

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:

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

define parseA:

parseA toklis = choose production based on hd toklis:

if $A \rightarrow \vec{X}$ chosen: handle \vec{X}
else if $A \rightarrow \vec{Y}$ chosen: handle \vec{Y} ,
else if *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”

- Need to formalize some things...
- Define “ \Rightarrow ”: $X_1 \dots X_n \Rightarrow X_1 \dots X_{i-1} \alpha X_{i+1} \dots X_n$ (for any $1 \leq i \leq n$) if the grammar has production $X_i \rightarrow \alpha$.
- \Rightarrow^+ and \Rightarrow^* are the transitive and reflexive-transitive closures of \Rightarrow . (Say \vec{X} derives α if $\vec{X} \Rightarrow^* \alpha$.)
- α is a *sentential form* of G if the start symbol of G derives α . If, furthermore, α consists solely of tokens, then it is a *sentence*. (These notions correspond to being the “frontier” of a syntax tree; some care is needed in defining “frontier” to account for ϵ -productions.)

“choose production based on hd toklis” (cont.)

- \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: $\text{FOLLOW}(A) = \{t \in T \mid \exists \text{ a sentential form } \alpha A t \beta\} \cup \{\text{eof} \mid A \text{ can appear as the last symbol in a sentential form}\}$

There are well-known algorithms for calculating FIRST and FOLLOW sets, but we will consider only simple cases where they can be calculated by inspection.

symbol in a sentential form

$$L \rightarrow (M)$$

$$M \rightarrow \epsilon \mid id \ N$$

$$N \rightarrow \epsilon \mid , \ id \ N$$

$$L \Rightarrow (M)$$

$$\Rightarrow (id \ N)$$

$$\Rightarrow (id , \ id \ N)$$

$$FIRST((M)) = \{ (\}$$

$$FIRST(\epsilon) = \{ \cdot \}$$

$$FIRST(id \ N) = \{ id \}$$

$$FIRST(, \ id \ N) = \{ , \}$$

$$FOLLOW(L) = \{ eof \}$$

$$FOLLOW(M) = \{) \}$$

$$FOLLOW(N) = \{) \}$$

“choose production based on hd toklis” (cont.)

- 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 $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$:
 - (a) The sets $\text{FIRST}(\alpha_1), \dots, \text{FIRST}(\alpha_n)$ are pairwise disjoint.
($\forall i, j. i \neq j \Rightarrow \alpha_i \cap \alpha_j = \emptyset$)
 - (b) If A is nullable, then suppose α_i is the unique right-hand side such that $\bullet \in \text{FIRST}(\alpha_i)$. Then, for all $j \neq i$, $\text{FIRST}(\alpha_j) \cap \text{FOLLOW}(A) = \emptyset$.
should be: $\text{FIRST}(\alpha_i) \cap \text{FIRST}(\alpha_j) = \emptyset$

$\text{FIRST}(\alpha_j)$

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 ( $\bullet$  ∈ FIRST( $\vec{Z}$ ) and t ∈ FOLLOW(A))
    then handle  $\vec{Z}$  ( $\vec{Z}$  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:

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

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

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

- Left-recursion removal:

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

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

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

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

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:

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

B

Ambiguity

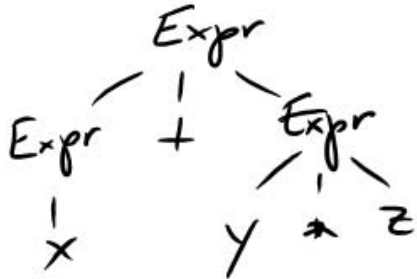
- No test for ambiguity
- Recursive descent and LR(1) parsing not applicable to ambiguous grammar (possible to “cheat” with LR parser - will see how next week)

Expression grammars

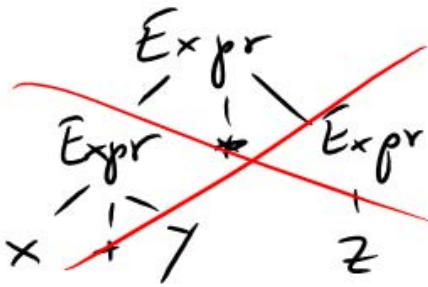
- **Expressions are challenging for several reasons:**
 - **Grammar should enforce precedence, if possible**
 - **Grammar should enforce associativity, if possible**
 - **Grammar shouldn't be ambiguous**
 - **Should be easy to construct *abstract* syntax tree**
- **Especially hard to write LL(1) parser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

Enforcing precedence

$x + y * z$

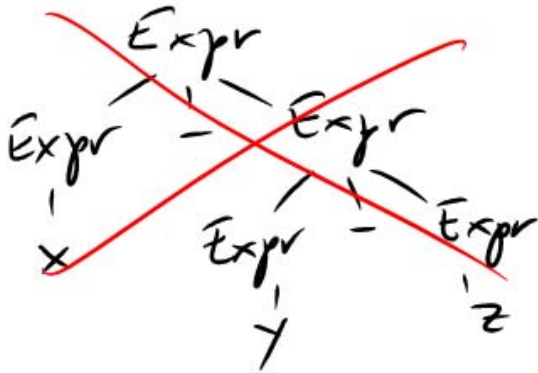


$x * y + z$

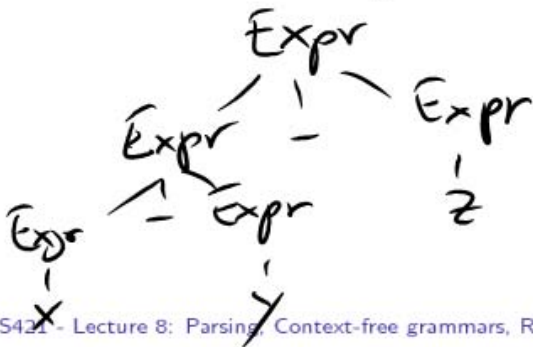
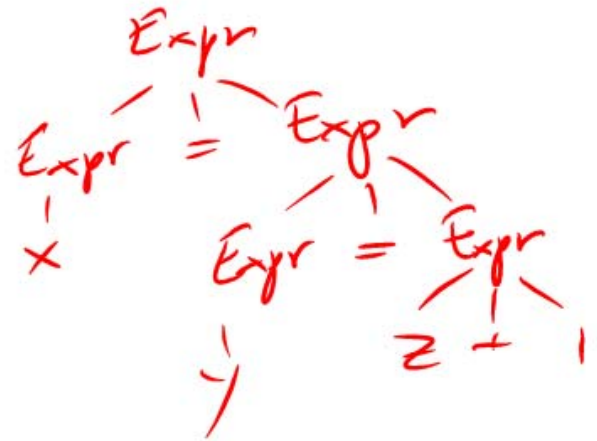


Enforcing associativity

$x - y - z$



$x = y = z + 1$



Some expression grammars

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

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

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

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

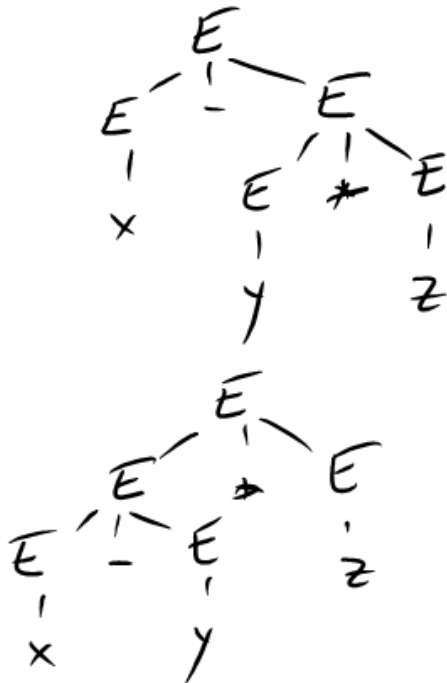
$G_E: E \rightarrow E - T \mid T$
 $T \rightarrow \text{id} \mid T * \text{id}$

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

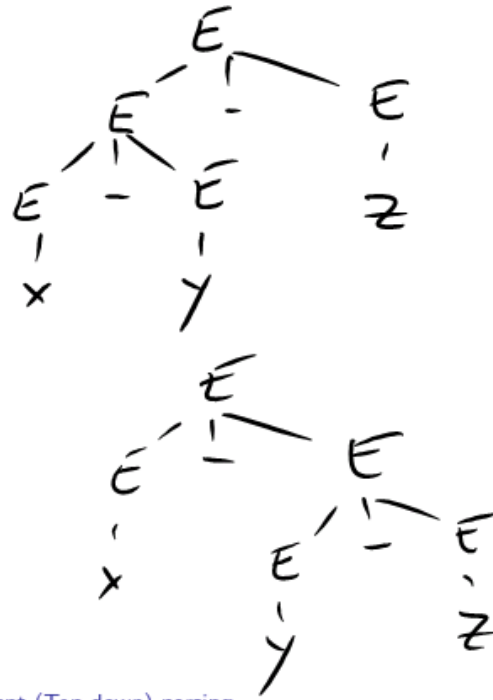
● $G_A: E \rightarrow id \mid E - E \mid E * E$

- Ambiguous
- Does not enforce anything
- Not LL(1) or LR(1)

● $x - y * z$



$x - y - z$



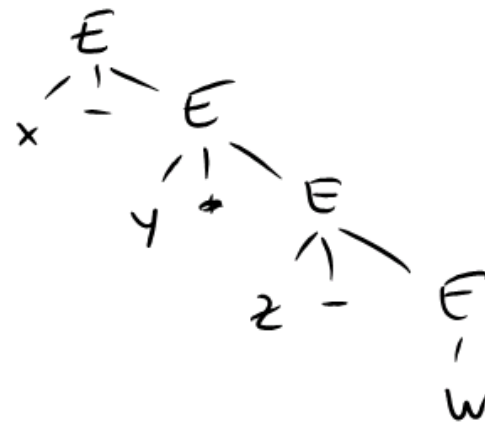
● $G_B: E \rightarrow id \mid id - E \mid id * E$

- Unambiguous
- Does not enforce anything
- Not LL(1) - but can left factor

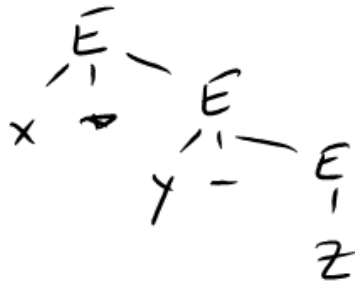
● $x - y * z$



$x - y * z - w$



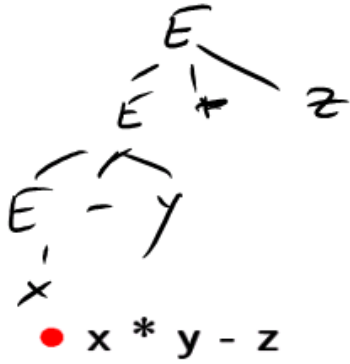
● $x * y - z$



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

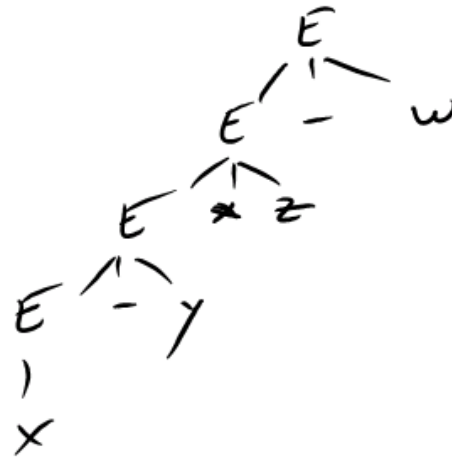
- Unambiguous
- Does not enforce anything
- Not LL(1) (left recursive)

● $x - y * z$



● $x * y - z$

$x - y * z - w$



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

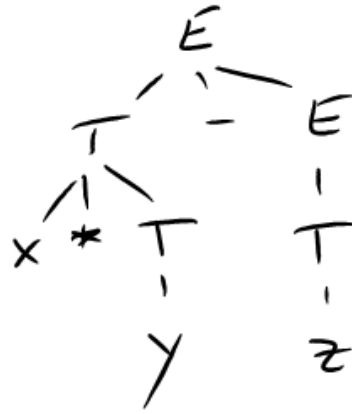
- Unamb.
- Enforces
- Enforce
- LL(1)?

precedence
 right assoc.
 No - can left-factor.

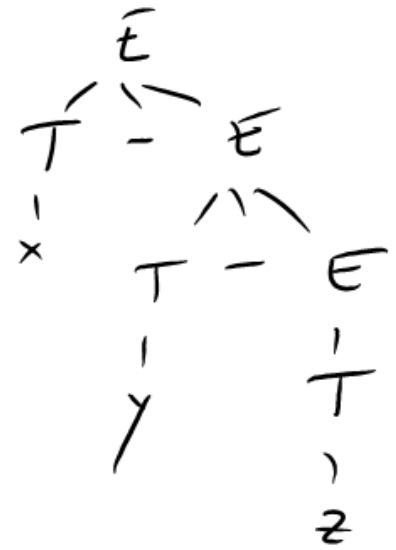
● $x - y * z$



$x * y - z$



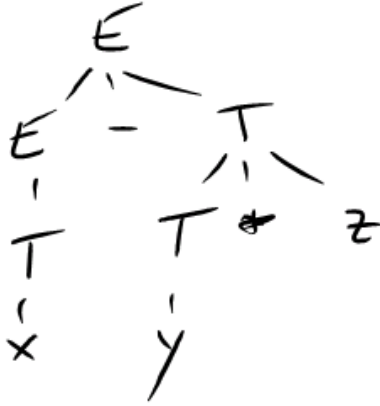
$x - y - z$



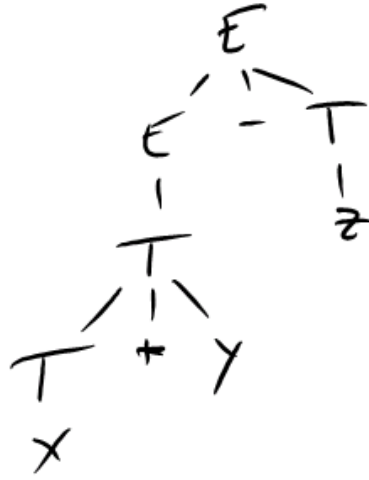
● $G_E: E \rightarrow E - T \mid T$
 $T \rightarrow \text{id} \mid T * \text{id}$

- Unamb.
- Enforces prec. & left assoc.
- LL(1)? No-left recursive

● $x - y * z$



$x * y - z$



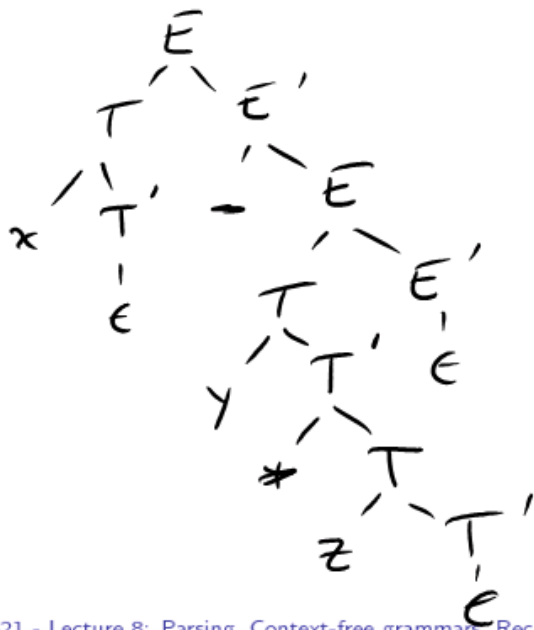
$x - y - z$



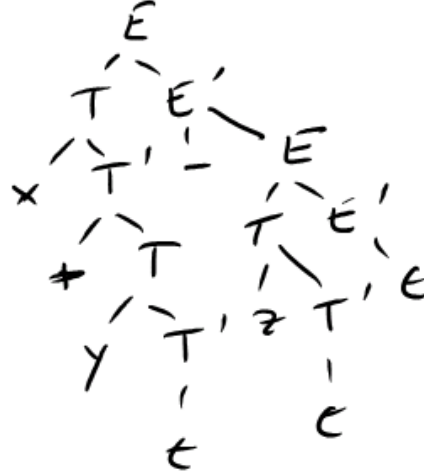
- $G_F: E \rightarrow T E'$
 $E' \rightarrow \epsilon \mid - E$
 $T \rightarrow \text{id } T'$
 $T' \rightarrow \epsilon \mid * T$

- Unamb.
- Enforce prec, right assoc.
- LL(1)? Yes

• $x - y * z$



$x * y - z$



$x - y - z$

