

## Homework & Announcements

- Reading: Chapter
- Homework: Exercises at end
- Due: 11/5

Copyright © 2002–2018 UMaine Computer Science Department – 1 / 31

# COS 140: Foundations of Computer Science

## Specifying Programming Languages: Backus–Naur Form (BNF)

Fall 2018

<b>Introduction</b>	<b>3</b>
Problem . . . . .	3
Syntax and Semantics . . . . .	4
Syntax Formalisms . . . . .	5
Backus–Naur Form . . . . .	7
<b>Languages</b>	<b>8</b>
Language Terminology . . . . .	8
Chomsky Hierarchy . . . . .	9
<b>Backus-Naur Form</b>	<b>10</b>
BNF . . . . .	10
Example . . . . .	12
Derivation . . . . .	13
Recursion . . . . .	14
<b>Parsing</b>	<b>15</b>
Parsing . . . . .	15
Example . . . . .	16
Ambiguity . . . . .	17
Unambiguous grammar . . . . .	19
Precedence . . . . .	20
Example Grammar . . . . .	22
Parsing . . . . .	23
Example Parse . . . . .	24
Parsers . . . . .	31

## Problem

- Problem: how to *specify* a programming language?
- Have to have a way to describe its syntax (and, possibly, its semantics)

## Syntax and Semantics

- *Syntax*:
  - Form of something (e.g., language)
  - Describes relationships between components
- *Semantics*:
  - Meaning of something
  - Describes what the statements will do
- Need to capture formally so it's clear how language is to be used.

Copyright © 2002–2018 UMaine Computer Science Department – 4 / 31

## Syntax Formalisms

- What is needed in a syntax formalism?
  - Able to specify language's grammar for users and compiler designers
  - Easy to write and understand
  - Expressive enough to capture all programming languages
  - Easy to translate into algorithms for machine translation of language

Copyright © 2002–2018 UMaine Computer Science Department – 5 / 31

## Syntax Formalisms

- Formalisms often expressed as *production rules*:
  - $LHS \rightarrow RHS$
  - Match LHS with what you have, produce RHS
  - E.g.:  $S \rightarrow NP \ VP$
- As we'll see: can also think of going backwards
  - If you have NP and VP  $\Rightarrow$  have a valid sentence S.

Copyright © 2002–2018 UMaine Computer Science Department – 6 / 31

## Backus–Naur Form

- BNF is the most common way of specifying *grammar* of a programming language
- Also important for:
  - Computability theory
  - Natural language processing
  - Many other places in CS where you need to specify a grammar

Copyright © 2002–2018 UMaine Computer Science Department – 7 / 31

### Language Terminology

- *String*: any combination of characters in the language – may be the whole program
- *Lexemes*: lowest-level unit of language – analogous to words in English
- *Syntactic category*: classes of lexemes that play the same role in the structure of a statement – analogous to parts of speech
- *Tokens*: Low-level syntactic categories corresponding to the lexemes
- *Constituents*: tokens or groups of tokens that are put together in ways specified by the grammar – analogous to noun phrases, etc.
- *Grammar*: specification of the syntax; a set of rules describing the legal strings of the language
- *Terminal* and non-terminal symbols

Copyright © 2002–2018 UMaine Computer Science Department – 8 / 31

### Types of Languages: The Chomsky Hierarchy

- As go down the hierarchy, languages increase in complexity, more machinery needed to recognize them
- *Regular languages* (regular expressions) – tokens
  - Single symbol on the LHS; RHS has at most one non-terminal:
  - E.g.:  $S \rightarrow aS \mid b$  – generates  $a^+b$
- *Context-free languages* – programming languages
  - LHS: single symbol; terminals, non-terminals in RHS
  - E.g.:  $S \rightarrow aSb \mid \text{nil}$  – generates  $a^n b^n$
- *Context-sensitive languages*
  - RHS has at least as many symbols as LHS
  - E.g.:  $aSc \rightarrow aSbSc$
- *Recursively-enumerable languages* – Turing-equivalent representation
  - No restrictions on LHS, RHS
  - E.g.:  $aaaS \rightarrow aaaTQ$

Copyright © 2002–2018 UMaine Computer Science Department – 9 / 31

## Backus-Naur Form (BNF)

- Grammar formalism for context-free languages (only one symbol on LHS)
- Terminology:
  - *Rewrite (production) rules* – rules that rewrite current pattern into a new one
  - *LHS, RHS* – parts of rule: LHS  $\rightarrow$  RHS
  - *Terminal symbols* – symbols that appear only on RHS – i.e., do not get replaced by anything
  - *Non-terminal symbols* – symbols that appear in some LHS
  - *Alternatives (|)* – different RHS's that can be used with the LHS

Copyright © 2002–2018 UMaine Computer Science Department – 10 / 31

## Backus-Naur Form (BNF)

- E.g.:  $\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle \mid \langle \text{var} \rangle + \langle \text{var} \rangle$   
 $\langle \text{var} \rangle \rightarrow A \mid B \mid C$
- Often written with  $::=$  instead of  $\rightarrow$ :  
 $\langle \text{var} \rangle ::= A \mid B \mid C$

Copyright © 2002–2018 UMaine Computer Science Department – 11 / 31

### Example: Producing a String from BNF

$$\begin{aligned}\langle \text{expr} \rangle &\rightarrow \langle \text{var} \rangle \mid \langle \text{var} \rangle + \langle \text{var} \rangle \\ \langle \text{var} \rangle &\rightarrow A \mid B \mid C\end{aligned}$$

- Start with the kind of constituent you want to produce

$\langle \text{expr} \rangle$

- Replace each non-terminal in LHS with a RHS that appears in the rule for which that non-terminal is in the LHS

$\langle \text{var} \rangle + \langle \text{var} \rangle$

- Keep doing this until there are only non-terminals in the string

$A + \langle \text{var} \rangle \Rightarrow \dots \Rightarrow A + B$

Copyright © 2002–2018 UMaine Computer Science Department – 12 / 31

### Derivation

- *Start symbol*: where derivations begin to derive all possible strings (programs)
- *Sentential form*: any string derived from the start symbols using the rewrite rules
- *Derivation*: rewrite sentential forms until have all terminal symbols
- Options:
  - Order of rewriting – left-most, right-most derivations
  - Alternative used for rewriting

Copyright © 2002–2018 UMaine Computer Science Department – 13 / 31



## Recursion

- In BNF: LHS appears in the RHS
- Allows infinite language from finite grammar
- Cannot specify the number of times rule will be applied
- E.g.: to allow any number of additions in an expression:  
$$\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle \mid \langle \text{var} \rangle + \langle \text{expr} \rangle$$
- Need an alternative that can stop the recursion!

Copyright © 2002–2018 UMaine Computer Science Department – 14 / 31

## Parsing

15 / 31

### Parsing

- What is *parsing*?
  - Checking the legality of input against a grammar
  - Producing a *parse tree* that shows the relationships between the tokens
  - Parse trees record derivations
- Start with a string of only tokens (terminal symbols)
- Find a RHS pattern in the string, replace with the LHS.
- Continue until you have the start symbol
- If you can do this: the input was a valid sentence in the language
- Create a parse tree by showing how we replace terminal/non-terminal symbols using appropriate rules

Copyright © 2002–2018 UMaine Computer Science Department – 15 / 31

## Parsing Example

Input: A + B + C      Grammar:  $\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle \mid \langle \text{var} \rangle + \langle \text{expr} \rangle$   
 Processed so far:       $\langle \text{var} \rangle \rightarrow A \mid B \mid C$   
 Parse tree:

(1)

Input: A + B + C      Grammar:  $\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle \mid \langle \text{var} \rangle + \langle \text{expr} \rangle$   
 Processed so far: A       $\langle \text{var} \rangle \rightarrow A \mid B \mid C$   
 Parse tree:

A

(2)

Input: A + B + C      Grammar:  $\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle \mid \langle \text{var} \rangle + \langle \text{expr} \rangle$   
 Processed so far: A       $\langle \text{var} \rangle \rightarrow A \mid B \mid C$   
 Parse tree:

```

    <var>
    |
    A
    
```

(3)

Input: A + B + C      Grammar:  $\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle \mid \langle \text{var} \rangle + \langle \text{expr} \rangle$   
 Processed so far: A +       $\langle \text{var} \rangle \rightarrow A \mid B \mid C$   
 Parse tree:

```

    <expr>
   / | \
  <var> + <expr>
   |
   A
    
```

(4)

Input: A + B + C      Grammar:  $\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle \mid \langle \text{var} \rangle + \langle \text{expr} \rangle$   
 Processed so far: A + B + C       $\langle \text{var} \rangle \rightarrow A \mid B \mid C$   
 Parse tree:

```

    <expr>
   / | \
  <var> + <expr>
   |   / | \
   A  <var> + <expr>
       |   |
       B  <var>
           |
           C
    
```

(5)

Copyright © 2002–2018 UMaine Computer Science Department – 16 / 31

## Ambiguity

- More than one different, legal, correct parse trees for the same string
- Problem, since parse tree indicates relationship between the constituents
- Property of the grammar

Copyright © 2002–2018 UMaine Computer Science Department – 17 / 31

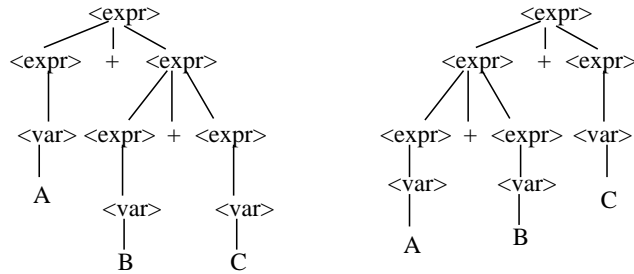
## Ambiguity Example

- Grammar:

$$\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle \mid \langle \text{expr} \rangle + \langle \text{expr} \rangle$$

$$\langle \text{var} \rangle \rightarrow A \mid B \mid C$$

- Parse trees:

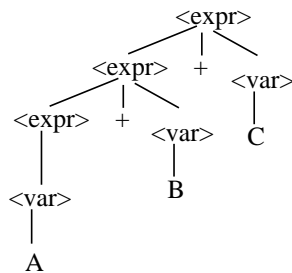


- Sub-expressions that are more deeply nested in the tree are evaluated first

Copyright © 2002–2018 UMaine Computer Science Department – 18 / 31

## An Unambiguous Grammar for Addition Expressions

- Need to ensure that the expression can be evaluated in only one way
- Want the expression to be evaluated left to right  $\Rightarrow$  first expression formed (lowest in tree) needs to be to the left
- Replace rule in previous with:

$$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{var} \rangle$$


Copyright © 2002–2018 UMaine Computer Science Department – 19 / 31

## Operator Precedence

- Need a constituent for each set of operators with the same precedence
- Operators with the highest precedence need to be built at the lowest level in the tree
- Get left associativity if we parse left  $\rightarrow$  right, keep recursion to the left  
 $\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{var} \rangle$
- Right associativity: recursion to the left.

Copyright © 2002–2018 UMaine Computer Science Department – 20 / 31

## Levels of Precedence

- Items in parentheses and identifiers, numbers: *factor*  
A, B, (C + D)
- Multiplication and division: *term*  
A\*B, A/B
- Addition and subtraction: *expression*  
A+B, A-B

Copyright © 2002–2018 UMaine Computer Science Department – 21 / 31

## Example Grammar

```
<expr> → <expr> + <term> |  
        <expr> - <term> |  
        <term>  
<term> → <term> * <factor> |  
        <term> / <factor> |  
        <factor>  
<factor> → (<expr>) | <id>  
<id> → A | B | C | D
```

Copyright © 2002–2018 UMaine Computer Science Department – 22 / 31

## Parsing

- Can think of there being a “state” of the parse
  - Records where we are in the process
  - E.g.:  $\langle \text{expr} \rangle + \langle \text{term} \rangle$  would mean that we’ve turned the tokens we’ve read so far into the non-terminal  $\langle \text{expr} \rangle$ , the terminal symbol (token)  $+$ , and the non-terminal  $\langle \text{term} \rangle$
- Idea of state can be used to parse using a *state machine* or automata – we’ll talk about this later in course
- Can also think of there being some tokens that have not yet been read
- Start with an empty state, and with all terminals unread

Copyright © 2002–2018 UMaine Computer Science Department – 23 / 31

# Parse of A + B \* C - D

State:  
(empty)  
Input left: A + B \* C - D

No grammar rules apply to empty state.  
Read next token (A).

A + B \* C - D

State:  
A  
Input left: + B \* C - D

Only rule that applies is:  
<id> ::= A

A + B \* C - D

State:  
<id>  
Input left: + B \* C - D

Only rule that applies is:  
<factor> ::= <id>

<id>  
A + B \* C - D

State:  
<factor>  
Input left: + B \* C - D

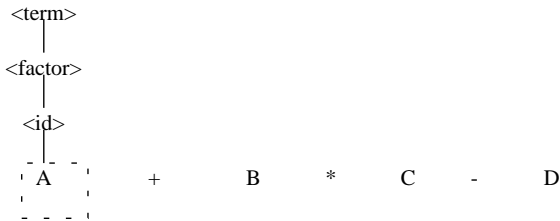
Only rule that applies is:  
<term> ::= <factor>

<factor>  
|  
<id>  
|  
A + B \* C - D

# Parse of $A + B * C - D$

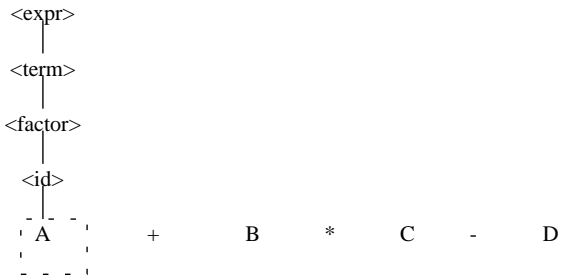
State:  
 <term>  
 Input left: + B \* C - D

Only rule that applies is:  
 <expr> ::= <term>



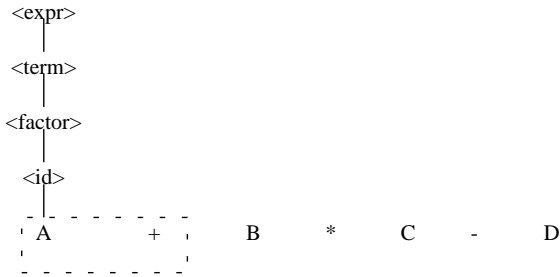
State:  
 <expr>  
 Input left: + B \* C - D

Although <expr> is the current state, there are still input tokens left, so we're not done.  
 Read next token (+).



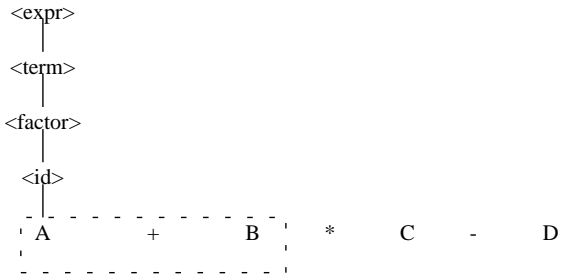
State:  
 <expr> +  
 Input left: B \* C - D

No rule applies.  
 Read next token (B).



State:  
 <expr> + B  
 Input left: \* C - D

Can't replace whole thing.  
 Use <id> ::= B



# Parse of $A + B * C - D$

State:  
 $\langle \text{expr} \rangle + \langle \text{id} \rangle$   
 Input left: \* C - D

Can't replace whole thing.  
 Use  $\langle \text{factor} \rangle ::= \langle \text{id} \rangle$

State:  
 $\langle \text{expr} \rangle + \langle \text{factor} \rangle$   
 Input left: \* C - D

Can't replace whole thing.  
 Use  $\langle \text{term} \rangle ::= \langle \text{factor} \rangle$

State:  
 $\langle \text{expr} \rangle + \langle \text{term} \rangle$   
 Input left: \* C - D

Can apply  $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{term} \rangle$

State:  
 $\langle \text{expr} \rangle$   
 Input left: \* C - D

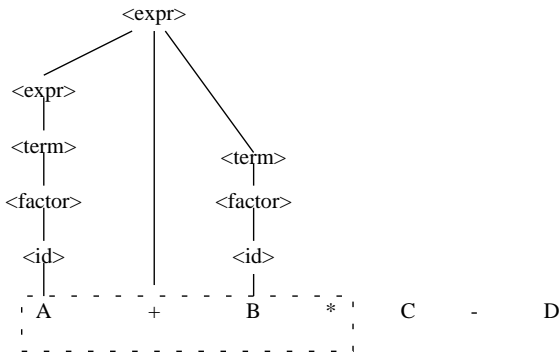
Can't stop -- still have input tokens left.  
 Read next token (\*)



# Parse of A + B \* C - D

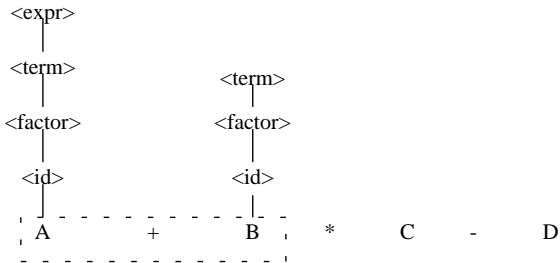
State:  
 <expr> \*  
 Input left: C - D

No grammar rule produces anything beginning with <expr> \* -- dead end.  
 Backtrack to previous state



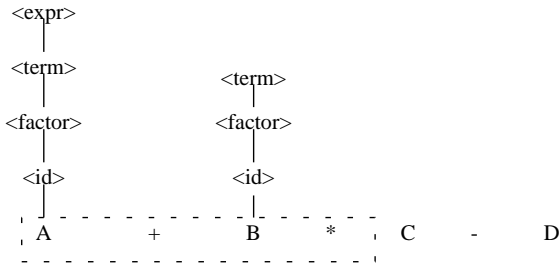
State:  
 <expr> + <term>  
 Input left: \* C - D

Although we could apply <expr> ::= <expr> + <term>, we've already tried that with no luck -- so read next input token (\*).



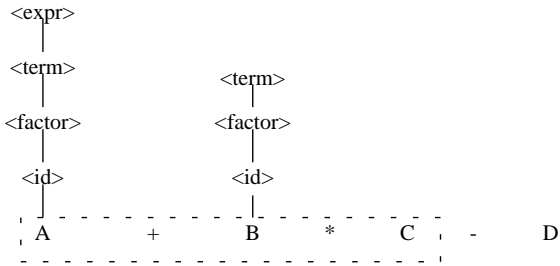
State:  
 <expr> + <term> \*  
 Input left: C - D

Nothing matches directly with <term> \*  
 Read next token (C).



State:  
 <expr> + <term> \* C  
 Input left: - D

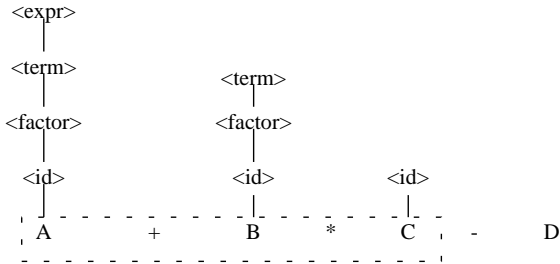
Nothing matches directly with <term> \* C  
 Use grammar rule: <id> ::= C



# Parse of $A + B * C - D$

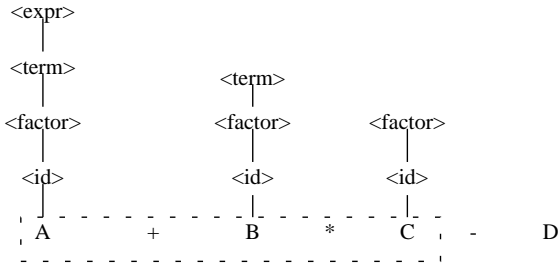
State:  
 $\langle \text{expr} \rangle + \langle \text{term} \rangle * \langle \text{id} \rangle$   
 Input left: - D

Nothing matches directly with  $\langle \text{term} \rangle * \langle \text{id} \rangle$   
 Use grammar rule:  $\langle \text{factor} \rangle ::= \langle \text{id} \rangle$



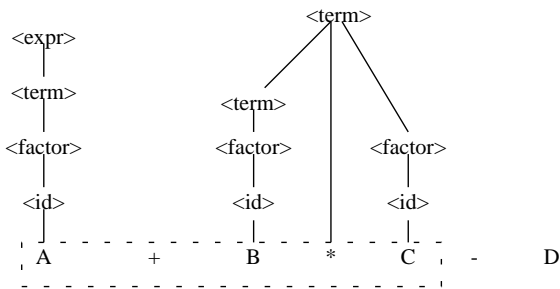
State:  
 $\langle \text{expr} \rangle + \langle \text{term} \rangle * \langle \text{factor} \rangle$   
 Input left: - D

Use grammar rule:  
 $\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{factor} \rangle$



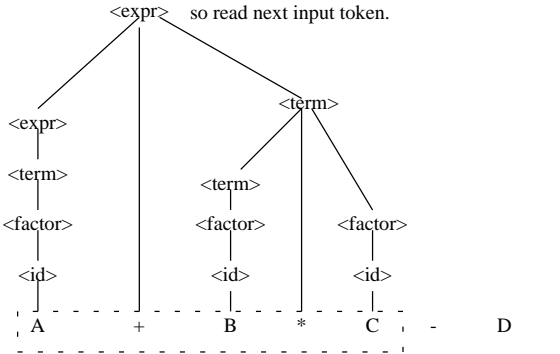
State:  
 $\langle \text{expr} \rangle + \langle \text{term} \rangle$   
 Input left: - D

Use grammar rule:  
 $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{term} \rangle$

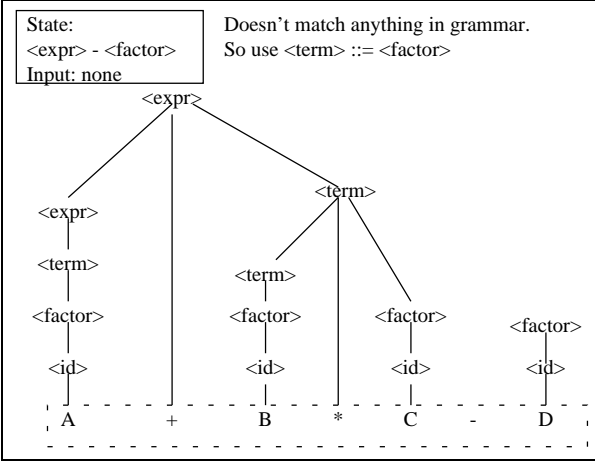
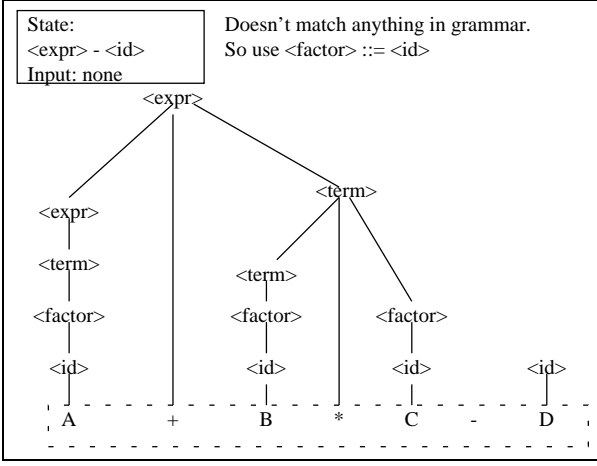
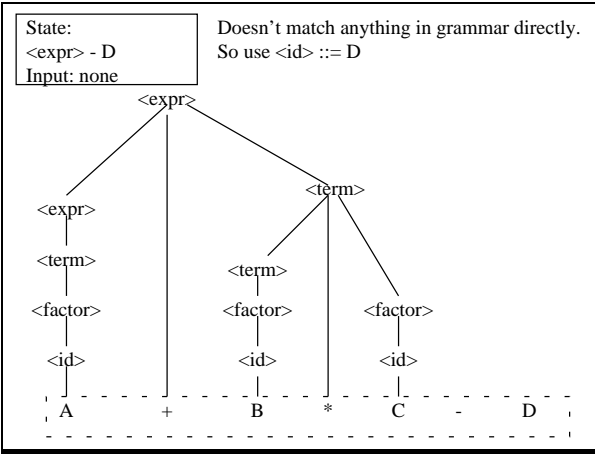
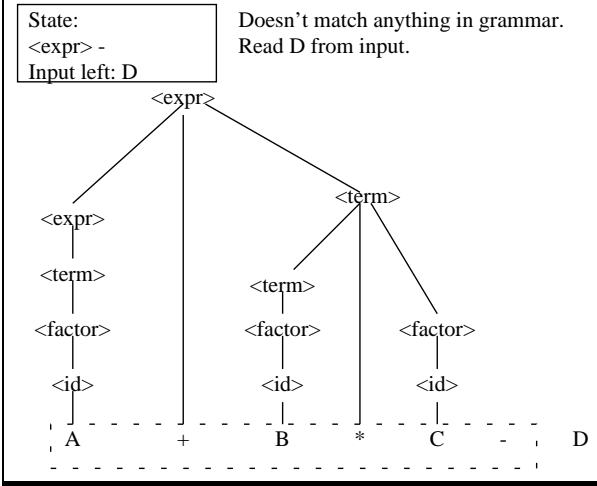


State:  
 $\langle \text{expr} \rangle$   
 Input left: - D

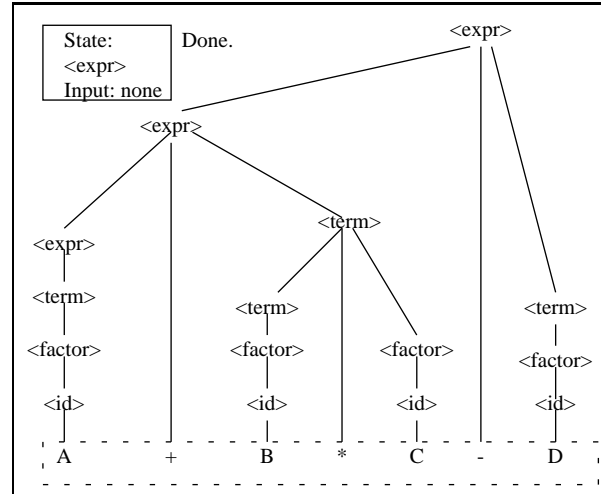
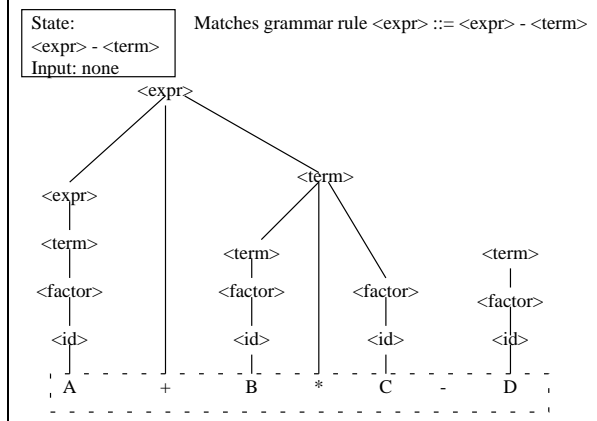
Although  $\langle \text{expr} \rangle$  is state, there is still input, so we're not done.  
 No matching rules in grammar; so read next input token.



# Parse of $A + B * C - D$



## Parse of $A + B * C - D$



Copyright © 2002–2018 UMaine Computer Science Department – 30 / 31

## Parsers

- Kinds of grammars/kinds of parsing:
  - LL: left-to-right, leftmost derivation
  - LR: left-to-right, rightmost derivation
  - Different amounts of *lookahead*: LR(1), e.g.
- LL parsing: e.g., *recursive-descent parsers*
- LR parsing: e.g., *shift-reduce parsers* – typically table-driven

Copyright © 2002–2018 UMaine Computer Science Department – 31 / 31