



Top-Down Parsing

Dragon ch. 4.4

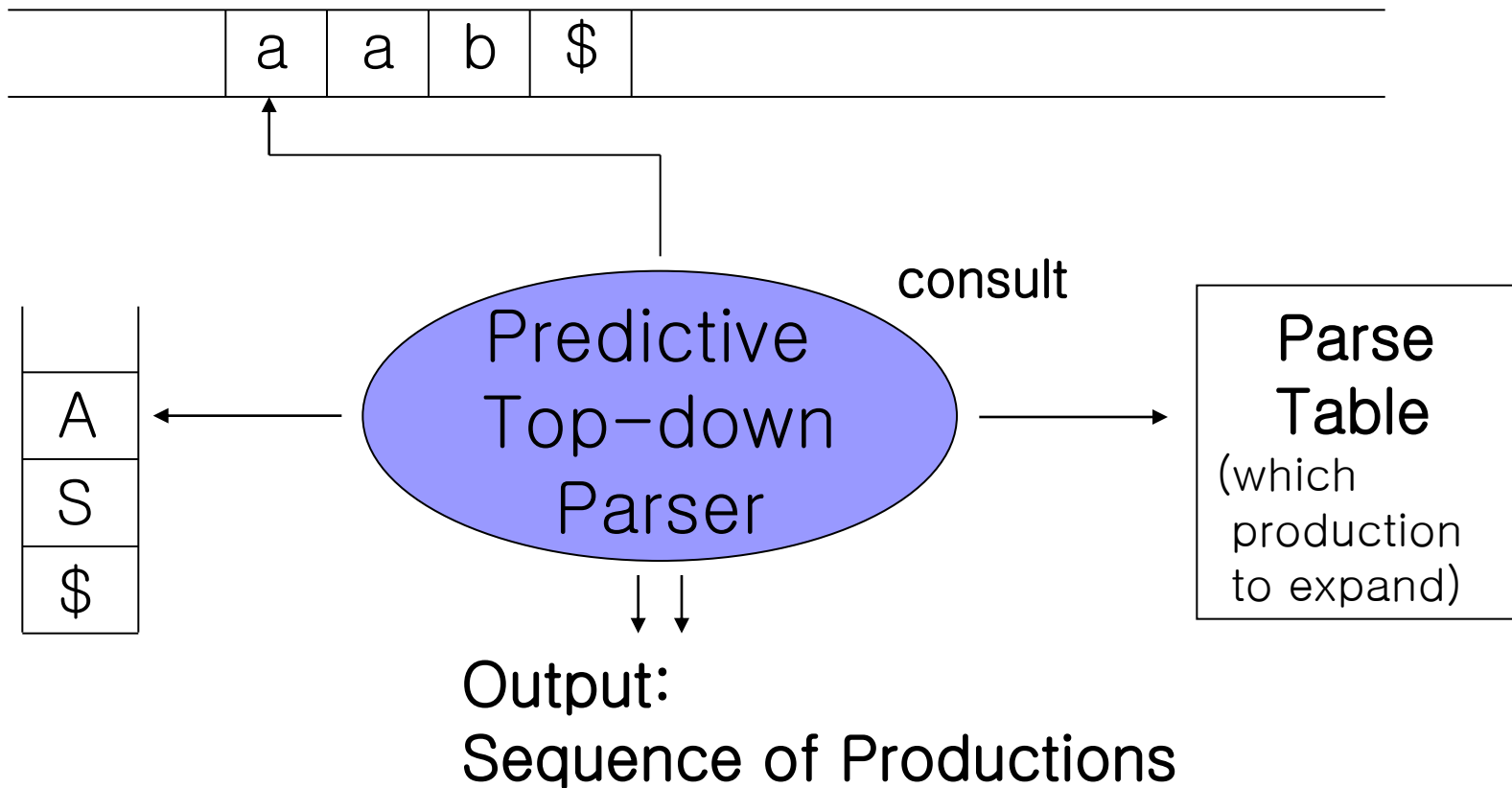
Recognizer and Parser

- A **recognizer** is a machine (system) that can **accept** a terminal string for some grammar and determine whether the string is in the language accepted by the grammar
- A **parser**, in addition, finds a **derivation** for the string
 - For grammar G and string x , find a derivation $S \Rightarrow^* x$ if one exists
 - Construct a parse tree corresponding to this derivation
 - Input is read (scanned) from left to right
 - Two types of the parser: **top-down** vs. **bottom-up**

Top-down Parsing

- Top-down parsing **expands a tree** from the **top** (start symbol) using a **stack**
 - Put the start symbol on the stack top
 - Repeat
 - **Expand a nonterminal** on the stack top
 - Match stack tops with input terminal symbols
- Problem: which production to expand?
 - If there are multiple productions for a given nonterminal
 - One way: **guess!**

Structure of top-down parsing



Example of Parsing by Guessing

- P of an Example Grammar

- $S \rightarrow AS|B, A \rightarrow a, B \rightarrow b$

- Parsing process

input	stacktop	action
a a b \$	S \$	S \rightarrow AS
a a b \$	A S \$	A \rightarrow a
a a b \$	a S \$	Match a
a b \$	S \$	S \rightarrow AS
a b \$	A S \$	A \rightarrow a
a b \$	a S \$	Match a
b \$	S \$	S \rightarrow B
b \$	B \$	B \rightarrow b
b \$	b \$	Match b
\$	\$	End

- In reality, computers do not guess very well

- So we use **lookahead** for correct expansion

- Before we do this, we must “condition” the grammar

Removal of Left Recursion

- Problem: infinite regression
 - $A \rightarrow A\alpha \mid \beta$, (the corresponding language is $\beta\alpha^*$)
 - Remove of immediate left recursion
 $A \rightarrow \beta B, B \rightarrow \alpha B \mid \epsilon$
- More generally,
 - $A \rightarrow A\alpha_1, A \rightarrow A\alpha_2$
 - $A \rightarrow \beta_1, A \rightarrow \beta_2$
 - $A \rightarrow (\beta_1 \mid \beta_2) B, B \rightarrow (\alpha_1 \mid \alpha_2) B \mid \epsilon$

Example of Removing Left Recursion

- Example of removing left immediate recursion

$$E \rightarrow E + T$$
$$E \rightarrow T$$

$$E \rightarrow TB$$
$$B \rightarrow +TB \mid e$$

- Can remove all left recursions
- Refer to Dragon Ch. 4.1 page 177

Left Factoring

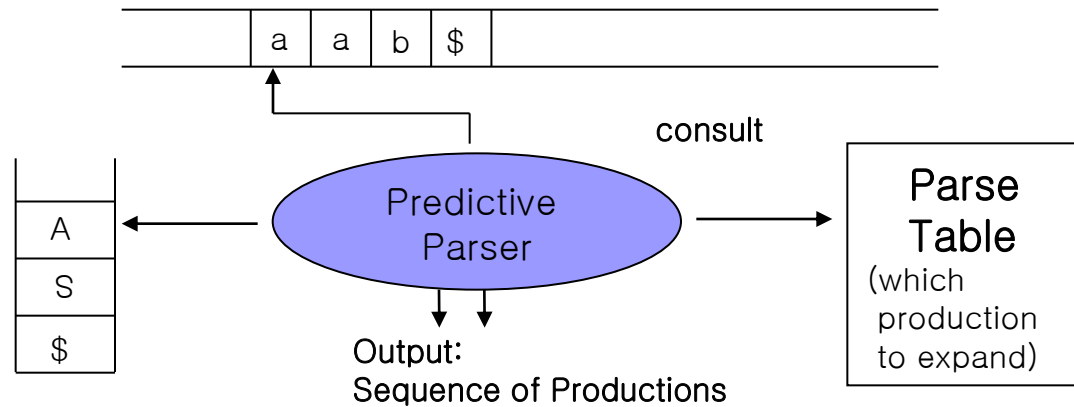
- Not have sufficient information right now
 - $A \rightarrow \alpha\beta \mid \alpha\gamma$
- Left factoring: turn two alternatives into one so that we match α first and hope it helps
 - $A \rightarrow \alpha B, B \rightarrow \beta \mid \gamma$
 - Example:

$$\begin{array}{l} E \rightarrow T + E \\ E \rightarrow T \\ \hline E \rightarrow TB \\ B \rightarrow +E \mid e \end{array}$$

Predictive Top-Down Parsing

- Perform educated guess
 - Do not blindly guess productions that cannot even get the first symbol right
 - If the current input symbol is a and the top stack symbol is S , which of the two productions ($S \rightarrow bS$, $S \rightarrow a$) should be expanded?
- Two versions
 - Non-Recursive version with a stack
 - Recursive version: recursive descent parsing

Table-Driven Non-Recursive Parsing



- Input buffer: the string to be parsed followed by \$
- Stack: a sequence of grammar symbols with \$ at the bottom
 - Initially, the stack has the start symbol on top of \$
- Parsing table: two dimensional array $M[A, a]$, where A is a non-terminal and a is a terminal or \$; it has productions
- Output: a sequence of productions expanded

Action of the Parser

When X is a symbol on top of the stack and a is the current input symbol

- If $X = a = \$$, a successful completion of parsing
- If $X = a \neq \$$, pops X off the stack and advances the input pointer to the next input symbol
- If X is a nonterminal, consult $M[X, a]$ which will be either an X -production or an error;
 - If $M[X, a] = \{X \rightarrow UVW\}$, X on top of stack is replaced by WVU (with U on top) and print its production number
 - If $[X, a] = \text{error}$ means a parsing error

An Example Grammar

Original grammar

$$E \rightarrow E + T \mid T$$
$$T \rightarrow T * F \mid F$$
$$F \rightarrow (E) \mid \text{id}$$

After removing left recursion

$$E \rightarrow TE'$$
$$E' \rightarrow +TE' \mid e$$
$$T \rightarrow FT'$$
$$T' \rightarrow *FT' \mid e$$
$$F \rightarrow (E) \mid \text{id}$$

An Example Parsing Table

	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow e$	$E' \rightarrow e$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow e$	$T' \rightarrow *F$		$T' \rightarrow e$	$T' \rightarrow e$
F	$F \rightarrow id$			$F \rightarrow (E)$		

How is **id + id * id** parsed?

How to Construct the Parse Table?

- For this, we use three functions
 - *Nullable()*: can it be a **null**?
 - Predicate, $V^* \rightarrow \{\text{true}, \text{false}\}$
 - Telling if a string of nonterminals is nullable, i.e., can derive an empty string
 - *FNE()*: **first** but **no epsilon**
 - Terminals that can appear at the beginning of a derivation from a string of grammar symbols
 - *Follow()*: what can **follow** after a nonterminal?
 - Terminals (or \$) that can appear after a nonterminal in some sentential form

Nullable()

- **Nullable(α)** = true if $\alpha \Rightarrow * \epsilon$
= false, otherwise
 - Start with the obvious ones, e.g., $A \rightarrow \epsilon$
 - Add new ones when $A \rightarrow \alpha$ and Nullable(α)
 - Keep going until there is no change
- More formally,
 - Nullable(ϵ) = true
 - Nullable($X_1 X_2 \dots X_n$) = true iff Nullable(X_i) $\forall i$
 - Nullable(A) = true if $A \rightarrow \alpha$ and Nullable(α)

FNE()

- Definition: $FNE(\alpha) = \{a \mid \alpha \Rightarrow^* aX\}$
- FNE() is computed as in Nullable()
 - $FNE(a) = \{a\}$
 - $FNE(X_1X_2\dots X_n) =$
if(!Nullable(X_1)) then $FNE(X_1)$
else $FNE(X_1) \cup FNE(X_2X_3\dots X_n)$
 - if $A \rightarrow \alpha$ then $FNE(A) \supseteq FNE(\alpha)$

FNE() Computation Example

- For our example grammar

- $E \rightarrow TE'$
- $E' \rightarrow +TE' \mid e$
- $T \rightarrow FT'$
- $T' \rightarrow *FT' \mid e$
- $F \rightarrow (E) \mid id$

- We can compute FNE() as follows

Nullable(T) = false	FNE(E) = FNE(T) = {(,id}
	FNE(E') = {+}
Nullable(F) = false	FNE(T) = FNE(F) = {(,id}
	FNE(T') = {*}
	FNE(F) = {(, id}

First()

- The Dragon book uses First(), which is a combination of Nullable() and FNE()
 - If α is nullable $\text{First}(\alpha) = \{\mathbf{a} \mid \alpha \Rightarrow^* \mathbf{a}X\} \cup \{\epsilon\}$
else $\text{First}(\alpha) = \{\mathbf{a} \mid \alpha \Rightarrow^* \mathbf{a}X\}$
- First() can be computed from Nullable() and FNE(), or directly (see Dragon book)

Follow()

- $\text{Follow}(A) = \{a \mid S \Rightarrow^* \alpha A a \beta\}$, where a might be $\$$
 - $\text{Follow}()$ is needed if there is an ϵ -production
- To compute $\text{Follow}()$,
 - $\$ \in \text{Follow}(S)$
 - When $A \rightarrow \alpha B \beta$,
 $\text{Follow}(B) \supseteq \text{FNE}(\beta)$
 - When $A \rightarrow \alpha B \beta$ and $\text{Nullable}(\beta)$,
 $\text{Follow}(B) \supseteq \text{Follow}(A)$

Follow() Computation Example

- For our example grammar

- $E \rightarrow TE'$
- $E' \rightarrow +TE' \mid e$
- $T \rightarrow FT'$
- $T' \rightarrow *FT' \mid e$
- $F \rightarrow (E) \mid id$

- When $A \rightarrow \alpha B \beta$,
Follow(B) \supseteq FNE(β)
- When $A \rightarrow \alpha B \beta$ and Nullable(β),
Follow(B) \supseteq Follow(A)

- We can compute Follow() as follows

$FNE(E) = FNE(T) = \{ (, id \}$

$FNE(E') = \{ + \}$

$FNE(T) = FNE(F) = \{ (, id \}$

$FNE(T') = \{ * \}$

$FNE(F) = \{ (, id \}$

$Follow(E) = \{ \$,) \}$

$Follow(E') = \{ \$,) \}$

$Follow(T) = \{ +, \$,) \}$

$Follow(T') = \{ +, \$,) \}$

$Follow(F) = \{ *, +, \$,) \}$

Predictive Parsing Table

- How to construct the parsing table
 - Mapping $N \times T \rightarrow P$
 - $A \rightarrow \alpha \in M[A, a]$ for each $a \in FNE(\alpha Follow(A))$
 - $a \in FNE(\alpha)$, or
 - $Nullable(\alpha)$ and $a \in FOLLOW(A)$
- Meaning of “ $Nullable(\alpha)$ and $a \in FOLLOW(A)$ ”
 - Since the stack has (part of) a **sentential form** with A at the top, we can remove A (by expanding $A \rightarrow \alpha$) then try to match a with a symbol below A in the stack
 - Why? The symbol below A must be in $Follow(A)$, so there is a chance that it can be a (✖ or is this always guaranteed?)

Predictive Parsing Table

- For our example grammar
 - $E \rightarrow TE'$
 - $E' \rightarrow +TE' \mid e$
 - $T \rightarrow FT'$
 - $T' \rightarrow *FT' \mid e$
 - $F \rightarrow (E) \mid id$
- The parsing table is as follows:

	FNE()	Follow()	id	+	*	()	\$
E	(, id	\$,)	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'	+	\$,)		$E' \rightarrow +TE'$			$E' \rightarrow e$	$E' \rightarrow e$
T	(, id	+, \$,)	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'	*	+, \$,)		$T' \rightarrow e$	$T' \rightarrow *F$		$T' \rightarrow e$	$T' \rightarrow e$
F	(, id	*, +, \$,)	$F \rightarrow id$			$F \rightarrow (E)$		

LL(1) Grammar

- Definition: a grammar **G is LL(1)** if there is at most one production for any entry in the table
 - So we can do top-down parsing with one lookahead
- **LL(1)** means **left-to-right scan**, performing **leftmost** derivation, with **one** symbol lookahead

LL(1) Conditions

- **G is LL(1)** iff whenever $A \rightarrow \alpha$ and $A \rightarrow \beta$ are distinct productions of G , the following holds
 - α and β do not both derive strings beginning with a ($a \in \Sigma$)
 - α and β do not both derive ϵ
 - if $\beta \Rightarrow^* \epsilon$ then $FNE(\alpha) \cap \text{Follow}(A)$ is empty
- In other words, **G is LL(1)** if
 - if G is ϵ -free and unambiguous, $FNE(\alpha) \cap FNE(\beta) = \emptyset$
 - If an ϵ -production is present,
 $FNE(\alpha \text{Follow}(A)) \cap FNE(\beta \text{Follow}(A)) = \emptyset$

Testing for non-LL(1)ness

- In practice, for LL(1) testing, it is easiest to construct the parse table and check
- Some shortcuts to test if G is **not LL(1)**
 - G is left-recursive (e.g., $A \rightarrow A\alpha \mid \beta$)
 - Common left factors (e.g., $A \rightarrow \alpha\beta \mid \alpha\gamma$)
 - G is ambiguous (e.g., $S \rightarrow Aa \mid a, A \rightarrow \varepsilon$)

Non-LL(1) Grammar

- Consider the following grammar G1, which is not LL(1)

$S \rightarrow Bbc$

$B \rightarrow \epsilon | b | c$

- $FNE(B) = FNE(S) = \{b, c\}$,
- $FOLLOW(S) = \{\$ \}$, $FOLLOW(B) = \{b\}$

	FNE	FOLLOW	b	c
S	{b,c}	{\\$}	$S \rightarrow Bbc$	$S \rightarrow Bbc$
B	{b,c}	{b}	$B \rightarrow \epsilon$ $B \rightarrow b$	$B \rightarrow c$

- Since $FNE(\epsilon FOLLOW(B)) = FNE(b FOLLOW(B)) = \{b\}$
- We want consider a larger class of LL parsing, **LL(k)**, which look-ahead more symbols

LL(K) Parsing

- Begin by extending the definition of FNE() and FOLLOW()
 - Definitions of $FNE_k()$ and $FOLLOW_k()$

$$FNE_k(\alpha) = \{w \mid (|w| < k \text{ and } \alpha \Rightarrow^* w) \text{ or} \\ (|w| = k \text{ and } \alpha \Rightarrow^* wx \text{ for some } x)\}$$
$$FOLLOW_k(A) = \{w \mid S \Rightarrow^* \alpha A \beta \text{ and } w \in FNE_k(\beta)\}$$

- As with FOLLOW(), we will implicitly augment the grammar with $S' \rightarrow S\k so that our definitions are:
 $FOLLOW_k(a) = \{w \mid S \Rightarrow^* \alpha A \beta \text{ and } \omega \in FNE_k(\beta \$^k)\}$

LL(K) Parsing Definition

- G is LL(k) for some fixed k if, whenever there are two leftmost derivations,

$S \Rightarrow^* wA\alpha \Rightarrow w\beta\alpha \Rightarrow^* wx$, and
 $S \Rightarrow^* wA\alpha \Rightarrow w\gamma\alpha \Rightarrow^* wy$ and $\beta \neq \gamma$,
then $FNE_k(x) \neq FNE_k(y)$

Strong-LL(k) Parsing

- Simplest way to implementing LL(k) parsing table
 - Insert $A \rightarrow \alpha \in M[A, x]$ for each $x \in FNE_k(\alpha \text{Follow}_k(A))$
 - A grammar G is **strong-LL(k)** if there is at most one production for any entry in the table
 - If $FNE_k(\beta \text{FOLLOW}_k(A)) \cap FNE_k(\gamma \text{FOLLOW}_k(A)) = \Phi$ for all $A \rightarrow \beta$ and $A \rightarrow \gamma$ in G

Non-LL(1), but Strong-LL(2) Grammar

- Consider our non-LL(1) grammar G1 again

$S \rightarrow Bbc$

$B \rightarrow \epsilon | b | c$

- $FNE_2(BbcFOLLOW_2(S)) = \{bc, bb, cb\}$
- $FNE_2(\epsilon FOLLOW_2(B)) = \{bc\}$, $FNE_2(bFOLLOW_2(B)) = \{bb\}$,
 $FNE_2(cFOLLOW_2(B)) = \{cb\}$

	bc	bb	cb
S	$S \rightarrow Bbc$	$S \rightarrow Bbc$	$S \rightarrow Bbc$
B	$B \rightarrow \epsilon$	$B \rightarrow b$	$B \rightarrow c$

- So, G1 is Strong-LL(2) even though it is not LL(1)

LL(2) but Non-Strong LL(2) Grammar

- Consider the following grammar G2

$S \rightarrow Bbc|aBcb$

$B \rightarrow \epsilon|b|c$

- $FNE_2()$ and $FOLLOW_2()$ functions:

- $FNE_2(S) = \{ab, ac, bb, bc, cb\}$, $FNE_2(B) = \{b,c\}$

- $FOLLOW_2(S) = \{\$\$ \}$, $FOLLOW_2(B) = \{bc,cb\}$

- $FNE_2(\epsilon FOLLOW_2(B)) = \{bc,cb\}$

- $FNE_2(bFOLLOW_2(B)) = \{bb,bc\}$, so not strong-LL(2)

- But isn't G LL(2), either?

- Check with the LL(k) definition

- $S \Rightarrow Bbc \Rightarrow$

- $S \Rightarrow aBcb \Rightarrow$

Modified Grammar G2'

- G2 is indeed LL(2), then what's wrong with strong-LL(2) algorithm? Why can't it generate a LL(2) parsing table?
 - Because of **Follow()**, which does not always tell the truth
- Let us rewrite G2 with two new nonterminals, B_{bc} and B_{cb} , to keep track of local lookahead (context) information
 - $S \rightarrow B_{bc}bc | aB_{cb}cb$
 - $B_{bc} \rightarrow \epsilon | b | c$
 - $B_{cb} \rightarrow \epsilon | b | c$
- Now, in place of $FNE_2(\beta FOLLOW_2(B))$ to control putting $B \rightarrow \beta$ into table, use $FNE_2(\beta R)$ to control $B_R \rightarrow \beta$, where R is local lookahead

- For $S \rightarrow B_{bc}bc$, $FNE_2(B_{bc}bc\$\$) = \{bc,bb,cb\}$
- For $S \rightarrow aB_{cb}cb$, $FNE_2(aB_{cb}cb\$\$) = \{ac,ab\}$
- For $B_{bc} \rightarrow \epsilon$, $FNE_2(\epsilon\{bc\}) = \{bc\}$
- For $B_{bc} \rightarrow b$, $FNE_2(b\{bc\}) = \{bb\}$
- For $B_{bc} \rightarrow c$, $FNE_2(c\{bc\}) = \{cb\}$
- For $B_{cb} \rightarrow \epsilon$, $FNE_2(\epsilon\{cb\}) = \{cb\}$
- For $B_{cb} \rightarrow b$, $FNE_2(b\{cb\}) = \{bc\}$
- For $B_{cb} \rightarrow c$, $FNE_2(c\{cb\}) = \{cc\}$

■ Corresponding LL(2) Table: G2' is strong-LL(2)

	bc	bb	cb	ab	ac	cc
S	$S \rightarrow B_{bc}bc$	$S \rightarrow B_{bc}bc$	$S \rightarrow B_{bc}bc$	$S \rightarrow aB_{cb}cb$	$S \rightarrow aB_{cb}cb$	
B_{bc}	$B_{bc} \rightarrow \epsilon$	$B_{bc} \rightarrow b$	$B_{bc} \rightarrow c$			
B_{cb}	$B_{cb} \rightarrow b$		$B_{cb} \rightarrow \epsilon$			$B_{cb} \rightarrow c$

LL(k) vs. Strong-LL(k)

- **LL(k)** definition says

- $\omega A \alpha \Rightarrow \omega \beta \alpha, \omega A \alpha \Rightarrow \omega \gamma \alpha$
 $FNE_k(\beta \alpha) \cap FNE_k(\gamma \alpha) = \Phi$
- $x A \delta \Rightarrow x \beta \delta, x A \delta \Rightarrow x \gamma \delta$
 $FNE_k(\beta \delta) \cap FNE_k(\gamma \delta) = \Phi$

- **Strong-LL(k)** definition adds additional constraint

- $FNE_k(\beta \alpha) \cap FNE_k(\gamma \delta) = \Phi$
- $FNE_k(\beta \delta) \cap FNE_k(\gamma \alpha) = \Phi$

- **Why?** Because it relies on **Follow(A)** to get the context information, which always includes both α and δ

LL(1) = Strong LL(1) ?

- One question:
 - We saw an example grammar that is LL(2), yet not strong-LL(2)
 - Then, are there any example grammars that are LL(1), yet not Strong-LL(1)?
- The issue is the granularity of the lookahead
- The lookahead of LL(2) is finer than LL(1) since it look aheads more
- A nice exam question

Recursive-Descent Parsing

- Instead of stack, use recursive procedures
 - Sequence of production calls implicitly define parse tree
 - Given a parse table $M[A,a]$, it is easy to write one

```
extern token lookahead;
void match(token tok) {
    if (lookahead != tok) error();
    else lookahead = get_next_token();
}
void E() {
    switch (lookahead) {
        case 'id':
        case '(' : T(); Ep(); break;
        default : error();
    }
}
```

```
void Ep() {
    switch (lookahead) {
        case '+' : match('+'); T(); Ep(); break;
        case ')' :
        case '$' : break;
        default : error();
    }
}
...
main() {
    lookahead = get_next_token();
    E();
}
```

LL(1) Summary

- LL(1) parsing
 - Stack, lookahead, parsing table
 - Parsing table construction
 - Nullable(), FNE(), Follow()
- LL(1) grammar
 - Actually represent limited class of languages
 - i.e., many programming languages are not LL(1)
 - So, consider a larger class: LR bottom-up parsing