrow \epsilon \mid - E$
 $T \rightarrow \text{id } T'$
 $T' \rightarrow \epsilon \mid * T$

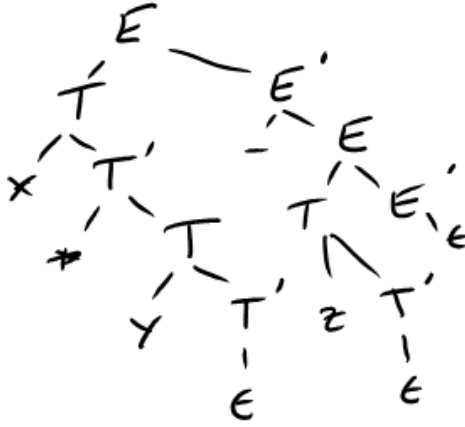
Left-factored version of G_D

- Unambiguous
- LR(1), LL(1)
- Enforces precedence
- Enforces right-asso

• $x - y * z$



$x * y - z$



$x - y - z$

