# 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, ocamlyacc)

# Top-down parsing

- **For each non-terminal with productions:**

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

**define parseA:**

$\texttt{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 $X_1 X_2 \ldots X_n$ : handle $X_1$; handle $X_2$; $\ldots$; handle $X_n$

    where handle $\texttt{t}$ : $\texttt{if hd toklis = t}$
                $\textbf{then}$ remove $\texttt{t}$ and continue
                $\textbf{else}$ error

    handle $B$ : $\texttt{parseB toklis}$

# "choose production based on hd toklis"

- **Need to formalize some things...**

- **Define "$\Rightarrow$":** $X_1 \ldots X_n \Rightarrow X_1 \ldots X_{i-1} \alpha X_{i+1} \ldots X_n$ **(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* $\alpha$ **if** $\vec{X} \Rightarrow^* \alpha$**.)**

- $\alpha$ **is a** *sentential form* **of G if the start symbol of G derives** $\alpha$**. If, furthermore,** $\alpha$ **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** $\epsilon$**-productions.)**

# "choose production based on hd toklis" (cont.)

- $\vec{X}$ **is** *nullable* **if it can derive** $\epsilon$**.**

- **Define:** **FIRST(**$\vec{X}$**)** $= \{t \in T \,|\, \vec{X} \Rightarrow^* t\alpha$ **for some** $\alpha\}$ $\cup$ $\{\bullet \,|\, \vec{X}\ nullable\}$**.**

- **Define: FOLLOW(**$A$**)** $= \{t \in T \,|\, \exists$ **a sentential form** $\alpha A t \beta\}$

    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.

# "choose production based on hd toklis" (cont.)

- **Define: G is *left-recursive* if $\exists A : A \Rightarrow^+ A\alpha$ for some $\alpha$.**

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

# Top-down parsing revisited

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

$$A \to \vec{X} \mid \vec{Y} \mid \ldots \mid \vec{Z}$$

construct `parseA`:

```
parseA toklis = let t = hd toklis in
     if t ∈ FIRST(X⃗) then handle X⃗
     else if t ∈ FIRST(Y⃗) then handle Y⃗
     ...else if t ∈ FIRST(Z⃗) or (• ∈ FIRST(Z⃗) and t ∈ FOLLOW(A))
            then handle Z⃗ (Z⃗ the unique nullable right-hand side of A, if any)
            else error
```

handle $X_1, X_2, \ldots, X_n$ : handle $X_1$; handle $X_2$; $\ldots$; 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:**

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

$$\Rightarrow \quad A \rightarrow \alpha B$$
$$B \rightarrow \beta \mid \gamma$$

- **Left-recursion removal:**

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

$$\Rightarrow \quad A \rightarrow \beta B$$
$$B \rightarrow \epsilon \mid \alpha B$$

# Example

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

**A → id | '(' B ')'**
**B → A | A '+' B**

● **Grammar 3 transformed to LL(1) form:**

**A → id | '(' B ')'**
**B → A C**
**C → '+' B | $\epsilon$**

# 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

# Enforcing associativity

# Some expression grammars

$G_A$: E → id | E - E | E * E

$G_B$: E → id | id - E | id * E

$G_C$: E → id | E - id | E * id

$G_D$:  E → T - E | T
T → id | id * T

$G_E$:  E → E - T | T
T → id | T * id

$G_F$:  E → T E′
E′ → ε | - E
T → id T′
T′ → ε | * T

● $G_A$: **E → id | E - E | E * E**

● **x - y * z**                                      **x - y - z**

- $G_B$: **E → id | id - E | id * E**

- **x - y * z**                                    **x - y * z - w**

- **x * y - z**

- $G_C$: **E $\rightarrow$ id | E - id | E * id**

- **x - y * z**                                        **x - y * z - w**

- **x * y - z**

● $G_D$: **E → T - E | T**
**T → id | id \* T**

● **x - y \* z**                    **x \* y - z**                    **x - y - z**

● $G_E$: **E → E - T | T**
    **T → id | T * id**

● **x - y * z**                    **x * y - z**                    **x - y - z**

● $G_F$: **E → T E′**

         **E′ → $\epsilon$ | - E**

         **T → id T′**

         **T′ → $\epsilon$ | * T**

● **x - y * z**                          **x * y - z**                          **x - y - z**