# CS 536 Announcements for Wednesday, March 6, 2024
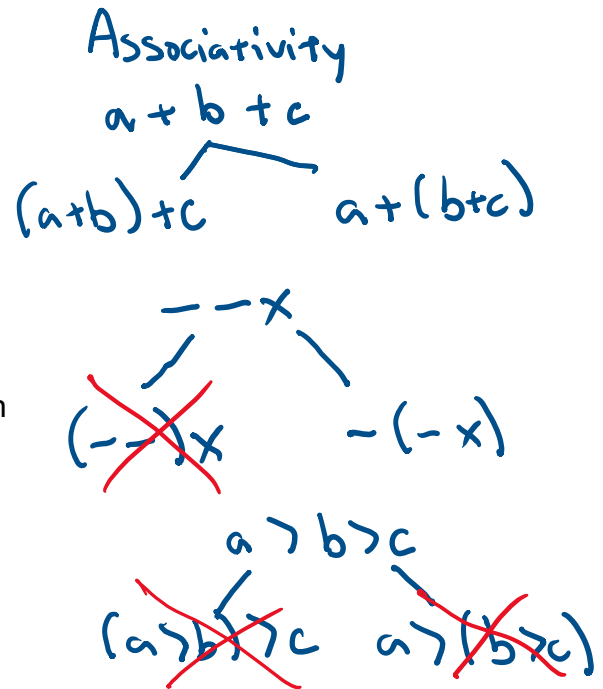
**Last Time**
- wrap up CYK
- classes of grammars
- top-down parsing

**Today**
- review grammar transformations
- building a predictive parser
- FIRST and FOLLOW sets

**Next Time**
- predictive parsing and syntax-directed translation

Associativity

$$a + b + c$$

$$(a+b)+c \qquad a+(b+c)$$

$$--x$$

$$(--x)x \qquad \sim(-x)$$

$$a > b > c$$

$$(a>b)>c \qquad a>(b>c)$$

# LL(1) Predictive Parser

**Predict the parse tree top-down**

**Parser structure**

- 1 token lookahead

- parse/selector table

- stack tracking current parse tree's frontier

**Necessary conditions**

- left-factored

- free of left-recursion

# Review of LL(1) grammar transformations

**Necessary (but not sufficient conditions) for LL(1) parsing**

- free of left recursion – no left-recursive rules
- left-factored – no rules with a common prefix, for any nonterminal
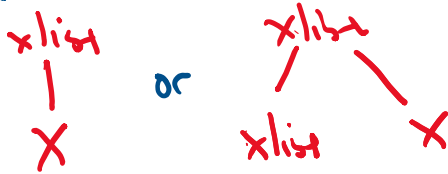
**Why left-recursion is a problem**

<u>Outside/high-level view</u>

CFG snippet: xlist → xlist X | X

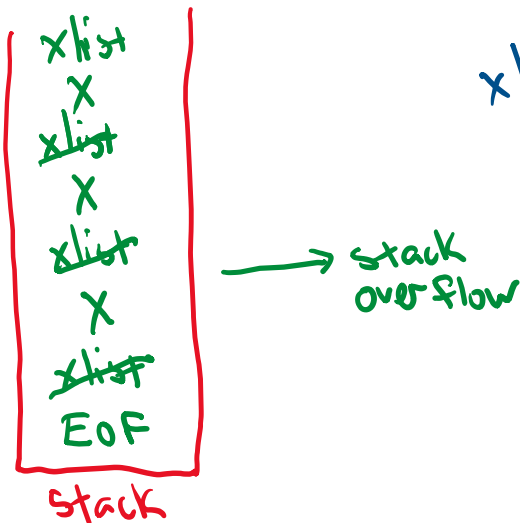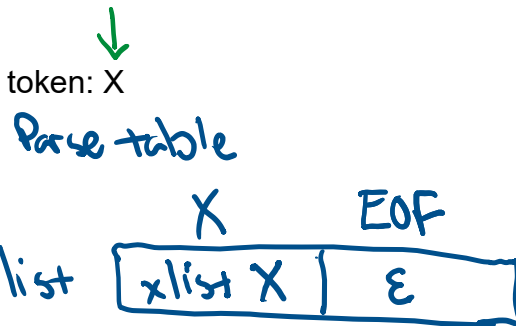Current parse tree:  [ xlist ]          Current token: X

How to grow parse tree?

xlist
  |
  X

or

xlist
 /  \
xlist  X

Depends on if there
are more X's
→ need more lookahead

<u>Inside/algorithmic-level view</u>

CFG snippet: xlist → xlist X | ~~X~~ ε

Current parse tree:  [ xlist ]          Current token: X

Parse table

|  | X | EOF |
|---|---|---|
| xlist | xlist X | ε |

xlist
X
xlist
X
xlist
X
xlist
EOF

→ stack overflow

stack

# Removing left-recursion (review)

Replace

$$A \rightarrow A\,\alpha \mid \boxed{\beta} \quad \leftarrow \text{head of list}$$

with

$$A \rightarrow \boxed{\beta}\, A'$$
$$A' \rightarrow \alpha\, A' \mid \varepsilon$$

where β does not start with A (or may be ε)



Preserves the language (as a list of α's, starting with a β), but uses ==right recursion==

**Example**

$$\text{xlist} \rightarrow \text{xlist } X \mid \boxed{\varepsilon}$$

==$\text{xlist} \rightarrow \varepsilon\ \text{xlist}'$==

$\text{xlist}' \rightarrow X\ \text{xlist}' \mid \varepsilon$

$\underline{\text{xlist} \rightarrow \text{xlist}'}$

remove
&
change $\text{xlist}'$
to $\text{xlist}$

$\Rightarrow \text{xlist} \rightarrow X\ \text{xlist} \mid \varepsilon$

# Left factoring (review)

**Removing a common prefix from a grammar**

Replace

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid ... \mid \alpha\beta_n \mid \gamma_1 \mid \gamma_2 \mid ... \mid \gamma_m$$

with

$$A \rightarrow \alpha A' \mid \gamma_1 \mid \gamma_2 \mid ... \mid \gamma_m$$
$$A' \rightarrow \beta_1 \mid \beta_2 \mid ... \mid \beta_n$$

where $\beta_i$ and $\gamma_i$ are sequence of symbols with no common prefix

Note: $\gamma_i$ may not be present, and one of the $\beta_i$ may be $\varepsilon$

**Idea**: combine all "problematic" rules that start with $\alpha$ into one rule $\alpha A'$
  $A'$ now represents the suffix of the problematic rules

## Example 1

exp $\rightarrow$ < A > | < B > | < C > | D

$$exp \rightarrow\ < exp' \mid D$$
$$exp' \rightarrow A> \mid B> \mid C>$$

## Example 2

stmt $\rightarrow$ ID ASSIGN exp | ID ( elist ) | return
exp $\rightarrow$ INTLIT | ID
elist $\rightarrow$ exp | exp COMMA elist

$$stmt \rightarrow ID\ stmt' \mid return$$
$$stmt' \rightarrow ASSIGN\ exp \mid (elist)$$
$$exp \rightarrow INTLIT \mid ID$$
$$elist \rightarrow exp\ elist'$$
$$elist' \rightarrow \varepsilon \mid COMMA\ elist$$

# Building the parse table

**Goal**: given production *lhs* → *rhs*, determine what terminals would lead us to choose that production

*ie, figure out T such that table[lhs][T] = rhs*

*— also what terminals could indicate an error at this point?*

- what terminals could *rhs possibly start with*?
- What terminals could possibly come after *lhs*?

**Idea:** FIRST(*rhs*) = set of terminals that begin sequences derivable from *rhs*

Suppose top-of-stack symbol is nonterminal $p$ and the current token is **A** and we have
- Production 1: $p \rightarrow \alpha$
- Production 2: $p \rightarrow \beta$

FIRST lets us disambiguate:

- if **A** ∈ FIRST(α), then *production 1 is a viable choice*

- if **A** ∈ FIRST(β), then *production 2 is a viable choice*

- if **A** is in just one of them, then *we can predict which production to use*

# FIRST sets

FIRST(α) is the set of terminals that begin the strings derivable from α, and also, if α can derive ε, then ε is in FIRST(α).

Formally,

*terminals*

$$\text{FIRST}(\alpha) = \{ T \mid (T \in \Sigma \wedge \alpha \Rightarrow^* T\beta) \vee (T = \varepsilon \wedge \alpha \Rightarrow^* \varepsilon) \}$$

**For a symbol X**

- if X is terminal: FIRST(X) = {X}

- if X is ε : FIRST(X) = {ε}

- if X is nonterminal : for each production X → $Y_1 Y_2 Y_3 .. Y_n$
  - put FIRST($Y_1$) − ε into FIRST(X)
  - if ε is in FIRST($Y_1$), put FIRST($Y_2$) − ε into FIRST(X)
  - if ε is in FIRST($Y_2$), put FIRST($Y_3$) − ε into FIRST(X)
  - ...
  - if ε is in FIRST($Y_i$) for all i, put ε into FIRST(X)

*repeat until there are no changes in any nonterminal's FIRST set*

Original CFG
expr →    expr + term
        | term
term →    term * factor
        | factor
factor →    exponent ^ factor
        | exponent
exponent → INTLIT
        | ( expr )

Transformed CFG
expr → term expr'
expr' → + term expr' | ε
term → factor term'
term' → * factor term' | ε
factor → exponent factor'
factor' → ^ factor | ε
exponent → INTLIT | ( expr )

| | FIRST | FOLLOW |
|---|---|---|
| expr | INTLIT ( | EOF ) |
| expr' | + ε | = FOLLOW(expr)  EOF ) |
| term | INTLIT ( | + EOF ) |
| term' | * ε | = FOLLOW(term)  + EOF ) |
| factor | INTLIT ( | * + EOF ) |
| factor' | ∧ ε | * + EOF ) |
| exponent | INTLIT ( | ∧ * + EOF ) |

| | | FIRST |
|---|---|---|
| expr | → term expr' | INTLIT ( |
| expr' | → + term expr' | + |
| expr' | → ε | ε |
| term | → factor term' | INTLIT ( |
| term' | → * factor term' | * |
| term' | → ε | ε |
| factor | → exponent factor' | INTLIT ( |
| factor' | → ^ factor | ∧ |
| factor' | → ε | ε |
| exponent | → INTLIT | INTLIT |
| exponent | → ( expr ) | ( |

# Computing FIRST(α) (continued)

**Extend FIRST to strings of symbols α**

— want to define FIRST for all RHS of productions

Let $\alpha = Y_1 Y_2 Y_3 .. Y_n$
- put FIRST($Y_1$) – ε into FIRST(α)
    - if ε is in FIRST($Y_1$), put FIRST($Y_2$) – ε into FIRST(α)
    - if ε is in FIRST($Y_2$), put FIRST($Y_3$) – ε into FIRST(α)
    - ...
    - if ε is in FIRST($Y_i$) for all i, put ε into FIRST(α)

Given two productions for nonterminal $p$
- Production 1: $p \rightarrow \alpha$    FIRST(α) ← ⎤
- Production 2: $p \rightarrow \beta$    FIRST(β) ← ⎦ — look for current token

If only 1 has it, pick that production
If both have it, grammar is <u>not</u> LL(1)
If neither have it, if one FIRST set has ε in it,
     look at what terminals can <u>follow</u> $p$

# FOLLOW sets

For single nonterminal $a$, FOLLOW($a$) is the set of terminals that can appear immediately to the right of $a$

Formally,

terminals

$$\textbf{FOLLOW}(a) = \{ T \mid (T \in \Sigma \wedge s \Rightarrow^* \alpha a T \beta) \vee (T = EOF \wedge s \Rightarrow^* \alpha a) \}$$
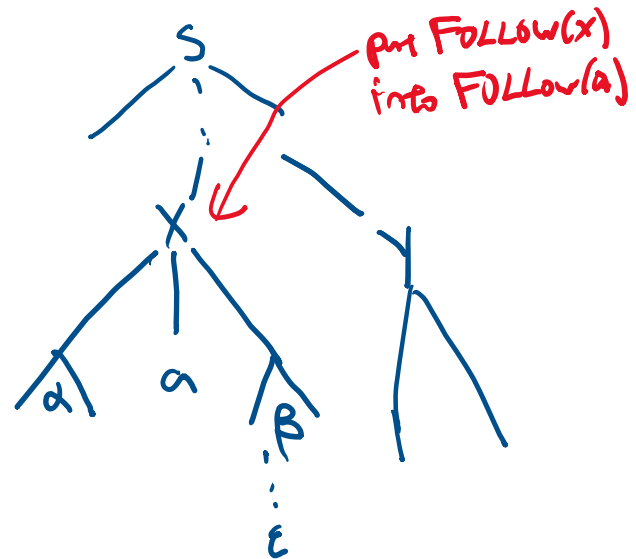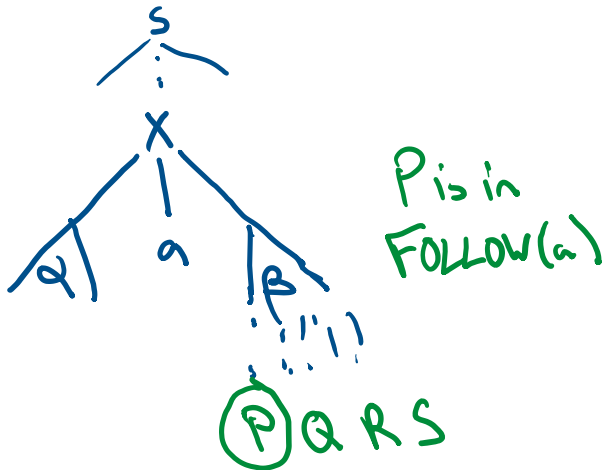
# Computing FOLLOW sets

**To build FOLLOW(a)**
- if a is the <u>start non-term</u>, put <u>EOF in FOLLOW(a)</u>
- for each production x → α a β
  - put FIRST(β) – ε into FOLLOW(a)
  - if ε is in FIRST(β), put FOLLOW(x) into FOLLOW(a)
- for each production x → α a
  - put FOLLOW(x) into FOLLOW(a)

*repeat until no changes*

*P is in FOLLOW(a)*

*P Q R S*

*put FOLLOW(x) into FOLLOW(a)*

# Building the parse table

```
for each production x → α {

    for each terminal T in FIRST(α) {
        put α in table[x][T]
    }

    if ε is in FIRST(α) {

        for each terminal T in FOLLOW(x) {
            put α in table[x][T]
        }
    }
}
```