

# Chap 3. Parsing

COMP321 컴파일러

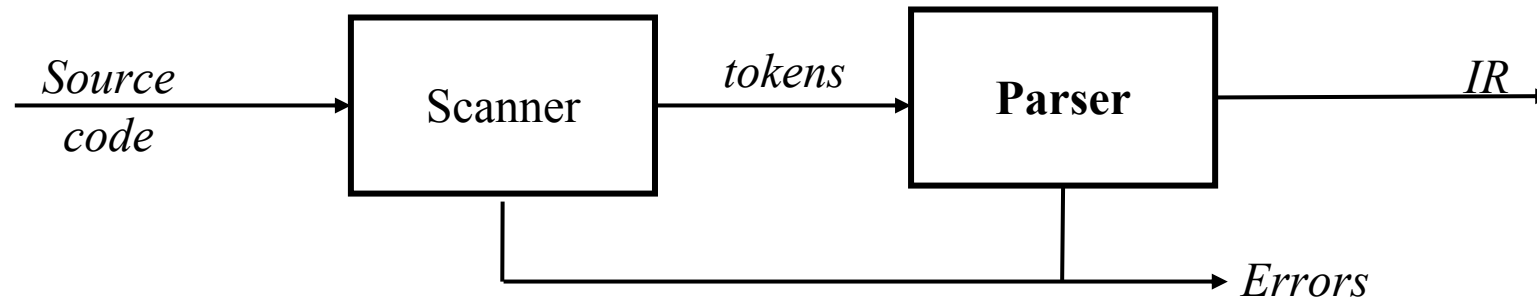
2007년 가을학기

경북대학교 전자전기컴퓨터학부

© 2004-7 N Baek @ GALab, KNU

## **3.1 Introduction**

# The Front End



## Parser

- Checks the stream of **words** and their parts of speech for grammatical correctness
- Determines if the input is **syntactically well formed**
- Guides checking at deeper levels than syntax
- Builds an **IR representation of the code**

# The Study of Parsing

- We need:
  - a grammar : **CFG** (context-free grammar)
    - $\rightarrow$  a language  $L(G)$
    - keep in mind that **our goal is building parsers**,  
not studying the mathematics of arbitrary languages
  - parsing methods
    - **top-down parsing**
      - hand-coded recursive descent parsers
    - **bottom-up parsing**
      - generated LR(1) parsers

## **3.2 Expressing Syntax**

# Specifying Syntax with a Grammar

- **parser**
  - **an engine** that determines whether or not the input program is a **syntactically valid sentence** in the programming language.
- **context-free grammar**
  - how to syntactically form the program
- example CFG for Sheep Noise in **Backus-Naur Form**

$$\begin{array}{l} \textit{SheepNoise} \rightarrow \underline{\textit{baa}} \textit{SheepNoise} \\ \quad \quad \quad | \quad \underline{\textit{baa}} \end{array}$$

– derivations:

- $\textit{SheepNoise} \rightarrow \underline{\textit{baa}} \textit{SheepNoise} \rightarrow \underline{\textit{baa}} \underline{\textit{baa}}$

# Context-Free Grammar

a grammar is a four tuple,  $G = (T, NT, s, P)$

- $T$  is a set of *terminal symbols* (words)
  - Lexical Analyzer (scanner) returns a terminal.
- $NT$  is a set of *non-terminal symbols* (syntactic variables)
- $s$  is the *start symbol* (set of strings in  $L(G)$ )
- $P$  is a set of *productions* or *rewrite rules*
  - $P: NT \rightarrow (NT \cup T)^+$

## BNF (Backus-Naur Form)

- the traditional notation used to represent CFG.
- original form:  $\langle \text{SheepNoise} \rangle ::= \text{baa} \langle \text{SheepNoise} \rangle$   
| baa

# Deriving Syntax

We can use the *SheepNoise* grammar to create sentences

- use the productions as *rewriting rules*
- 1:             $SheepNoise \rightarrow \underline{baa} SheepNoise$
- 2:                        | baa

Rule	Sentential Form
—	<i>SheepNoise</i>
2	<u>baa</u>

Rule	Sentential Form
—	<i>SheepNoise</i>
1	<u>baa</u> <i>SheepNoise</i>
2	<u>baa</u> <u>baa</u>

Rule	Sentential Form
—	<i>SheepNoise</i>
1	<u>baa</u> <i>SheepNoise</i>
1	<u>baa</u> <u>baa</u> <i>SheepNoise</i>
2	<u>baa</u> <u>baa</u> <u>baa</u>

*And so on ...*



# A More Useful Grammar

- To explore the uses of CFGs, we need a more complex grammar

1	$Expr$	$\rightarrow$	$Expr\ Op\ Expr$
2			<u>num</u>
3			<u>id</u>
4	$Op$	$\rightarrow$	<u>+</u>
5			<u>=</u>
6			<u>*</u>
7			<u>/</u>

- $G = (T, NT, s, P)$ 
  - $T = \{ \underline{\text{num}}, \underline{\text{id}}, \underline{+}, \underline{=}, \underline{*}, \underline{/} \}, NT = \{ Expr, Op \}, s = Expr$

# A More Useful Grammar

To explore the uses of CFGs, we need a more complex grammar

1	$Expr$	$\rightarrow$	$Expr Op Expr$
2			<u>num</u>
3			<u>id</u>
4	$Op$	$\rightarrow$	<u>+</u>
5			<u>=</u>
6			<u>*</u>
7			<u>/</u>

Rule	Sentential Form
—	$Expr$
1	$Expr Op Expr$
2	$\langle id, \underline{x} \rangle Op Expr$
5	$\langle id, \underline{x} \rangle - Expr$
1	$\langle id, \underline{x} \rangle - Expr Op Expr$
2	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle Op Expr$
6	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * Expr$
3	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$

We denote this derivation:  $Expr \rightarrow^* \underline{id} - \underline{num} * \underline{id}$

- Such a sequence of rewrites is called a ***derivation***
- Process of discovering a derivation is called ***parsing***

# Derivations

- At each step, we choose a non-terminal to replace
- Different choices can lead to different derivations

Two derivations are of interest

- **Leftmost derivation** — replace leftmost NT at each step
- **Rightmost derivation** — replace rightmost NT at each step

These are the two *systematic* derivations

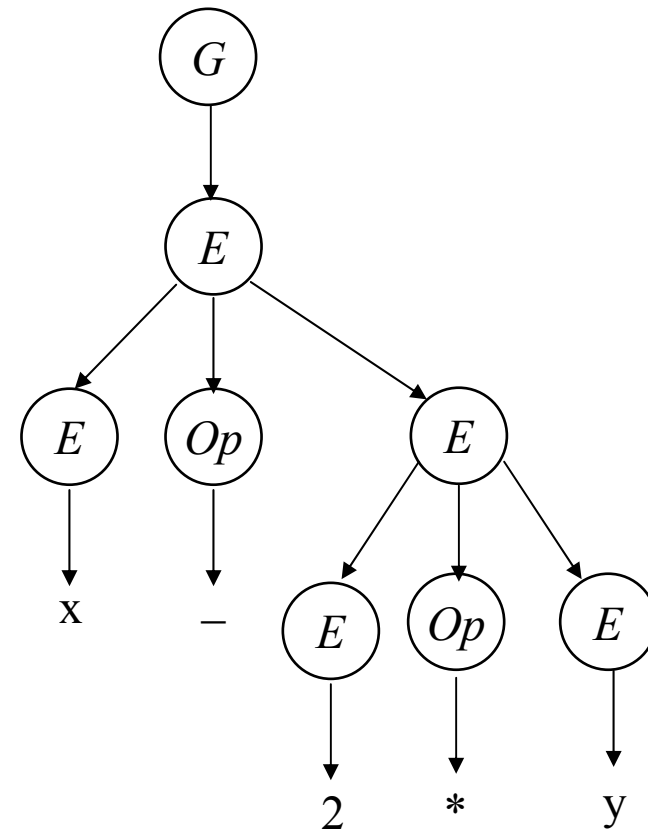
- We don't care about randomly-ordered derivations!

# Leftmost Derivation

1	$Expr$	$\rightarrow$	$Expr Op Expr$
2			<u>num</u>
3			<u>id</u>
4	$Op$	$\rightarrow$	<u>+</u>
5			<u>=</u>
6			<u>*</u>
7			<u>/</u>

Rule	Sentential Form
—	$Expr$
1	$Expr Op Expr$
3	$\langle id, \underline{x} \rangle Op Expr$
5	$\langle id, \underline{x} \rangle - Expr$
1	$\langle id, \underline{x} \rangle - Expr Op Expr$
2	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle Op Expr$
6	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * Expr$
3	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$

This evaluates as  $x - (2 * y)$

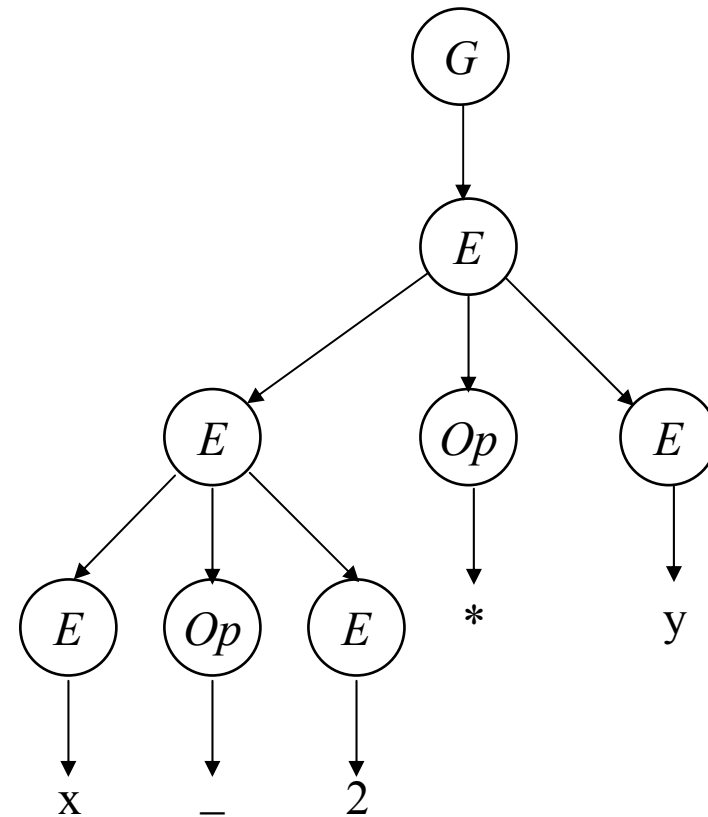


# Rightmost Derivation

1	$Expr$	$\rightarrow$	$Expr Op Expr$
2			<u>num</u>
3			<u>id</u>
4	$Op$	$\rightarrow$	$\pm$
5			$=$
6			$*$
7			$/$

Rule	Sentential Form
—	$Expr$
1	$Expr Op Expr$
3	$Expr Op <id, \underline{y}>$
6	$Expr * <id, \underline{y}>$
1	$Expr Op Expr * <id, \underline{y}>$
2	$Expr Op <num, \underline{2}> * <id, \underline{y}>$
5	$Expr - <num, \underline{2}> * <id, \underline{y}>$
3	$<id, \underline{x}> - <num, \underline{2}> * <id, \underline{y}>$

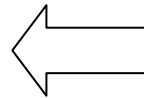
This evaluates as  $(x - 2) * y$



# Derivations and Precedence

- derivation에 따라, operator precedence가 다르게 나옴
  - compiler 입장에서는 심각한 문제
  - 해결책? grammar 자체의 수정
    - 새로운 non-terminal 도입

level two	1	<i>Goal</i>	→	<i>Expr</i>
	2	<i>Expr</i>	→	<i>Expr</i> + <i>Term</i>
	3			<i>Expr</i> - <i>Term</i>
	4			<i>Term</i>
level one	5	<i>Term</i>	→	<i>Term</i> * <i>Factor</i>
	6			<i>Term</i> / <i>Factor</i>
	7			<i>Factor</i>
	8	<i>Factor</i>	→	<u>number</u>
	9			<u>id</u>

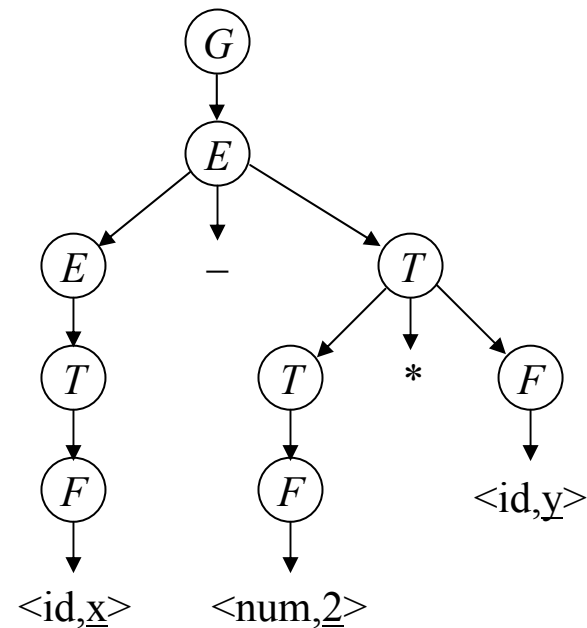


1	<i>Expr</i>	→	<i>Expr</i> <i>Op</i> <i>Expr</i>
2			<u>num</u>
3			<u>id</u>
4	<i>Op</i>	→	+
5			-
6			*
7			/

# Derivations and Precedence

- 수정된 grammar에서는 leftmost derivation, rightmost derivation 모두 같은 parse tree를 생성
  - 모두  $x - (2 * y)$  생성
  - 왜? **grammar** 자체에서 다른 **derivation**을 금지

<i>Rule</i>	<i>Sentential Form</i>
—	<i>Goal</i>
1	<i>Expr</i>
3	<i>Expr</i> – <i>Term</i>
5	<i>Expr</i> – <i>Term</i> * <i>Factor</i>
9	<i>Expr</i> – <i>Term</i> * $\langle \text{id}, \underline{y} \rangle$
7	<i>Expr</i> – <i>Factor</i> * $\langle \text{id}, \underline{y} \rangle$
8	<i>Expr</i> – $\langle \text{num}, \underline{2} \rangle$ * $\langle \text{id}, \underline{y} \rangle$
4	<i>Term</i> – $\langle \text{num}, \underline{2} \rangle$ * $\langle \text{id}, \underline{y} \rangle$
7	<i>Factor</i> – $\langle \text{num}, \underline{2} \rangle$ * $\langle \text{id}, \underline{y} \rangle$
9	$\langle \text{id}, \underline{x} \rangle$ – $\langle \text{num}, \underline{2} \rangle$ * $\langle \text{id}, \underline{y} \rangle$



*Its parse tree*

# Ambiguous Grammar

- original grammar에는 또다른 문제점 !
  - leftmost derivation만 적용해도 다른 parse tree 가능

1	$Expr$	$\rightarrow$	$Expr Op Expr$
2			<u>num</u>
3			<u>id</u>
4	$Op$	$\rightarrow$	$\pm$
5			$=$
6			$*$
7			$/$

Rule	Sentential Form
—	$Expr$
1	$Expr Op Expr$
3	$\langle id, \underline{x} \rangle Op Expr$
5	$\langle id, \underline{x} \rangle - Expr$
1	$\langle id, \underline{x} \rangle - Expr Op Expr$
2	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle Op Expr$
6	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * Expr$
3	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$

*Original choice*

Rule	Sentential Form
—	$Expr$
1	$Expr Op Expr$
1	$Expr Op Expr Op Expr$
3	$\langle id, \underline{x} \rangle Op Expr Op Expr$
5	$\langle id, \underline{x} \rangle - Expr Op Expr$
2	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle Op Expr$
6	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * Expr$
3	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$

*New choice*



# Ambiguous Grammars

## Definitions

- If a grammar has **more than one leftmost derivation** for a single *sentential form*, the grammar is *ambiguous*
- If a grammar has **more than one rightmost derivation** for a single sentential form, the grammar is *ambiguous*
- The leftmost and rightmost derivations for a sentential form may differ, even in an unambiguous grammar

Ambiguous grammar는 **automatic parsing 불가능**

# Ambiguous Grammar: If-Then-Else

- $Stmt \rightarrow \underline{\text{if}} \ Expr \ \underline{\text{then}} \ Stmt$   
|  $\underline{\text{if}} \ Expr \ \underline{\text{then}} \ Stmt \ \underline{\text{else}} \ Stmt$   
|  $AssignStmt$
- input:  $\underline{\text{if}} \ Expr_1 \ \underline{\text{then}} \ \underline{\text{if}} \ Expr_2 \ \underline{\text{then}} \ Stmt_1 \ \underline{\text{else}} \ Stmt_2$
- output1:  $\underline{\text{if}} \ Expr_1 \ \underline{\text{then}} \ \{ \underline{\text{if}} \ Expr_2 \ \underline{\text{then}} \ Stmt_1 \ \underline{\text{else}} \ Stmt_2 \}$
- output2:  $\underline{\text{if}} \ Expr_1 \ \underline{\text{then}} \ \{ \underline{\text{if}} \ Expr_2 \ \underline{\text{then}} \ Stmt_1 \} \ \underline{\text{else}} \ Stmt_2$
- 해결책
  - in Shell, Perl, etc.: introduce **elif, endif**
  - in most programming languages:
    - use **inner-most unmatched-if rule**

# Ambiguous Grammar: If-Then-Else

Removing the ambiguity

- Must rewrite the grammar to avoid generating the problem
- **Match each else to innermost unmatched if**
- With this grammar, the example has only one derivation

- $Stmt \rightarrow \underline{\text{if}} \ Expr \ \underline{\text{then}} \ Stmt$   
          |  $\underline{\text{if}} \ Expr \ \underline{\text{then}} \ WithElse \ \underline{\text{else}} \ Stmt$   
          |  $AssignStmt$
- $WithElse \rightarrow \underline{\text{if}} \ Expr \ \underline{\text{then}} \ WithElse \ \underline{\text{else}} \ WithElse$   
              |  $AssignStmt$

# An Example: If-Then-Else

- input:  $\underline{\text{if}}\ Expr_1\ \underline{\text{then}}\ \underline{\text{if}}\ Expr_2\ \underline{\text{then}}\ Stmt_1\ \underline{\text{else}}\ Stmt_2$

1	$Stmt$	$\rightarrow$	$\underline{\text{if}}\ Expr\ \underline{\text{then}}\ Stmt$
2			$\underline{\text{if}}\ Expr\ \underline{\text{then}}\ WithElse\ \underline{\text{else}}\ Stmt$
3			$\underline{AssignStmt}$
4	$WithElse$	$\rightarrow$	$\underline{\text{if}}\ Expr\ \underline{\text{then}}\ WithElse\ \underline{\text{else}}\ WithElse$
5			$\underline{AssignStmt}$

Rule	Sentential Form
—	$Stmt$
1	$\underline{\text{if}}\ Expr\ \underline{\text{then}}\ Stmt$
2	$\underline{\text{if}}\ Expr\ \underline{\text{then}}\ \underline{\text{if}}\ Expr\ \underline{\text{then}}\ WithElse\ \underline{\text{else}}\ Stmt$
3	$\underline{\text{if}}\ Expr\ \underline{\text{then}}\ \underline{\text{if}}\ Expr\ \underline{\text{then}}\ WithElse\ \underline{\text{else}}\ AssignStmt$
5	$\underline{\text{if}}\ Expr\ \underline{\text{then}}\ \underline{\text{if}}\ Expr\ \underline{\text{then}}\ AssignStmt\ \underline{\text{else}}\ AssignStmt$

# CFG and RE

- Regular Expression : Regular Grammar 와 equivalent
- **RG** (regular grammar) : **CFG with only left-linear rules.**
  - All rules are:  $A \rightarrow a$  or  $A \rightarrow a B$ 
    - $A, B \in NT, a \in T$
- 즉, CFG 만으로 RE 까지 표현 가능
  - 그러나, 여전히 **RE**를 쓴다
    - DFA-based scanner가 more efficient
    - comment 처리 ? CFG로는 처리 곤란
- CFG의 계층 구조
  - **$RG \subset LL(1) \subset LR(1) \subset CFG$**

## **3.3 Top-Down Parsing**

# Parsing Techniques

## *Top-down parsers (LL(1), recursive descent)*

- Start at the root of the parse tree and grow toward leaves
- Pick a production & try to match the input
- Bad “pick”  $\Rightarrow$  may need to backtrack
- Some grammars are backtrack-free (*predictive parsing*)

## *Bottom-up parsers (LR(1), operator precedence)*

- Start at the leaves and grow toward root
- As input is consumed, encode possibilities in an internal state
- Start in a state valid for legal first tokens
- Bottom-up parsers handle a large class of grammars

# Top-down Parsing Algorithm

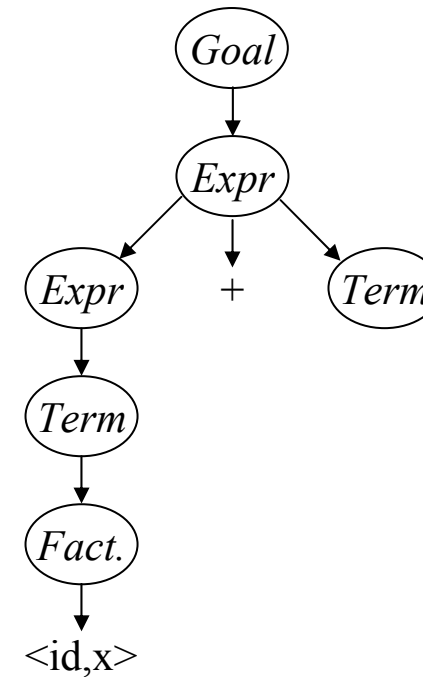
- Construct the **root** node of the parse tree
- Repeat until the leaf nodes of the parse tree matches the input string
  - At a node labeled  $A$ ,  
**select a production** with  $A$  on its LHS and,  
for each symbol on its RHS,  
construct the appropriate child
  - When a terminal symbol is added to the leaf node and it  
doesn't match the leaf node, **backtrack**
  - Find the next node to be expanded
- 어떤 **production rule**을 고르느냐 ?



# An Example

- input:  $x - 2 * y$   
 –  $\uparrow$  : current scanning pointer

1	<i>Goal</i>	$\rightarrow$	<i>Expr</i>
2	<i>Expr</i>	$\rightarrow$	<i>Expr</i> + <i>Term</i>
3			<i>Expr</i> – <i>Term</i>
4			<i>Term</i>
5	<i>Term</i>	$\rightarrow$	<i>Term</i> * <i>Factor</i>
6			<i>Term</i> / <i>Factor</i>
7			<i>Factor</i>
8	<i>Factor</i>	$\rightarrow$	<u>number</u>
9			<u>id</u>

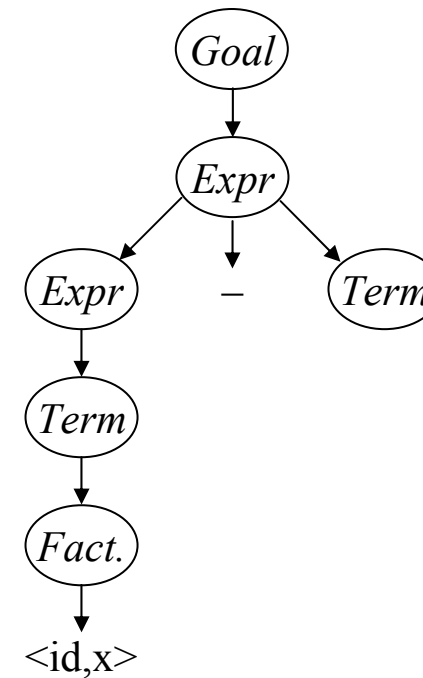


Rule	Sentential Form	Input
—	<i>Goal</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
1	<i>Expr</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
2	<i>Expr</i> + <i>Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
4	<i>Term</i> + <i>Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
7	<i>Factor</i> + <i>Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
9	$\langle \text{id}, x \rangle$ + <i>Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
9	$\langle \text{id}, x \rangle$ + <i>Term</i>	$\underline{x} \uparrow - \underline{2} * \underline{y}$

# An Example

- input:  $x - 2 * y$
- “-” doesn’t match “+”
- We need **backtracking** !

Rule	Sentential Form	Input
—	<i>Goal</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
1	<i>Expr</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
2	<i>Expr</i> + <i>Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
4	<i>Term</i> + <i>Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
7	<i>Factor</i> + <i>Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
9	$\langle \text{id}, x \rangle$ + <i>Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
9	$\langle \text{id}, x \rangle$ + <i>Term</i>	$\underline{x} \uparrow - \underline{2} * \underline{y}$



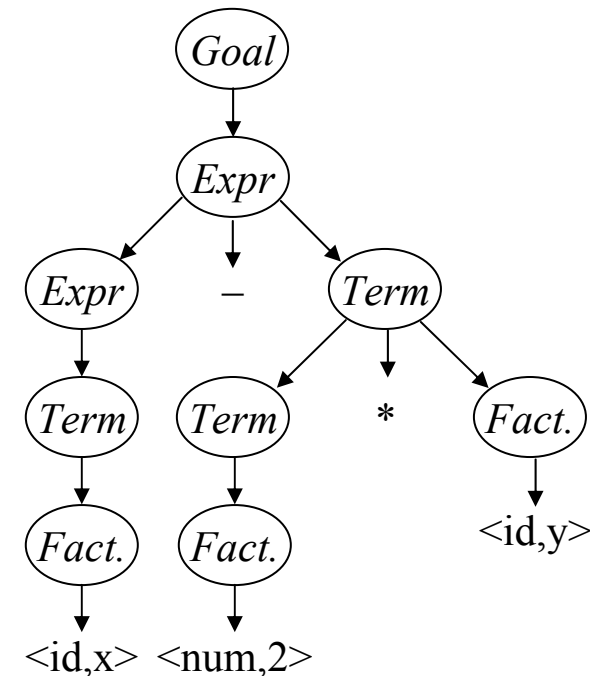
Rule	Sentential Form	Input
—	<i>Goal</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
1	<i>Expr</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
3	<i>Expr</i> - <i>Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
4	<i>Term</i> - <i>Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
7	<i>Factor</i> - <i>Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
9	$\langle \text{id}, x \rangle$ - <i>Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
9	$\langle \text{id}, x \rangle$ - <i>Term</i>	$\underline{x} \uparrow - \underline{2} * \underline{y}$
—	$\langle \text{id}, x \rangle$ - <i>Term</i>	$\underline{x} - \uparrow \underline{2} * \underline{y}$

# An Example

- continue from  $x - 2 * y$  : one more **backtracking**

Rule	Sentential Form	Input
—	$\langle \text{id}, x \rangle - \text{Term}$	$\underline{x} - \uparrow \underline{2} * \underline{y}$
7	$\langle \text{id}, x \rangle - \text{Factor}$	$\underline{x} - \uparrow \underline{2} * \underline{y}$
9	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle$	$\underline{x} - \uparrow \underline{2} * \underline{y}$
—	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle$	$\underline{x} - \underline{2} \uparrow * \underline{y}$

Rule	Sentential Form	Input
—	$\langle \text{id}, x \rangle - \text{Term}$	$\underline{x} - \uparrow \underline{2} * \underline{y}$
5	$\langle \text{id}, x \rangle - \text{Term} * \text{Factor}$	$\underline{x} - \uparrow \underline{2} * \underline{y}$
7	$\langle \text{id}, x \rangle - \text{Factor} * \text{Factor}$	$\underline{x} - \uparrow \underline{2} * \underline{y}$
8	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle * \text{Factor}$	$\underline{x} - \uparrow \underline{2} * \underline{y}$
—	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle * \text{Factor}$	$\underline{x} - \underline{2} \uparrow * \underline{y}$
—	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle * \text{Factor}$	$\underline{x} - \underline{2} * \uparrow \underline{y}$
9	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle * \langle \text{id}, y \rangle$	$\underline{x} - \underline{2} * \uparrow \underline{y}$
—	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle * \langle \text{id}, y \rangle$	$\underline{x} - \underline{2} * \underline{y} \uparrow$



# Another Possible Parse

Other choices for expansion are possible : no input consume !

<i>Rule</i>	<i>Sentential Form</i>	<i>Input</i>
—	<i>Goal</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
1	<i>Expr</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
2	<i>Expr + Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
2	<i>Expr + Term + Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
2	<i>Expr + Term + Term + Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$
2	<i>Expr + Term + Term + ... + Term</i>	$\uparrow \underline{x} - \underline{2} * \underline{y}$

This doesn't terminate

- Wrong choice of expansion leads to **non-termination**
- Non-termination is a **bad property** for a parser to have
- Parser must make the right choice

# Left Recursion

- 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

# Eliminating Left Recursion

To remove left recursion, we can  
**transform the grammar**

Consider a grammar fragment of the  
form

$$\begin{array}{l} Fee \rightarrow Fee \alpha \\ \quad | \beta \end{array}$$

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

We can rewrite this as

$$\begin{array}{l} Fee \rightarrow \beta Fie \\ Fie \rightarrow \alpha Fie \\ \quad | \epsilon \end{array}$$

where *Fie* is a new non-terminal

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

# Example

- Original Grammar

1	<i>Goal</i>	$\rightarrow$	<i>Expr</i>
2	<i>Expr</i>	$\rightarrow$	<i>Expr</i> + <i>Term</i>
3			<i>Expr</i> - <i>Term</i>
4			<i>Term</i>
5	<i>Term</i>	$\rightarrow$	<i>Term</i> * <i>Factor</i>
6			<i>Term</i> / <i>Factor</i>
7			<i>Factor</i>
8	<i>Factor</i>	$\rightarrow$	<u>number</u>
9			<u>id</u>

- New Grammar

1	<i>Goal</i>	$\rightarrow$	<i>Expr</i>
2	<i>Expr</i>	$\rightarrow$	<i>Term Expr'</i>
3	<i>Expr'</i>	$\rightarrow$	+ <i>Term Expr'</i>
4			- <i>Term Expr'</i>
5			$\epsilon$
6	<i>Term</i>	$\rightarrow$	<i>Factor Term'</i>
7	<i>Term'</i>	$\rightarrow$	* <i>Factor Term'</i>
8			/ <i>Factor Term'</i>
9			$\epsilon$
10	<i>Factor</i>	$\rightarrow$	<u>number</u>
11			<u>id</u>
12			( <i>Expr</i> )

# Eliminating Left Recursion

- general algorithm:

arrange the NTs into some order  $A_1, A_2, \dots, A_n$

**for**  $i \leftarrow 1$  **to**  $n$

**for**  $s \leftarrow 1$  **to**  $i - 1$

        replace each production  $A_i \rightarrow A_s \gamma$

        with  $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma,$

        where  $A_s \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$  are

        all the current productions for  $A_s$

**eliminate any immediate left recursion** on  $A_i$

    using the direct transformation

- 이제, termination은 보장... backtracking 해결 필요



# Backtrack-Free Grammar

- parser가 backtrack을 안 하려면,  
항상 correct choice 를 하면 된다.
- key idea
  - $A \rightarrow \alpha \mid \beta$  이고,
  - $\alpha \rightarrow^* \underline{a} \gamma$  and  $\beta \rightarrow^* \underline{b} \lambda$
  - scanner에서 look-ahead 한 terminal 이  
 $\underline{a}, \underline{b}$  중 어느 것이냐에 따라, correct choice 가능 !
  - $\underline{a}$  이면,  $A \rightarrow \alpha \rightarrow^* \underline{a} \gamma$  로 진행
  - $\underline{b}$  이면,  $A \rightarrow \beta \rightarrow^* \underline{b} \lambda$  로 진행
- **LL(1) grammar property !**

# FIRST( )

- **FIRST( $\alpha$ )** is the set of words that can appear as the first symbol in some string derived from  $\alpha$ .
  - $\alpha \in T \cup NT \cup \{ \varepsilon \}$
  - $\underline{x} \in \text{FIRST}(\alpha)$  *iff*  $\alpha \Rightarrow^* \underline{x} \gamma$ , for some  $\gamma$
  - $\varepsilon \in \text{FIRST}(\alpha)$  *iff*  $\alpha \Rightarrow^* \varepsilon$
- example
  - $Expr' \rightarrow + Term Expr'$ 
    - |  $- Term Expr'$
    - |  $\varepsilon$
  - $\text{FIRST}(Expr') = \{ +, -, \varepsilon \}$
- see page 99 for the detailed algorithm.

# FOLLOW( )

- **FOLLOW( $A$ )** is the set of symbols that can occur immediately after some non-terminal  $A$  in a valid sentence.
  - $B \rightarrow A \underline{a}$   
 $A \rightarrow \underline{c} \mid \epsilon$
  - $\text{FIRST}(A) = \{ \underline{c}, \epsilon \}$
  - $\text{FOLLOW}(A) = \{ \underline{a} \}$
  - $\epsilon \in \text{FIRST}(A)$  일 때 의미를 가짐
  - see page 99 for the detailed algorithm
- $\text{FIRST}^+(\alpha) = \text{FIRST}(\alpha)$  if  $\epsilon \notin \text{FIRST}(\alpha)$
- $\text{FIRST}^+(\alpha) = \text{FIRST}(\alpha) \cup \text{FOLLOW}(\alpha)$

# LL(1) Property

- LL(1) property: **backtrack-free grammar**의 조건
  - for any non-terminal  $A$   
with the rule  $A \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$ ,
  - it must be true that  
 $\text{FIRST}^+(\beta_i) \cap \text{FIRST}^+(\beta_j) = \emptyset, \quad \forall 1 \leq i < j \leq n$
- examples

$Factor \rightarrow \underline{Identifier}$   
 $\quad \mid \underline{Identifier} [ ExprList ]$   
 $\quad \mid \underline{Identifier} ( ExprList )$

no backtrack-free

$Factor \rightarrow \underline{Identifier} Arguments$   
 $Arguments \rightarrow [ ExprList ]$   
 $\quad \mid ( ExprList )$   
 $\quad \mid \epsilon$

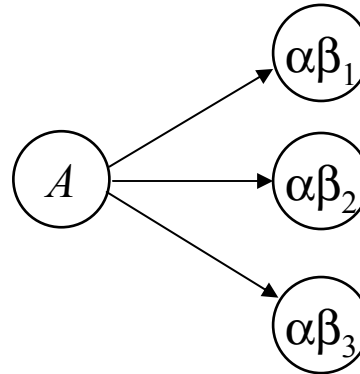
backtrack-free

# Left Factoring

What if my grammar does not have the LL(1) property?

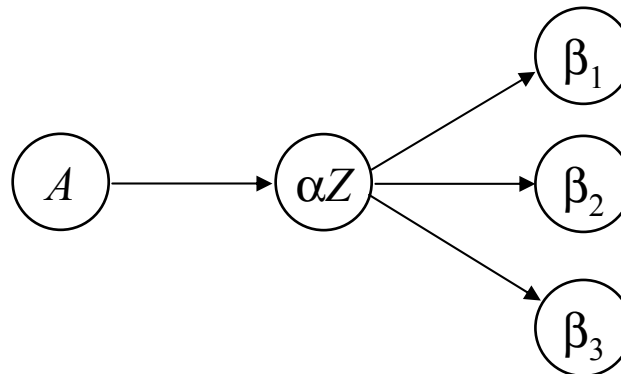
⇒ Sometimes, we can **transform the grammar**

$$\begin{array}{l} A \rightarrow \alpha\beta_1 \\ \quad | \alpha\beta_2 \\ \quad | \alpha\beta_3 \end{array}$$



becomes

$$\begin{array}{l} A \rightarrow \alpha Z \\ Z \rightarrow \beta_1 \\ \quad | \beta_2 \\ \quad | \beta_3 \end{array}$$



# Left Factoring

The Algorithm:

$\forall A \in NT,$

find **the longest prefix  $\alpha$**  that occurs in two  
or more right-hand sides of  $A$

if  $\alpha \neq \varepsilon$  then replace all of the  $A$  productions,

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \gamma,$$

with

$$A \rightarrow \alpha Z \mid \gamma$$

$$Z \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

where  $Z$  is a new element of  $NT$

Repeat until no common prefixes remain

# Predictive Parsing

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

- Can **write a simple routine** to recognize each  $lhs$
- Code is both simple & fast

Consider  $A \rightarrow \beta_1 \mid \beta_2 \mid \beta_3$ , with

$$\text{FIRST}^+(\beta_1) \cap \text{FIRST}^+(\beta_2) \cap \text{FIRST}^+(\beta_3) = \emptyset$$

```
/* find an A */  
if (current_word  $\in$  FIRST( $\beta_1$ ))  
    find a  $\beta_1$  and return true  
else if (current_word  $\in$  FIRST( $\beta_2$ ))  
    find a  $\beta_2$  and return true  
else if (current_word  $\in$  FIRST( $\beta_3$ ))  
    find a  $\beta_3$  and return true  
else  
    report an error and return false
```

직접 C code 작성 !

# Recursive-Descent Parser

- one kind of predicted parser
  - 대표적인 hand-written LL(1) parser

1	<i>Goal</i>	$\rightarrow$	<i>Expr</i>
2	<i>Expr</i>	$\rightarrow$	<i>Term Expr'</i>
3	<i>Expr'</i>	$\rightarrow$	$+ \textit{Term Expr'}$
4		$ $	$- \textit{Term Expr'}$
5		$ $	$\epsilon$
6	<i>Term</i>	$\rightarrow$	<i>Factor Term'</i>
7	<i>Term'</i>	$\rightarrow$	$* \textit{Factor Term'}$
8		$ $	$/ \textit{Factor Term'}$
9		$ $	$\epsilon$
10	<i>Factor</i>	$\rightarrow$	<u>number</u>
11		$ $	<u>id</u>

This produces a parser with six mutually recursive routines:

- *Goal*
- *Expr*
- *EPrime*
- *Term*
- *TPrime*
- *Factor*

Each recognizes one *NT* or *T*

The term descent refers to the direction in which the parse tree is built.



# Recursive-Descent Parser

*Goal*( ) : Goal  $\rightarrow$  Expr

token  $\leftarrow$  *next\_token*( );

**if** (*Expr*( ) = **true** & token = **EOF**)

**then** next compilation step;

**else**

report syntax error;

**return** false;

*Expr*( ) : Expr  $\rightarrow$  Term Expr'

**if** (*Term*( ) = **false**)

**then return** false;

**else return** *Eprime*( );

*Factor*( ) : Factor  $\rightarrow$  number | id

**if** (token = *Number*) **then**

token  $\leftarrow$  *next\_token*( );

**return true**;

**else if** (token = *Identifier*) **then**

token  $\leftarrow$  *next\_token*( );

**return true**;

**else**

report syntax error;

**return false**;

- see page 105 for all routines

# Recursive Descent (Summary)

1. Build FIRST (and FOLLOW) sets
2. Massage grammar to have  $LL(1)$  condition
  - a. Remove left recursion
  - b. Left factor it
3. Define a procedure for each non-terminal
  - a. Implement a case for each right-hand side
  - b. Call procedures as needed for non-terminals
4. Add extra code, as needed
  - a. Perform context-sensitive checking
  - b. Build an IR to record the code

# Building Top-down Parsers

Given an  $LL(1)$  grammar, and its FIRST & FOLLOW sets ...

- Emit a routine for each non-terminal
  - Nest of if-then-else statements to check alternate RHS's
  - Each returns true on success and throws an error on false
  - Simple, working (*, perhaps ugly,*) code

Improving matters

- **What about a table to encode the options?**
  - Interpret the table with a skeleton, as we did in scanning

# Building Top-down Parsers

## Strategy

- **Encode knowledge in a table**
- Use a standard “skeleton” parser to interpret the table

## Example

- 10: Factor  $\rightarrow$  id
- 11: Factor  $\rightarrow$  number
- Table might look like:

- Table might look like:
 

	Terminal Symbols							
	+	-	*	/	Id.	Num.	EOF	
Non-terminal Symbols	<u>Factor</u>	⊖	—	—	—	10	11	—

Error on '+'

Reduce by rule 10 on 'Id'

© 2006 N. Baek @ GALab, KNU

44

# Building Top Down Parsers

Building the complete table

- Need a row for every  $NT$  & a column for every  $T$
- Need an algorithm to build the table

Filling in  $TABLE[X, y]$ ,  $X \in NT, y \in T$

1. entry is the rule  $X \rightarrow \beta$ , if  $y \in \mathbf{FIRST}(\beta)$
2. entry is the rule  $X \rightarrow \epsilon$   
if  $y \in \mathbf{FOLLOW}(X)$  and  $X \rightarrow \epsilon \in G$
3. entry is **error** if **neither 1 nor 2** define it

If any entry is defined multiple times,  $G$  is not  $LL(1)$

This is the  **$LL(1)$  table construction algorithm**

## **3.4 Bottom-Up Parsing**

# Parsing Techniques

## *Top-down parsers*

LL(1)

- Start at the root of the parse tree and grow toward leaves
- **Pick a production & try to match the input**
- Bad “pick”  $\Rightarrow$  may need to backtrack
- Some grammars are backtrack-free

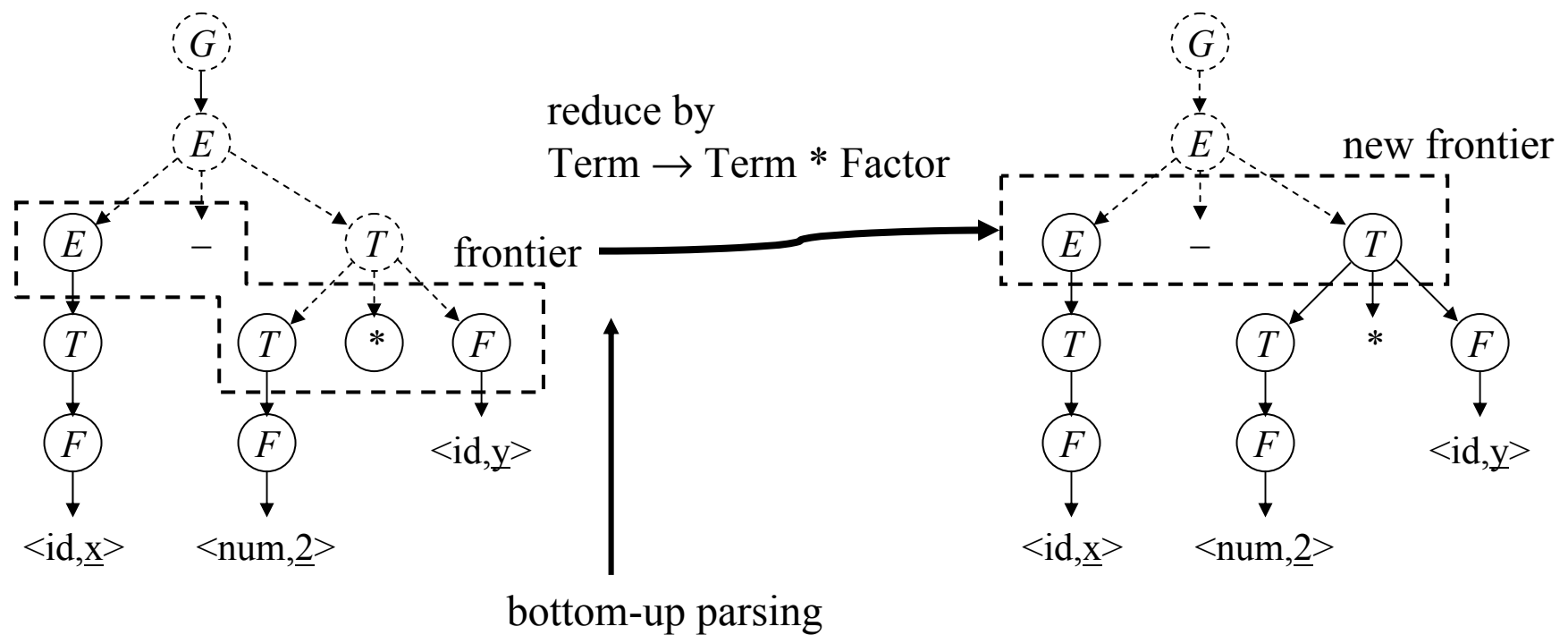
## *Bottom-up parsers*

LR(1)

- Start at the leaves and grow toward root
- As input is consumed, **encode possibilities in an internal state**
- Start in a state valid for legal first tokens
- Bottom-up parsers handle a large class of grammars

# Frontier & Reduce

- **frontier** (= upper frontier)
  - bottom-up parsing 중의, 현재 가장 upper node들





# Bottom-up Parsing

- parsing의 핵심은 derivation

$$S \Rightarrow \gamma_0 \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \dots \Rightarrow \gamma_{n-1} \Rightarrow \gamma_n \Rightarrow \text{sentence}$$

- **sentence** : a set of terminals (no non-terminal)

- bottom-up parsing

- input string : 완전히 derive된 sentence 라고 가정

- parsing 과정 : **sentence**에서 **start symbol S**까지 형성

$$S \Rightarrow \gamma_0 \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \dots \Rightarrow \gamma_{n-1} \Rightarrow \gamma_n \Rightarrow \text{sentence}$$

← bottom-up

- To reduce  $\gamma_i$  to  $\gamma_{i-1}$ ,  
find  $\beta$  in  $\gamma_i$ , then replace it using  $A \rightarrow \beta$  (handle)

- $\gamma_{i-1} = \alpha A \delta \Rightarrow \gamma_i = \alpha \beta \delta$

← bottom-up

# Handle

- rightmost derivation을 가정
  - 실제로는 backward로 처리  $\rightarrow$  leftmost  $NT$ 부터 처리
- A **handle** of a right-sentential form  $\gamma$  is a **pair**  $\langle A \rightarrow \beta, k \rangle$ 
  - $A \rightarrow \beta \in P$
  - $k$  is **the position** in  $\gamma$  of  $\beta$ 's rightmost symbol.
- handle을 선택하면,
  - **replacing  $\beta$  at  $k$  with  $A$**  produces the right sentential form from which  $\gamma$  is derived in the rightmost derivation.
  - right-sentential form이므로,  
**right of a handle은 only terminal symbols**

# Example: Grammar

- original grammar

2	$Expr$	$\rightarrow$	$Expr + Term$
3		$ $	$Expr - Term$
4		$ $	$Term$
5	$Term$	$\rightarrow$	$Term * Factor$
6		$ $	$Term / Factor$
7		$ $	$Factor$
8	$Factor$	$\rightarrow$	<u>number</u>
9		$ $	<u>id</u>
10		$ $	$(Expr)$

- augmented grammar
  - $S \rightarrow S'$  형태를 추가
  - 성공(accept) 판정을 쉽게

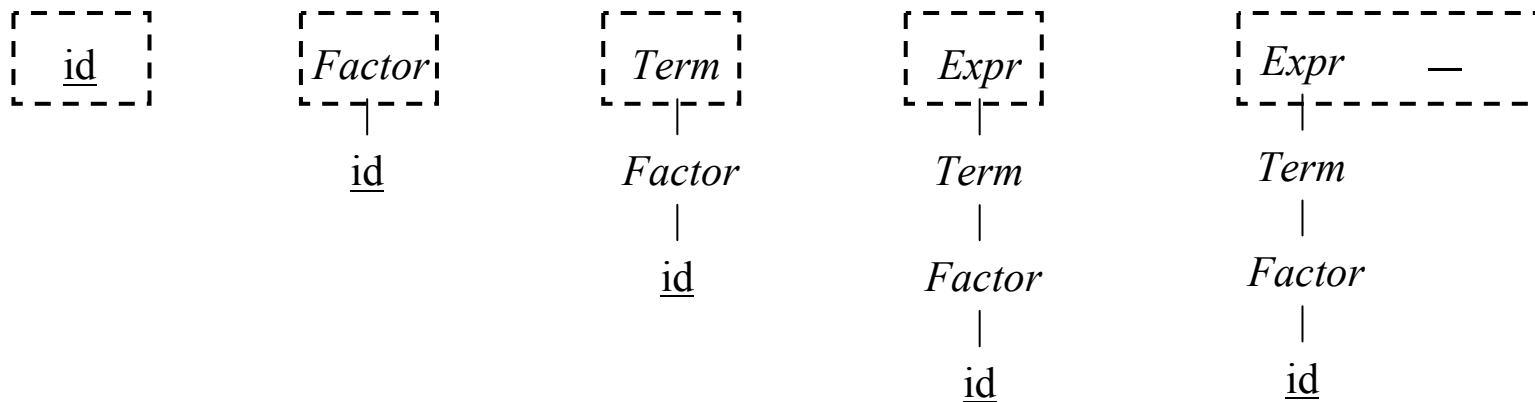
1	$Goal$	$\rightarrow$	$Expr$
2	$Expr$	$\rightarrow$	$Expr + Term$
3		$ $	$Expr - Term$
4		$ $	$Term$
5	$Term$	$\rightarrow$	$Term * Factor$
6		$ $	$Term / Factor$
7		$ $	$Factor$
8	$Factor$	$\rightarrow$	<u>number</u>
9		$ $	<u>id</u>
10		$ $	$(Expr)$

# Example: parsing states

frontier (stack)	next token	handle	action
	<u>id</u> – <u>num</u> * <u>id</u>		shift
<u>id</u>	– <u>num</u> * <u>id</u>	$\langle Factor \rightarrow \underline{id}, 1 \rangle$	reduce
<i>Factor</i>	– <u>num</u> * <u>id</u>	$\langle Term \rightarrow Factor, 1 \rangle$	reduce
<i>Term</i>	– <u>num</u> * <u>id</u>	$\langle Expr \rightarrow Term, 1 \rangle$	reduce
<i>Expr</i>	– <u>num</u> * <u>id</u>		shift
<i>Expr</i> –	<u>num</u> * <u>id</u>		shift
<i>Expr</i> – <u>num</u>	* <u>id</u>	$\langle Factor \rightarrow \underline{num}, 3 \rangle$	reduce
<i>Expr</i> – <i>Factor</i>	* <u>id</u>	$\langle Term \rightarrow Factor, 3 \rangle$	reduce
<i>Expr</i> – <i>Term</i>	* <u>id</u>		shift
<i>Expr</i> – <i>Term</i> *	<u>id</u>		shift
<i>Expr</i> – <i>Term</i> * <u>id</u>	<eof>	$\langle Factor \rightarrow \underline{id}, 5 \rangle$	reduce
<i>Expr</i> – <i>Term</i> * <i>Factor</i>	<eof>	$\langle Term \rightarrow Term * Factor, 5 \rangle$	reduce
<i>Expr</i> – <i>Term</i>	<eof>	$\langle Expr \rightarrow Expr - Term, 3 \rangle$	reduce
<i>Epxr</i>	<eof>	$\langle Goal \rightarrow Expr, 1 \rangle$	reduce
<i>Goal</i>	<eof>		accept

# Example: parse tree

frontier (stack)	next token	handle	action
	<u>id</u> – <u>num</u> * <u>id</u>		shift
<u>id</u>	– <u>num</u> * <u>id</u>	$\langle Factor \rightarrow \underline{id}, 1 \rangle$	reduce
<i>Factor</i>	– <u>num</u> * <u>id</u>	$\langle Term \rightarrow Factor, 1 \rangle$	reduce
<i>Term</i>	– <u>num</u> * <u>id</u>	$\langle Expr \rightarrow Term, 1 \rangle$	reduce
<i>Expr</i>	– <u>num</u> * <u>id</u>		shift
<i>Expr</i> –	<u>num</u> * <u>id</u>		shift
...	...	...	...



# Bottom-up Parser의 action

- data structure : a stack + scanner output
- **shift** : scanner  $\rightarrow$  push on the stack
  - how many shifts in a parser? # of tokens
- **reduce** : pop from the stack  $\rightarrow$  push a NT on the stack
  - how many reduces in a parser?
    - $O(\text{rules}) = O(\text{\# of tokens})$
- so, finite !
- reduce에서의 특징 : stack top 에서 pop
  - $\rightarrow$  more simple notation
  - **place holder** : •  $\rightarrow$  current stack top을 의미

# Example: new handle notation

frontier (stack)	next token	handle (new notation)	action
	<u>id</u> – <u>num</u> * <u>id</u>		shift
<u>id</u>	– <u>num</u> * <u>id</u>	$\langle Factor \rightarrow \underline{id} \bullet \rangle$	reduce
<i>Factor</i>	– <u>num</u> * <u>id</u>	$\langle Term \rightarrow Factor \bullet \rangle$	reduce
<i>Term</i>	– <u>num</u> * <u>id</u>	$\langle Expr \rightarrow Term \bullet \rangle$	reduce
<i>Expr</i>	– <u>num</u> * <u>id</u>		shift
<i>Expr</i> –	<u>num</u> * <u>id</u>		shift
<i>Expr</i> – <u>num</u>	* <u>id</u>	$\langle Factor \rightarrow \underline{num} \bullet \rangle$	reduce
<i>Expr</i> – <i>Factor</i>	* <u>id</u>	$\langle Term \rightarrow Factor \bullet \rangle$	reduce
<i>Expr</i> – <i>Term</i>	* <u>id</u>		shift
<i>Expr</i> – <i>Term</i> *	<u>id</u>		shift
<i>Expr</i> – <i>Term</i> * <u>id</u>	<eof>	$\langle Factor \rightarrow \underline{id} \bullet \rangle$	reduce
<i>Expr</i> – <i>Term</i> * <i>Factor</i>	<eof>	$\langle Term \rightarrow Term * Factor \bullet \rangle$	reduce
<i>Expr</i> – <i>Term</i>	<eof>	$\langle Expr \rightarrow Expr - Term \bullet \rangle$	reduce
<i>Epxr</i>	<eof>	$\langle Goal \rightarrow Expr \bullet \rangle$	reduce
<i>Goal</i>	<eof>		accept

# Potential Handle

frontier (stack)	next token	handle (new notation)	action
	<u>id</u> – <u>num</u> * <u>id</u>		shift
<u>id</u>	– <u>num</u> * <u>id</u>	$\langle Factor \rightarrow id \bullet \rangle$	reduce
<i>Factor</i>	– <u>num</u> * <u>id</u>	$\langle Term \rightarrow Factor \bullet \rangle$	reduce
<i>Term</i>	– <u>num</u> * <u>id</u>	$\langle Expr \rightarrow Term \bullet \rangle$	reduce
<i>Expr</i>	– <u>num</u> * <u>id</u>	why empty ?	shift
...	...	...	...

rules with *Expr* on RHS:

$Goal \rightarrow Expr$

$Expr \rightarrow Expr + Term$

$Expr \rightarrow Expr - Term$

potential handles:

$Goal \rightarrow Expr \bullet$

$Expr \rightarrow Expr \bullet + Term$

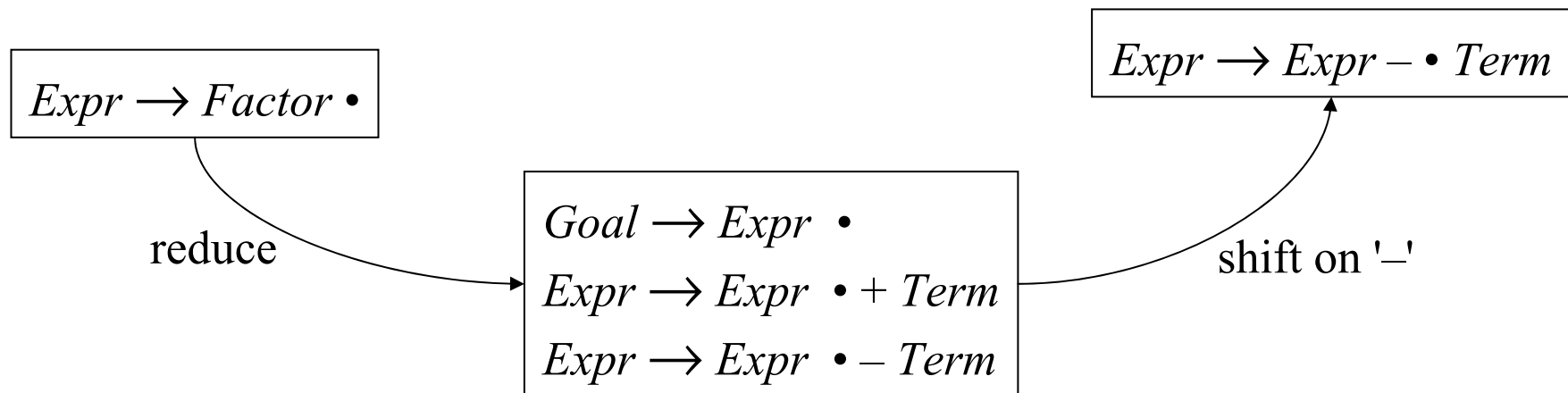
$Expr \rightarrow Expr \bullet - Term$

이 중의 하나이다 !



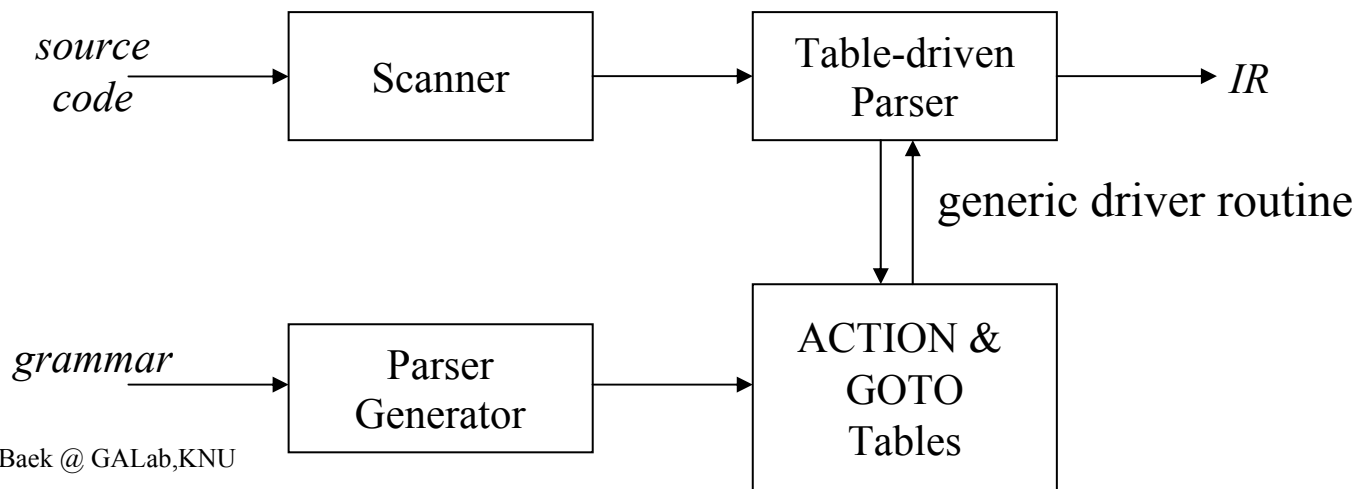
# LR(1) Parser Strategy

- rules  $\rightarrow$  all possible handles
  - rule with  $k$  symbols  $\rightarrow (k + 1)$  possible handles
- **handle recognizing DFA**
  - **states** : power set of all possible handles
  - **shift** : jump to another state with a token shift
  - **reduce** : jump to another state with a stack reduction



# LR(1) Parser

- LR(1)
  - **left-to-right scan** : scanner는 오직 한번 left-to-right scan
  - **rightmost derivation** : parsing tree는 rightmost derivation
  - **1 symbol of look-ahead**
- handle recognizing DFA + frontier stack 사용
  - DFA state + look-ahead symbol  $\rightarrow$  new DFA state
  - DFA transition function  $\rightarrow$  2차원 table



# LR(1) Parser의 4가지 action

- handle finding을 좀더 쉽게 하려면?
  - handle recognition **DFA**의 **state**도 **stack**에 기억하자 !
  - stack에 symbol, state가 교대로 push 된다.
- **shift** : get a token  $x$ , then push  $x, s_{\text{new}}$ 
  - $\dots (y, s_a) \rightarrow \dots (y, s_a), (x, s_{\text{new}})$
- **reduce** : apply a rule and go to a new state
  - $\dots (w, s_a), (x, s_b), (y, s_c), (z, s_d) \rightarrow \dots (w, s_a), (A, s_{\text{new}})$
- **accept** : stop parsing & report success
  - $(\$, s_0), (Goal, s_n), \text{token} = \text{EOF} \rightarrow \text{accept} !$
- **error** : call an error reporting / recovery routine

# Two Tables

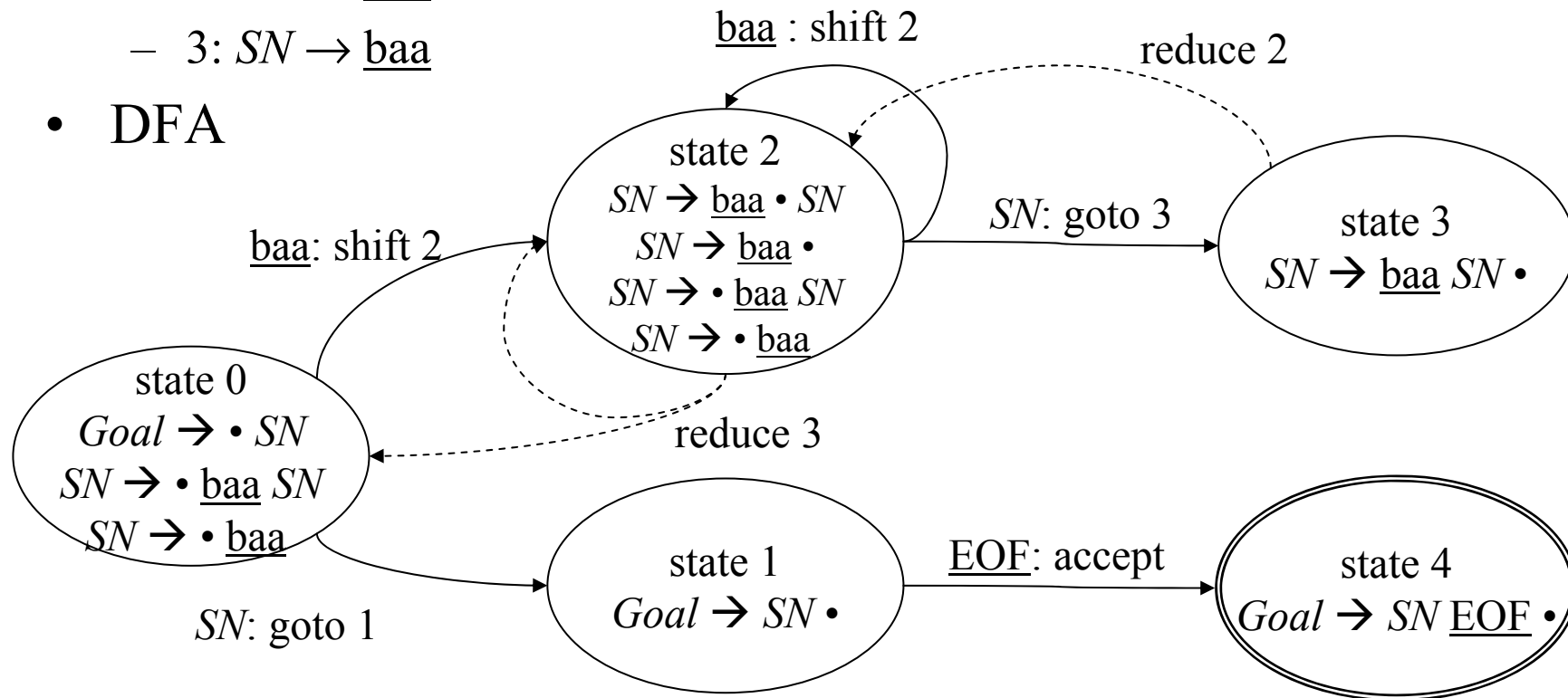
- **action table**
  - $(state, token)$   
→ shift  $state_{new}$
  - $(state, token)$   
→ reduce by the rule  $number$ 
    - reduce 후의  $state_{new}$  는 ?
  - $(state, EOF)$   
→ accept
  - undefined : error
- **goto table**
  - $(state, non-terminal)$   
→ go to  $state_{new}$ 
    - non-terminal은 reduce 후  
에 생긴 것

# Example : SheepNoise

- augmented grammar

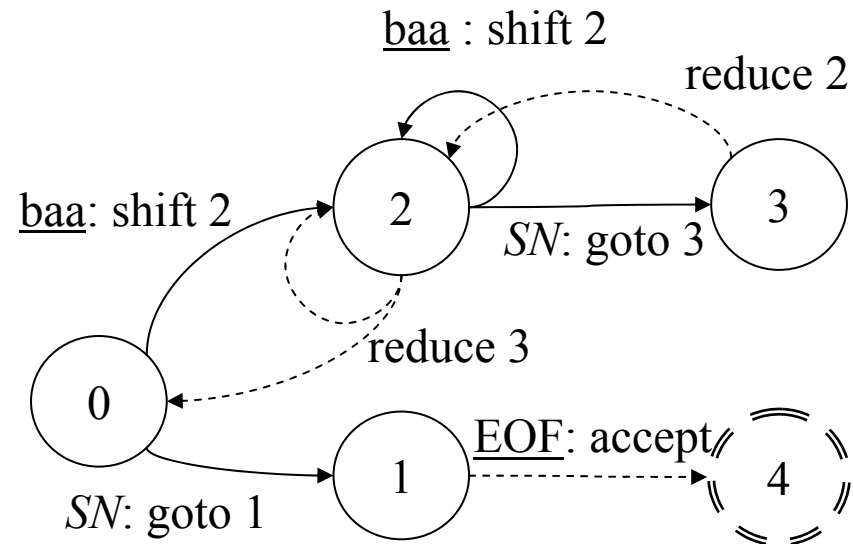
- 1:  $Goal \rightarrow SN$
- 2:  $SN \rightarrow \underline{baa} SN$
- 3:  $SN \rightarrow \underline{baa}$

- DFA



# Example : SheepNoise

- two tables



ACTION		
State	EOF	<u>baa</u>
0	—	shift 2
1	accept	—
2	reduce 3	shift 2
3	reduce 2	—

GOTO	
State	<i>SheepNoise</i>
0	1
1	—
2	3
3	—

# Example: "baa"

- input string "baa"

Stack	Input	Action
\$ 0	<u>baa</u> <u>EOF</u>	shift 2
\$ 0 <u>baa</u> 2	<u>EOF</u>	reduce 3
\$ 0 <i>SN</i> 1	<u>EOF</u>	accept

ACTION		
State	EOF	<u>baa</u>
0	—	shift 2
1	accept	—
2	reduce 3	shift 2
3	reduce 2	—

GOTO	
State	<i>SheepNoise</i>
0	1
1	—
2	3
3	—

# Example: "baa baa"

- input string "baa baa"

Stack	Input	Action
\$ 0	<u>b</u> aa <u>b</u> aa <u>EOF</u>	shift 2
\$ 0 <u>baa</u> 2	<u>b</u> aa <u>EOF</u>	shift 2
\$ 0 <u>baa</u> 2 <u>baa</u> 2	<u>EOF</u>	reduce 3
\$ 0 <u>baa</u> 2 <i>SN</i> 3	<u>EOF</u>	reduce 2
\$ 0 <i>SN</i> 1	<u>EOF</u>	accept

ACTION		
State	EOF	<u>baa</u>
0	—	shift 2
1	accept	—
2	reduce 3	shift 2
3	reduce 2	—

GOTO	
State	<i>SheepNoise</i>
0	1
1	—
2	3
3	—



# LR(1) Parser Again

- **parser driver routine** : see page 116 for more details  
push \$, 0 to the stack  
 $token \leftarrow$  from scanner  
**while** (true)  
     $state \leftarrow \text{top}(\text{stack})$   
    **if** Action[ $state, token$ ] = "shift  $i$ " **then** push token,  $i$ ;  
    **else if** Action[ $state, token$ ] = "reduce  $A \rightarrow \beta$ " **then**  
        pop  $2 * |\beta|$  symbols;  $state \leftarrow \text{top}(\text{stack})$ ;  
        push  $A$ , Goto[ $state, A$ ]  
    **else if** Action[ $state, token$ ] = "accept" **then** return  
    **else error endif**
- **remaining thing ?**
  - how to (automatically) construct the tables ?

## **3.5 Building LR(1) Tables**

# LR(1) Items

- a pair  $[A \rightarrow \beta \bullet \gamma, a]$ 
  - $A \rightarrow \beta \bullet \gamma$ : a handle or a potential handle
    - Caution:  $\beta, \gamma \equiv$  a sequence of symbols
  - $a \in T$ : look-ahead token : 모든 token을 cover해야 함!
  - handle-recognition DFA의 basic item
- possibility:  $[A \rightarrow \bullet \beta \gamma, a]$
- partially complete:  $[A \rightarrow \beta \bullet \gamma, a]$
- complete:  $[A \rightarrow \beta \gamma \bullet, a]$

# Example: Sheep Noise

- augmented grammar

1:  $Goal \rightarrow SN$

2:  $SN \rightarrow \underline{baa} SN$

3:  $SN \rightarrow \underline{baa}$

- all LR(1) items : 실제로 다 쓰이지는 않음...

$[Goal \rightarrow \bullet SN, EOF]$

$[Goal \rightarrow SN \bullet, EOF]$

$[SN \rightarrow \bullet \underline{baa}, EOF]$

$[SN \rightarrow \underline{baa} \bullet, EOF]$

$[SN \rightarrow \bullet \underline{baa}, \underline{baa}]$

$[SN \rightarrow \underline{baa} \bullet, \underline{baa}]$

$[SN \rightarrow \bullet \underline{baa} SN, EOF]$

$[SN \rightarrow \underline{baa} \bullet SN, EOF]$

$[SN \rightarrow \underline{baa} SN \bullet, EOF]$

$[SN \rightarrow \bullet \underline{baa} SN, \underline{baa}]$

$[SN \rightarrow \underline{baa} \bullet SN, \underline{baa}]$

$[SN \rightarrow \underline{baa} SN \bullet, \underline{baa}]$

# Closure(s) and Goto(s,x)

- closure(s) : DFA state 만들기
  - 같은 상황인 LR(1) item들  
→ DFA state

closure( $s$  : a set of states)

while ( $s$  is still changing)

$\forall$  item  $[A \rightarrow \beta \cdot C \delta, a] \in s$

$\forall$  rule  $C \rightarrow \gamma$

$\forall b \in \text{FIRST}(\delta a)$

$s \leftarrow s \cup \{ [C \rightarrow \cdot \gamma b] \}$

- Example: 아래 3개는 같은 상황
  - $[Goal \rightarrow \cdot SN, \text{EOF}]$
  - $[SN \rightarrow \cdot \underline{baa}, \text{EOF}]$
  - $[SN \rightarrow \cdot \underline{baa} SN, \text{EOF}]$

- goto( $s, x$ ) : DFA transition
  - 어디로 가야 하나

goto( $s, x$  : an NT or T)

$moved \leftarrow \emptyset$

$\forall$  item  $i \in s$

if  $i$  is  $[A \rightarrow \beta \cdot x \delta, a]$  then

$moved \leftarrow moved$

$\cup \{ [A \rightarrow \beta x \cdot \delta, a] \}$

return closure( $moved$ )

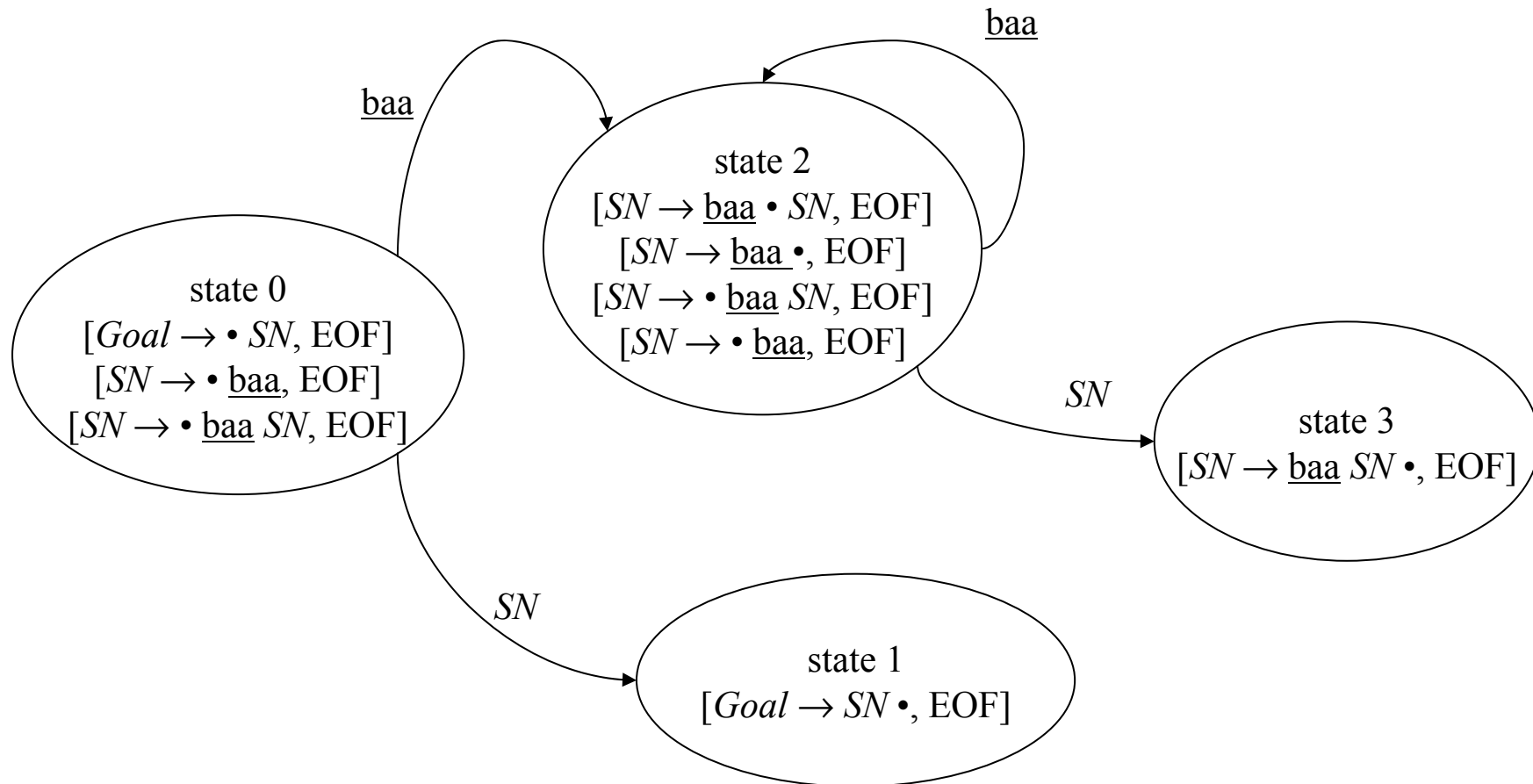
# Canonical Collection Construction

- $CC$ : canonical collection of sets of LR(1) items
  - final parser에서 하나의 state가 되는 items
  - $CC = \{ CC_0, CC_1, CC_2, \dots \}$
- $CC$  construction algorithm
  - $CC \leftarrow \{ CC_0 = \text{closure}(\{ [S' \rightarrow \bullet S, \text{EOF}] \}) \}$
  - while ( $CC$  is still changing)
    - $\forall$  set  $CC_j \in CC$
    - $\forall x$  following a  $\bullet$  in an item in  $CC_j$
    - $CC_k \leftarrow \text{goto}(CC_j, x)$  // 내부에서  $\text{closure}()$  계산 !
    - if  $CC_k \notin CC$  then add it to  $CC$
    - record the transition  $CC_j \rightarrow CC_k$  on  $x$

# Example: SheepNoise

- $CC_0 = \text{closure}([Goal \rightarrow \bullet SN, EOF])$   
 $= \{ [Goal \rightarrow \bullet SN, EOF], [SN \rightarrow \bullet \underline{baa} SN, EOF], [SN \rightarrow \bullet \underline{baa}, EOF] \}$
- $CC_1 = \text{goto}(CC_0, SN)$   
 $= \{ [Goal \rightarrow SN \bullet, EOF] \}$
- $CC_2 = \text{goto}(CC_0, baa)$   
 $= \{ [SN \rightarrow \underline{baa} \bullet SN, EOF], [SN \rightarrow \underline{baa} \bullet, EOF],$   
 $[SN \rightarrow \bullet \underline{baa} SN, EOF], [SN \rightarrow \bullet \underline{baa}, EOF] \}$
- $CC_3 = \text{goto}(CC_2, SN)$   
 $= \{ [SN \rightarrow \underline{baa} SN \bullet, EOF] \}$

# Example: Final DFA





# From DFA to the Table

- item  $[A \rightarrow \beta \bullet C \ \delta, a]$  : make a **shift** entry in Action Table
  - on C, shift to the state with  $[A \rightarrow \beta C \bullet \delta, a]$
- item  $[A \rightarrow \beta \bullet, a]$  : make a **reduce** entry in Action Table
  - reduce by the rule  $A \rightarrow \beta$
- special item  $[S \rightarrow S' \bullet, \text{EOF}]$  : **accept** in Action Table
- making **Goto** Table
  - for each  $n \in NT$ ,
  - if  $\text{goto}(CC_j, n) = CC_k$  then  $\text{Goto}[j, n] \leftarrow k$

# Final Table

- two table form

ACTION		
State	EOF	<u>baa</u>
0	—	shift 2
1	accept	—
2	reduce 3	shift 2
3	reduce 2	

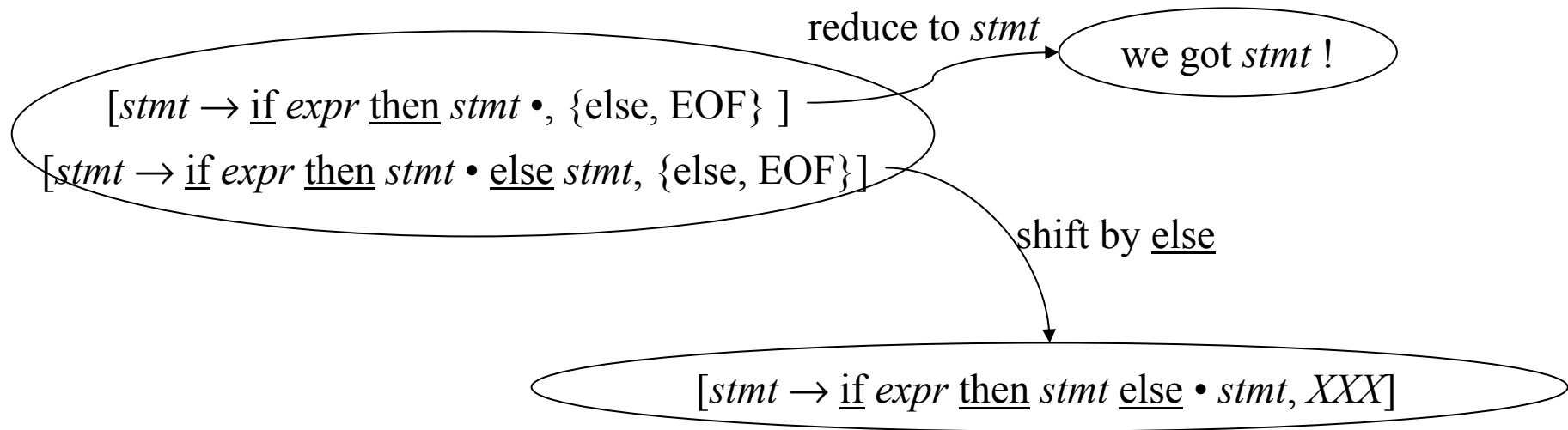
GOTO	
State	<i>SheepNoise</i>
0	1
1	—
2	3
3	—

- single table form

state	<u>EOF</u>	<u>baa</u>	<i>SN</i>
0		shift 2	goto 1
1	accept		
2	reduce 3	shift 2	goto 3
3	reduce 2		

# Shift-Reduce Conflict

- ambiguous if-then-else grammar
  - $stmt \rightarrow \underline{if} \ expr \ \underline{then} \ stmt$
  - $\quad \quad | \ \underline{if} \ expr \ \underline{then} \ stmt \ \underline{else} \ stmt$
  - $\quad \quad | \ assign$
- **CC** 생성 중에, 다음 상황 발생
  - Action table에 **conflict** 발생
- 해결책 ?
  - change your grammar
  - choose shift always
    - yacc의 해법



# Reduce-Reduce Conflict

- in Fortran and Ada,

$Factor \rightarrow FuncCall$

|  $ArrayRef$

| others...

$FuncCall \rightarrow \underline{iden} ( expr )$

$ArrayRef \rightarrow \underline{iden} ( expr )$

- CC** 생성 중에, 다음 상황 발생

- 해결책 ?

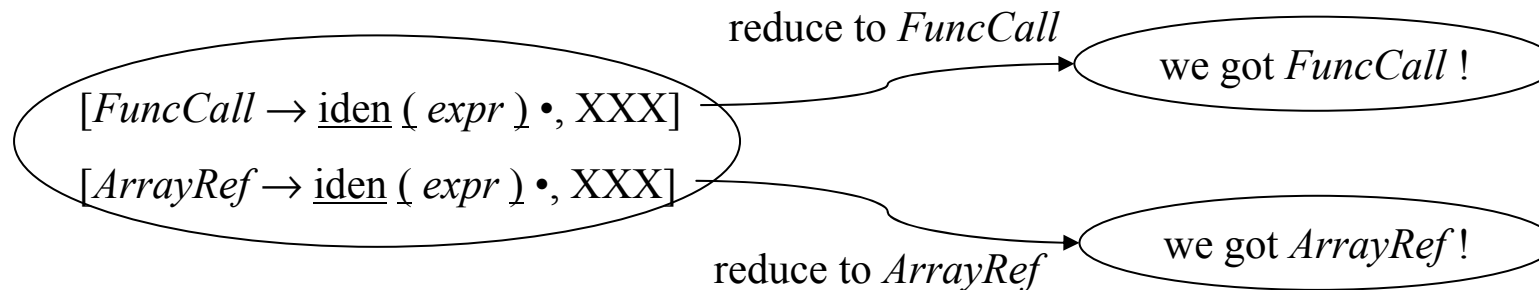
– change your language...

– scanner + symbol table 처리

- symbol table은  $iden$ 이  
function인지, array인지 알.

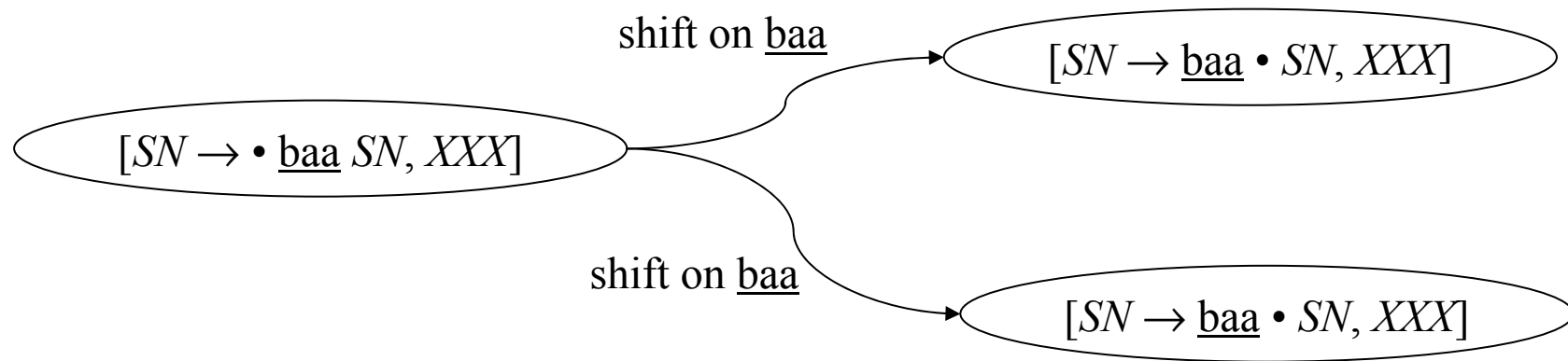
$FuncCall \rightarrow \underline{funcname} ( expr )$

$ArrayRef \rightarrow \underline{arrayname} ( expr )$



# Shift-Shift Conflict

- shift-shift conflict
  - **never occurs !**
  - why? we use DFA rather than NFA.



# LL(1) vs. LR(1)

	Advantages	Disadvantages
Top-down recursive descent	Fast Good locality Simplicity <b>Good error detection</b>	Hand-coded High maintenance <b>Right-Recursive Grammar</b>
LR(1)	Fast Deterministic langs. Automatable <b>Left or Right-Recursive Grammar</b>	Large working sets Poor error messages Large table sizes

## **3.6 Practical Issues**

# Error Recovery

- LL(1) 또는 LR(1) parser의 교과서적 구현  
→ single error detection and stop
- practical parser
  - **multiple error detection** 필요
  - how?
    - C/C++/Java 등: 모든 문장은 ';' 으로 끝
    - error 발생 시의 처리
      - scanner : ';' 까지 계속 eating
      - parser : error가 발생한 stmt 직전까지 pop



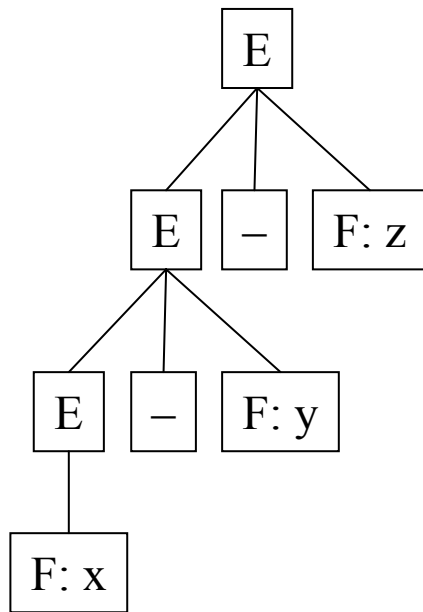
# Left and Right Recursion

- left-recursive grammar
  - LL(1), LR(1)에서 가능
- $list \rightarrow list \underline{elem}$   
           | elem
- example : elem elem elem
  - elem elem elem : shift
  - elem • elem elem : reduce
  - list • elem elem : shift
  - list elem • elem : reduce
  - list • elem : shift
  - list elem • : reduce
  - list • : accept
- 이쪽을 선호 !
- right recursive grammar
  - LR(1)에서 가능
- $list \rightarrow \underline{elem} list$   
           | elem
- example: elem elem elem
  - elem elem elem : shift
  - elem • elem elem : shift
  - elem elem • elem : shift
  - elem elem elem • : reduce
  - elem elem list • : reduce
  - elem list • : reduce
  - list • : accept
- **stack depth: 훨씬 깊어 진다 !**

# Associativity

- left-recursive grammar

- $expr \rightarrow expr '+' fact$   
 $\quad \quad \quad | \quad expr '-' fact$   
 $\quad \quad \quad | \quad fact$

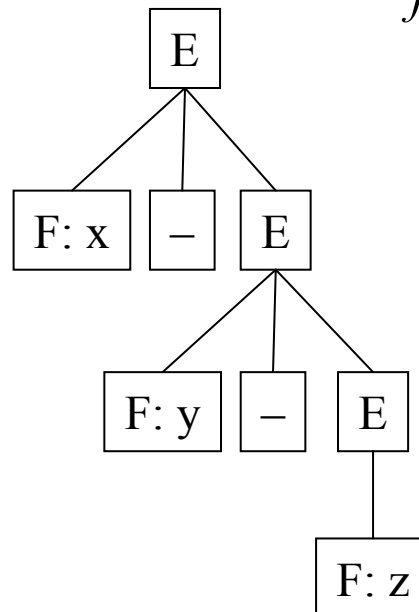


$(x - y) - z : \text{good !}$

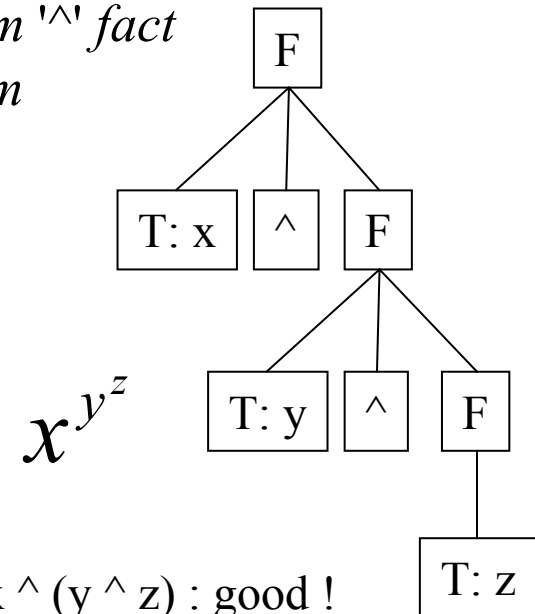
- right recursive grammar

- $expr \rightarrow fact '+' expr$   
 $\quad \quad \quad | \quad fact '-' expr$   
 $\quad \quad \quad | \quad fact$

- $fact \rightarrow term '^' fact$   
 $\quad \quad \quad | \quad term$



$x - (y - z) : \text{bad...}$



$x ^ (y ^ z) : \text{good !}$

## **3.7 Advanced Topics**

# Optimizing a Grammar

- original grammar

$$\begin{aligned} \text{Expr} &\rightarrow \text{Expr} \pm \text{Term} \\ &\quad | \text{Expr} = \text{Term} \\ &\quad | \text{Term} \end{aligned}$$
$$\begin{aligned} \text{Term} &\rightarrow \text{Term} * \text{Fact} \\ &\quad | \text{Term} / \text{Fact} \\ &\quad | \text{Fact} \end{aligned}$$
$$\begin{aligned} \text{Fact} &\rightarrow ( \text{Expr} ) \\ &\quad | \underline{\text{num}} \\ &\quad | \underline{\text{id}} \end{aligned}$$

- optimized grammar
  - parse tree depth 축소 효과

$$\begin{aligned} \text{Expr} &\rightarrow \text{Expr} \pm \text{Term} \\ &\quad | \text{Expr} = \text{Term} \\ &\quad | \text{Term} \end{aligned}$$
$$\begin{aligned} \text{Term} &\rightarrow \text{Term} * ( \text{Expr} ) \\ &\quad | \text{Term} * \underline{\text{num}} \\ &\quad | \text{Term} * \underline{\text{id}} \\ &\quad | \text{Term} / ( \text{Expr} ) \\ &\quad | \text{Term} / \underline{\text{num}} \\ &\quad | \text{Term} / \underline{\text{id}} \\ &\quad | ( \text{Expr} ) \\ &\quad | \underline{\text{num}} \\ &\quad | \underline{\text{id}} \end{aligned}$$

# Shrinking the Grammar

- original grammar
  - scanner의 return 형태:
    - $\langle +, \text{NULL} \rangle, \dots$
- reduced grammar
  - scanner의 return 형태:
    - $\langle \underline{\text{addop}}, + \rangle, \langle \underline{\text{mulop}}, * \rangle, \dots$
  - parser가 간단해짐

$$\begin{aligned} \text{Expr} &\rightarrow \text{Expr} \pm \text{Term} \\ &| \text{Expr} = \text{Term} \\ &| \text{Term} \end{aligned}$$

$$\begin{aligned} \text{Term} &\rightarrow \text{Term} * \text{Fact} \\ &| \text{Term} / \text{Fact} \\ &| \text{Fact} \end{aligned}$$

$$\begin{aligned} \text{Fact} &\rightarrow ( \text{Expr} ) \\ &| \underline{\text{num}} \\ &| \underline{\text{iden}} \end{aligned}$$

$$\begin{aligned} \text{Expr} &\rightarrow \text{Expr} \underline{\text{addop}} \text{Term} \\ &| \text{Term} \end{aligned}$$

$$\begin{aligned} \text{Term} &\rightarrow \text{Term} \underline{\text{mulop}} \text{Fact} \\ &| \text{Fact} \end{aligned}$$

$$\begin{aligned} \text{Fact} &\rightarrow ( \text{Expr} ) \\ &| \underline{\text{num}} \\ &| \underline{\text{iden}} \end{aligned}$$

# Other Construction Algorithms

- **LR(1)** : canonical LR(1) parsing
  - **SLR(1)** : simple LR(1) parsing
    - LR(1)에서 lookahead를 제거한 LR(1) item 사용
    - power가 떨어짐 → 그러나, 대부분 문제 없음
  - **LALR(1)** : lookahead LR(1) parsing
    - LR(1)에서, 모든 item이 아니라, 대표 item만 사용
    - 대부분의 programming language에서는 무난
    - yacc는 LALR(1) 사용
- **$LALR(1) \subset SLR(1) \subset LR(1) \subset CFG$**