# CS415 Compilers

# Syntax Analysis
# Part 2

These slides are based on slides copyrighted by
Keith Cooper, Ken Kennedy & Linda Torczon at Rice
University

- Midterm has been grades. Exams will be handed out in recitation on Wednesday. Please see canvas for grades.

- Third homework has been posted. NEW deadline:
  Due Thursday, March 10

- First project (local instruction scheduler) deadlines:
    code: March 9 @ 11:59pm – single tar file
    report: March 11 @ 11:59pm – single pdf file
  Late policy:
    → Grace period: 1 hour
    → 20% penalty for every started 24 hour period after the deadline.
    → Saturday/Sunday count as a single 24 hour period.
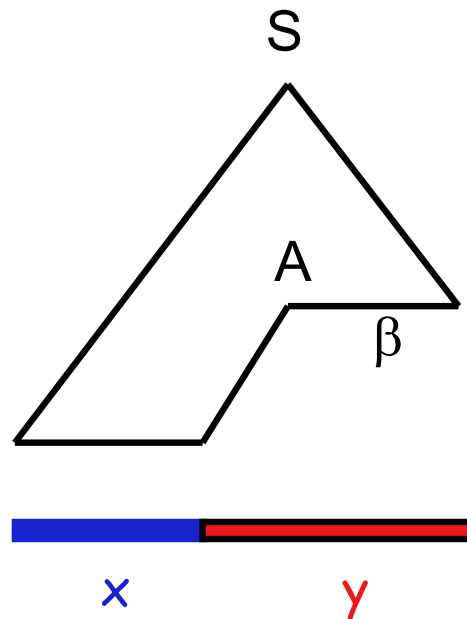
- Warning grades will be submitted by Friday, March 11

# Parsing
# (Syntax Analysis)

Top-Down Parsing

EAC Chapters 3.3

*LL(1), recursive descent*

1 input symbol lookahead

construct leftmost deriviation (forwards)

input: read left-to-right

$$S \Rightarrow^*_{lm} x A \beta \Rightarrow_{lm} x \delta \beta \Rightarrow^*_{lm} x y$$
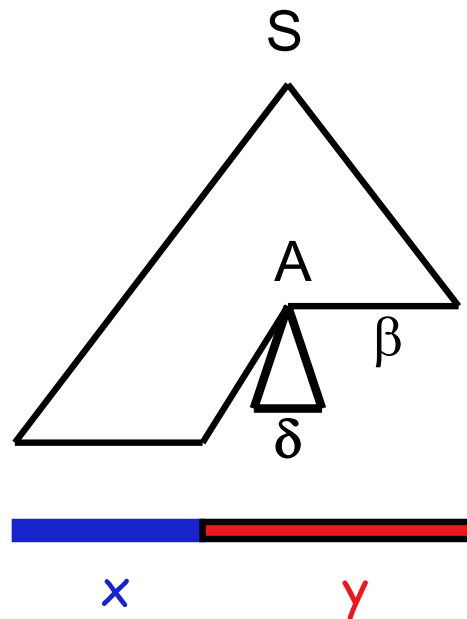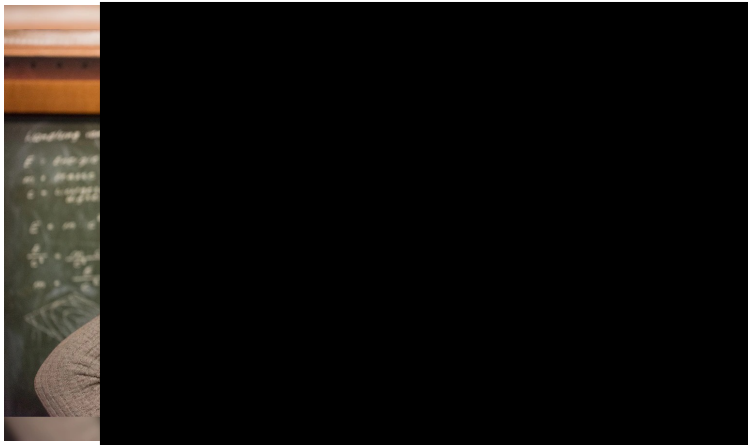


x        y

RUTGERS

*LL(1), recursive descent*

1 input symbol lookahead

construct leftmost deriviation (forwards)

input: read left-to-right

rule $A \rightarrow \delta$

$$S \Rightarrow^*_{lm} x \, A \, \beta \Rightarrow_{lm} x \, \delta \, \beta \Rightarrow^*_{lm} x \, y$$

S

A

$\beta$

$\delta$

x          y

*Scientist* → *Richard_Feynman | Albert_Einstein*

**Top-down**  This is what you see on the input before you make your rule decision:

**How much lookahead do you need?**

*Are we looking at either Richard Feynman or Albert Einstein?*

*Scientist* → *Richard_Feynman* | *Albert_Einstein*

**Top-down**   This is what you see on the input before you make your rule decision:



**How much lookahead do you need?**

*Are we looking at either Richard Feynman or Albert Einstein?*

**Scientist** → *Richard_Feynman | Albert_Einstein*

**Top-down**  This is what you see on the input before you make your rule decision:
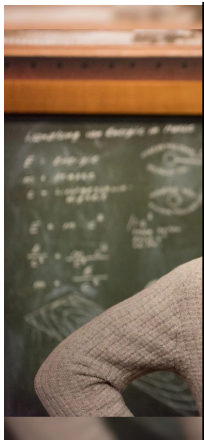


How much lookahead do you need?

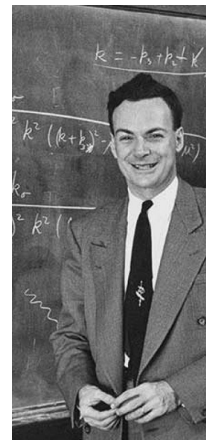*Are we looking at either Richard Feynman or Albert Einstein?*

Version with precedence

| 1 | Goal | → | Expr |
|---|---|---|---|
| 2 | Expr | → | Expr + Term |
| 3 | | | | Expr – Term |
| 4 | | | | Term |
| 5 | Term | → | Term * Factor |
| 6 | | | | Term / Factor |
| 7 | | | | Factor |
| 8 | Factor | → | number |
| 9 | | | | id |

*And the input x – 2 * y*

*Top-down parsers cannot handle left-recursive grammars*

Formally,

A grammar is *left recursive* if $\exists\ A \in NT$ such that

$\exists$ a derivation $A \Rightarrow^+ A\alpha$, for some string $\alpha \in (NT \cup T)^+$

Our expression grammar is left recursive

- This can lead to non-termination in a top-down parser
- For a top-down parser, any recursion must be right recursion
- We would like to convert the left recursion to right recursion

*Non-termination is a bad property in any part of a compiler*

To remove left recursion, we can transform the grammar

Consider a grammar fragment of the form

$$Fee \rightarrow Fee \; \alpha$$
$$| \; \beta$$

where neither $\alpha$ nor $\beta$ start with *Fee*

We can rewrite this as

$$Fee \rightarrow \beta \; Fie$$
$$Fie \rightarrow \alpha \; Fie$$
$$| \; \varepsilon$$

where *Fie* is a new non-terminal

*This accepts the same language, but uses only right recursion*

The expression grammar contains two cases of left recursion

| | | | | | | |
|---|---|---|---|---|---|---|
| *Expr* | → | *Expr + Term* | | *Term* | → | *Term * Factor* |
| | \| | *Expr – Term* | | | \| | *Term / Factor* |
| | \| | *Term* | | | \| | *Factor* |

Applying the transformation yields

| | | | | | | |
|---|---|---|---|---|---|---|
| *Expr* | → | *Term Expr'* | | *Term* | → | *Factor Term'* |
| *Expr'* | \| | + *Term Expr'* | | *Term'* | \| | * *Factor Term'* |
| | \| | – *Term Expr'* | | | \| | / *Factor Term'* |
| | \| | ε | | | \| | ε |

These fragments use only right recursion

Substituting them back into the grammar yields

| 1 | Goal | → | Expr |
|---|---|---|---|
| 2 | Expr | → | Term Expr' |
| 3 | Expr' | → | + Term Expr' |
| 4 | | \| | – Term Expr' |
| 5 | | \| | ε |
| 6 | Term | → | Factor Term' |
| 7 | Term' | → | * Factor Term' |
| 8 | | \| | / Factor Term' |
| 9 | | \| | ε |
| 10 | Factor | → | <u>number</u> |
| 11 | | \| | <u>id</u> |
| 12 | | \| | ( Expr ) |

- This grammar is correct, if somewhat non-intuitive.

- A top-down parser will terminate using it.

- A top-down parser may need to backtrack with it.

- General left recursion removal algorithm in EAC

*We set out to study parsing*

- Specifying syntax
  → Context-free grammars
  → Ambiguity

- Top-down parsers
  → Algorithm & its problem with left recursion
  → Left-recursion removal
  → Left factoring (will discuss later)

- Predictive top-down parsing
  → The LL(1) condition
  → Table-driven LL(1) parsers
  → Recursive descent parsers
    ▪ Syntax directed translation (example)

*If it picks the wrong production, a top-down parser may backtrack*

*Alternative is to look ahead in input & use context to pick correctly*

## How much lookahead is needed?

- In general, an arbitrarily large amount
- Use the Cocke-Younger, Kasami algorithm or Earley's algorithm

## Fortunately,

- Large subclasses of CFGs can be parsed with limited lookahead
- Most programming language constructs fall in those subclasses

Among the interesting subclasses are *LL(1)* and *LR(1)* grammars

## Basic idea

*Given $A \to \alpha \mid \beta$, the parser should be able to choose between $\alpha$ & $\beta$*

## FIRST sets

For some *rhs* $\alpha \in G$, define FIRST($\alpha$) as the set of tokens that appear as the first symbol in some string that derives from $\alpha$

That is, $a \in$ FIRST($\alpha$) *iff* $\alpha \Rightarrow^* a\, \gamma$, for some $\gamma$

## The LL(1) Property

If $A \to \alpha$ and $A \to \beta$ both appear in the grammar, we would like

$$\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \varnothing$$

This would allow the parser to make a correct choice with a lookahead of exactly one symbol !

This is almost correct, but not quite

$$a \in \text{FIRST}_1(\alpha) \text{ iff } \alpha \Rightarrow^* a\,\gamma, \text{ for some } \gamma$$

To build FIRST(X) for all grammar symbols X:
1.  if X is a terminal (token), FIRST(X) := { X }
2. if $X \rightarrow \varepsilon$, then $\varepsilon \in$ FIRST(X)

3. <u>iterate until</u> no more terminals or $\varepsilon$ can be added
                to any FIRST(X):
   if  $X \rightarrow Y_1 Y_2 \dots Y_k$ then
       $a \in$ FIRST(X) if  $a \in$ FIRST($Y_i$) and
                      $\varepsilon \in$ FIRST($Y_j$) for all $1 \leq j < i$
       $\varepsilon \in$ FIRST(X) if $\varepsilon \in$ FIRST($Y_i$) for all $1 \leq i \leq k$
   <u>end iterate</u>

Note: if $\varepsilon \notin$ FIRST($Y_1$), then FIRST($Y_i$) is irrelevant, for $1 < i$

$$a \in \text{FIRST}(\alpha) \; iff \; \alpha \Rightarrow^* \underline{a} \, \gamma, \; \text{for some } \gamma$$

To build **FIRST($\alpha$)** for $\alpha = X_1 X_2 \ldots X_n$ :

1. $a \in \text{FIRST}(\alpha)$ if $a \in \text{FIRST}(X_i)$ and
$$\varepsilon \in \text{FIRST}(X_j) \text{ for all } 1 \leq j < i$$

2. $\varepsilon \in \text{FIRST}(\alpha)$ if $\varepsilon \in \text{FIRST}(X_i)$ for all $1 \leq i \leq n$

For a non-terminal A, define FOLLOW(A) as

> *FOLLOW(A) := the set of terminals that can appear immediately to the right of A in some sentential form.*

Thus, a non-terminal's FOLLOW set specifies the tokens that can legally appear after it; a terminal has no FOLLOW set

$$\text{FOLLOW(A)} = \{\ a \in (T \cup \{eof\})\ |\ S\ eof \Rightarrow^* \alpha\ A\ a\ \gamma\ \}$$

To build FOLLOW(X) for all non-terminal X:

1. Place eof in FOLLOW( *<goal>* )

iterate <u>until</u> no more terminals or eof can be added
to any FOLLOW(X):

2. If $A \rightarrow \alpha B \beta$ then
put {FIRST($\beta$) - $\varepsilon$} in FOLLOW(B)

3. If $A \rightarrow \alpha B$ then
put FOLLOW(A) in FOLLOW(B)

4. If $A \rightarrow \alpha B \beta$ and $\varepsilon \in$ FIRST($\beta$) then
put FOLLOW(A) in FOLLOW(B)

If $A \to \alpha$ and $A \to \beta$ and $\varepsilon \in \textsc{First}(\alpha)$, then we need to ensure that $\textsc{First}(\beta)$ is disjoint from $\textsc{Follow}(A)$, too

Define $\textsc{First}^+(\delta)$ for rule $A \to \delta$ as

- $(\textsc{First}(\delta) - \{ \varepsilon \}) \cup \textsc{Follow}(A)$, if $\varepsilon \in \textsc{First}(\delta)$
- $\textsc{First}(\delta)$, otherwise

The LL(1) Property

A grammar is *LL(1)* iff $A \to \alpha$ and $A \to \beta$ implies
$$\text{FIRST}^+(\alpha) \cap \text{FIRST}^+(\beta) = \varnothing$$

This would allow the parser to make a correct choice with a lookahead of exactly one symbol !

<u>Question</u>: Can there be two rules $A \to \alpha$ and $A \to \beta$ in a *LL(1)* grammar such that $\varepsilon \in \text{FIRST}(\alpha)$ and $\varepsilon \in \text{FIRST}(\beta)$?

Given a grammar that has the *LL(1)* property

- Problem: NT A needs to be replaced in next derivation step
- Assume $A \rightarrow \beta_1 \mid \beta_2 \mid \beta_3$, with
  $\text{FIRST}^+(\beta_1) \cap \text{FIRST}^+(\beta_2) = \varnothing$, $\text{FIRST}^+(\beta_1) \cap \text{FIRST}^+(\beta_3) = \varnothing$, and
  $\text{FIRST}^+(\beta_2) \cap \text{FIRST}^+(\beta_3) = \varnothing$ (pair-wise disjoint sets)

```
/* find rule for A */
if (current token ∈ FIRST+ (β₁))

    select A → β₁
else if (current token ∈ FIRST+(β₂))

    select A → β₂
else if (current token ∈ FIRST+(β₃))

    select A → β₃
else
    report an error and return false
```

Grammars with the *LL(1)* property are called *predictive grammars* because the parser can "predict" the correct expansion at each point in the parse.

Parsers that capitalize on the *LL(1)* property are called *predictive parsers*.

One kind of predictive parser is the *recursive descent* parser. *The other is a table-driven parser* *table-driven parser.*

**More Syntax Analysis**

Top-down: Read EaC: Chapter 3.3

Bottom-up: Read EaC: Chapter 3.4