# Parsing — Part II
*(Ambiguity, Top-down parsing, Left-recursion Removal)*

## *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

Classic example — the *if-then-else* problem

> *Stmt* ® <u>if</u>  *Expr*  <u>then</u> *Stmt*
>
> | <u>if</u>  *Expr*  <u>then</u> *Stmt*  <u>else</u>  *Stmt*
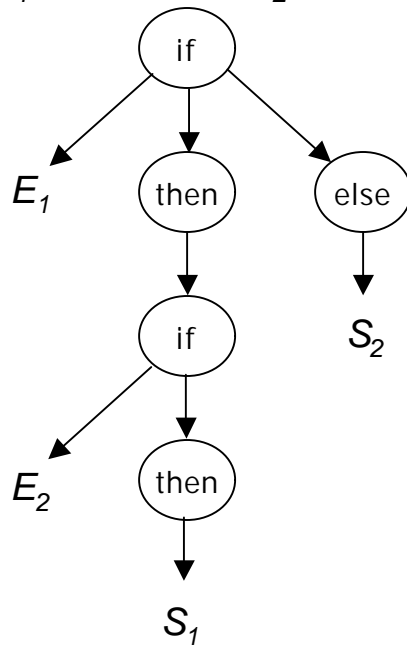>
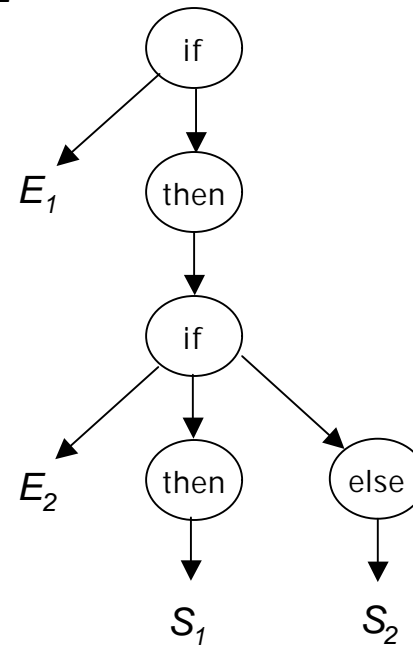> |  *… other stmts …*

*This ambiguity is entirely grammatical in nature*

## *Ambiguity*

This sentential form has two derivations (If a derivation has more than 1 parse tree, the grammar is ambiguous

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



production 2, then production 1

production 1, then production 2

## *Ambiguity*

Removing the ambiguity

- Must rewrite the grammar to avoid generating the problem

- Match each <u>else</u> to innermost unmatched <u>if</u>     *(common sense rule)*

| | | |
|---|---|---|
| 1 | *Stmt* → | *With Else* |
| 2 | \| | *NoElse* |
| 3 | *With Else* → | <u>If</u> *Expr* <u>then</u> *With Else* <u>else</u> *With Else* |
| 4 | \| | … other stmts … |
| 5 | *NoElse* → | <u>If</u> *Expr* <u>then</u> *Stmt* |
| 6 | \| | <u>If</u> *Expr* <u>then</u> *With Else* <u>else</u> *NoElse* |

With this grammar, the example has only one derivation

## *Ambiguity*

if *Expr$_1$* <u>then</u> <u>if</u> *Expr$_2$* <u>then</u> *Stmt$_1$* <u>else</u> *Stmt$_2$*

| Rule | Sentential Form |
|------|-----------------|
| — | Stmt |
| 2 | NoElse |
| 5 | <u>if</u> Expr <u>then</u> Stmt |
| ? | <u>if</u> E$_1$ <u>then</u> Stmt |
| 1 | <u>if</u> E$_1$ <u>then</u> WithElse |
| 3 | <u>if</u> E$_1$ <u>then</u> <u>if</u> Expr <u>then</u> WithElse <u>else</u> WithElse |
| ? | <u>if</u> E$_1$ <u>then</u> <u>if</u> E$_2$ <u>then</u> WithElse <u>else</u> WithElse |
| 4 | <u>if</u> E$_1$ <u>then</u> <u>if</u> E$_2$ <u>then</u> S$_1$ <u>else</u> WithElse |
| 4 | <u>if</u> E$_1$ <u>then</u> <u>if</u> E$_2$ <u>then</u> S$_1$ <u>else</u> S$_2$ |

This binds the <u>else</u> controlling *S$_2$* to the inner <u>if</u>

from Cooper & Torczon

## *Deeper Ambiguity*

Ambiguity usually refers to confusion in the CFG

Overloading can create deeper ambiguity

    a = f(17)

In some languages, f could be either a function or a subscripted variable

Disambiguating this one requires context

- Need values of declarations

- Really an issue of *type*, not context-free syntax

- Requires an extra-grammatical solution (not in CFG)

- Must to handle these with a different mechanism

  > Step outside grammar rather than use a more complex grammar

## *Ambiguity - the Final Word*

Ambiguity arises from two distinct sources

- Confusion in the context-free syntax        (<u>if</u>-<u>then</u>-<u>else</u>)

- Confusion that requires context to resolve     (overloading)


Resolving ambiguity

- To remove context-free ambiguity, rewrite the grammar

- To handle context-sensitive ambiguity takes cooperation

  - > Knowledge of declarations, types, ...

  - > Accept a superset of *L(G)* & check it with other means[†]

  - > This is a language design problem


Sometimes, the compiler writer accepts an ambiguous grammar

  - > Parsing techniques that "do the right thing"


[†]See Chapter 4

from Cooper & Torczon

## 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

from Cooper & Torczon

## Top-down Parsing

*A top-down parser starts with the root of the parse tree*

*The root node is labeled with the goal symbol of the grammar*

*Top-down parsing algorithm:*

*Construct the root node of the parse tree*

*Repeat until the fringe of the parse tree matches the input string*

1 *At a node labeled A, select a production with A on its lhs and, for each symbol on its rhs, construct the appropriate child*

2 *When a terminal symbol is added to the fringe and it doesn't match the fringe, backtrack*

3 *Find the next node to be expanded*          *(label Î NT)*

The key is picking the right production in step 1

> *That choice should be guided by the input string*

from Cooper & Torczon

# *Remember the expression grammar?*

Version with precedence derived last lecture

| 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> |

*And the input* <u>x</u> - <u>2</u> * <u>y</u>

# *Example*

Let's try <u>x</u> - <u>2</u> * <u>y</u> :

| Rule | Sentential Form | Input |
|------|-----------------|-------|
| — | *Goal* | ↑<u>x</u> - <u>2</u> * <u>y</u> |
| 1 | *Expr* | ↑<u>x</u> - <u>2</u> * <u>y</u> |
| 2 | *Expr + Term* | ↑<u>x</u> - <u>2</u> * <u>y</u> |
| 4 | *Term + Term* | ↑<u>x</u> - <u>2</u> * <u>y</u> |
| 7 | *Factor + Term* | ↑<u>x</u> - <u>2</u> * <u>y</u> |
| 9 | <id, x> + *Term* | ↑<u>x</u> - <u>2</u> * <u>y</u> |
| 9 | <id, x> + *Term* | <u>x</u> ↑- <u>2</u> * <u>y</u> |

## Example

Let's try x - 2 * y :

| Rule | Sentential Form | Input |
|------|-----------------|-------|
| — | Goal | ↑x - 2 * y |
| 1 | Expr | ↑x - 2 * y |
| 2 | Expr + Term | ↑x - 2 * y |
| 4 | Term + Term | ↑x - 2 * y |
| 7 | Factor + Term | ↑x - 2 * y |
| 9 | \<id, x\>+ Term | ↑x - 2 * y |
| 9 | \<id, x\>+ Term | x ↑- 2 * y |

This worked well, except that "-" doesn't match "+"

The parser must backtrack to here

# *Example*

Continuing with <u>x</u> - <u>2</u> * <u>y</u> :

| *Rule* | *Sentential Form* | *Input* |
|--------|-------------------|---------|
| — | *Goal* | ↑<u>x</u> - <u>2</u> * <u>y</u> |
| 1 | *Expr* | ↑<u>x</u> - <u>2</u> * <u>y</u> |
| 3 | *Expr - Term* | ↑<u>x</u> - <u>2</u> * <u>y</u> |
| 4 | *Term - Term* | ↑<u>x</u> - <u>2</u> * <u>y</u> |
| 7 | *Factor - Term* | ↑<u>x</u> - <u>2</u> * <u>y</u> |
| 9 | <id, x> *- Term* | ↑<u>x</u> - <u>2</u> * <u>y</u> |
| 9 | <id, x> *- Term* | <u>x</u> ↑- <u>2</u> * <u>y</u> |
| — | <id, x> *- Term* | <u>x</u> - ↑<u>2</u> * <u>y</u> |

## *Example*

Continuing with x - 2 * y :

| Rule | Sentential Form | Input |
|---|---|---|
| — | Goal | ↑x - 2 * y |
| 1 | Expr | ↑x - 2 * y |
| 3 | Expr - Term | ↑x - 2 * y |
| 4 | Term - Term | ↑x - 2 * y |
| 7 | Factor - Term | ↑x - 2 * y |
| 9 | <id, x> - Term | ↑x - 2 * y |
| 9 | <id, x> - Term | x ↑- 2 * y |
| — | <id, x> - Term | x - ↑2 * y |

This time, "-" and "-" matched

We can advance past "-" to look at "2"

⇒ Now, we need to expand *Term* - the last *NT* on the fringe

## *Example*

Trying to match the "2" in x - 2 * y :

| *R u le* | *S en te n tial  For m* | *In pu t* |
|---|---|---|
| — | \<id, x\> - *Ter m* | x - ↑2 * y |
| 7 | \<id, x\> - *Fact  o r* | x - ↑2 * y |
| 9 | \<id, x\> - \<n um ,2\> | x - ↑2 * y |
| — | \<id, x\> - \<n um ,2\> | x - 2 ↑* y |

Goal → Expr → (Expr, -, Term); Expr → Term → Fact. → \<id,x\>; Term → Fact. → \<num,2\>

## *Example*

Trying to match the "2" in <u>x</u> - <u>2</u> * <u>y</u> :

| Rule | Sentential Form | Input |
|------|-----------------|-------|
| — | $\langle id, x \rangle$ - *Term* | <u>x</u> - ↑<u>2</u> * <u>y</u> |
| 7 | $\langle id, x \rangle$ - *Factor* | <u>x</u> - ↑<u>2</u> * <u>y</u> |
| 9 | $\langle id, x \rangle$ - $\langle num, 2 \rangle$ | <u>x</u> - ↑<u>2</u> * <u>y</u> |
| — | $\langle id, x \rangle$ - $\langle num, 2 \rangle$ | <u>x</u> - <u>2</u> ↑* <u>y</u> |

Where are we?

- "2" matches "2"

- We have more input, but no *NT*s left to expand

- The expansion terminated too soon

⇒ Need to backtrack
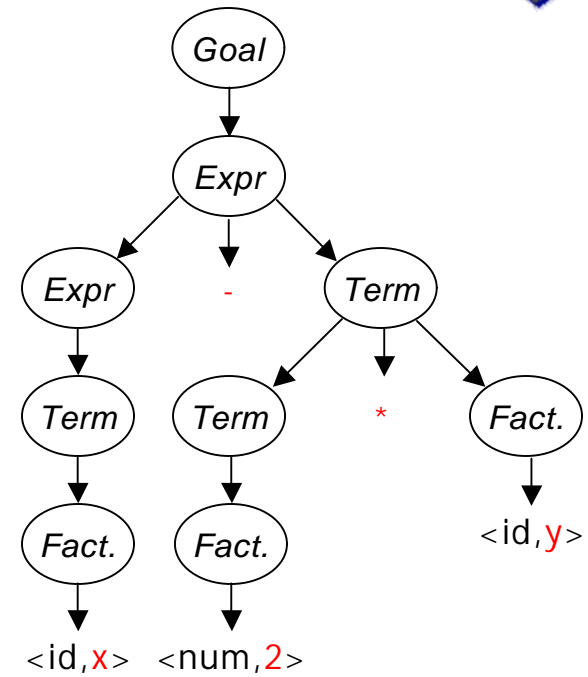
# Example

Trying again with "2" in <u>x</u> - <u>2</u> * <u>y</u> :

| Rule | Sentential Form | Input |
|------|-----------------|-------|
| — | <id, x> - *Term* | <u>x</u> - ↑<u>2</u> * y |
| 5 | <id, x> - *Term * Factor* | <u>x</u> - ↑<u>2</u> * y |
| 7 | <id, x> - *Factor * Factor* | <u>x</u> - ↑<u>2</u> * y |
| 8 | <id, x> - <num,2> * *Factor* | <u>x</u> - ↑<u>2</u> * y |
| — | <id, x> - <num,2> * *Factor* | <u>x</u> - <u>2</u> ↑* y |
| — | <id, x> - <num,2> * *Factor* | <u>x</u> - <u>2</u> * ↑y |
| 9 | <id, x> - <num,2> * <id, y> | <u>x</u> - <u>2</u> * ↑y |
| — | <id, x> - <num,2> * <id, y> | <u>x</u> - <u>2</u> * <u>y</u>↑ |

This time, we matched & consumed all the input

$\Rightarrow$ Success!

## *Another possible parse*

Other choices for expansion are possible

| Rule | Sentential Form | Input |
|------|-----------------|-------|
| — | Goal | ↑x - 2 * y |
| 1 | Expr | ↑x - 2 * y |
| 2 | Expr + Term | ↑x - 2 * y |
| 2 | Expr + Term + Term | ↑x - 2 * y |
| 2 | Expr + Term + Term + Term | ↑x - 2 * y |
| 2 | Expr + Term + Term + … + Term | ↑x - 2 * y |

consuming no input !

This doesn't terminate                    *(obviously)*

• 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

from Cooper & Torczon

## *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 $A \Rightarrow^+ A\boldsymbol{a}$, for some string $\boldsymbol{a} \in (NT \cup T)^+$

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*

from Cooper & Torczon

## Eliminating Left Recursion

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*

from Cooper & Torczon

## Eliminating Left Recursion

The expression grammar contains two cases of left recursion

$$Expr \rightarrow Expr + Term$$
$$| \quad Expr - Term$$
$$| \quad Term$$

$$Term \rightarrow Term * Factor$$
$$| \quad Term \,/\, Factor$$
$$| \quad Factor$$

Applying the transformation yields

$$Expr \rightarrow Term \; Expr'$$
$$Expr' \quad | \quad + Term \; Expr'$$
$$| \quad - Term \; Expr'$$
$$| \quad \varepsilon$$

$$Term \rightarrow Factor \; Term'$$
$$Term' \quad | \quad * Factor \; Term'$$
$$| \quad / Factor \; Term'$$
$$| \quad \varepsilon$$

These fragments use only right recursion

They retains the original left associativity

## *Eliminating Left Recursion*

Substituting back into the grammar yields

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

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

- It is left associative, as was the original

- A top-down parser will terminate using it.

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

## Eliminating Left Recursion

The transformation eliminates immediate left recursion

What about more general, indirect left recursion

The general algorithm:

*arrange the NTs into some order $A_1$, $A_2$, ..., $A_n$*

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

*replace each production $A_i \rightarrow A_s \, g$ with*
*$A_i \rightarrow d_1 \, g \mid d_2 \, g \mid \ldots \mid d_k \, g$, where $A_s \rightarrow d_1 \mid d_2 \mid \ldots \mid d_k$*
*are all the current productions for $A_s$*

*eliminate any immediate left recursion on $A_i$*
*using the direct transformation*

This assumes that the initial grammar has no cycles ($A_i \Rightarrow^+ A_i$),

and no epsilon productions

How does this algorithm work?

1. Impose arbitrary order on the non-terminals

2. Outer loop cycles through NT in order

3. Inner loop ensures that a production expanding $A_i$ has no non-terminal $A_s$ in its *rhs,* for $s < I$

4. Last step in outer loop converts any direct recursion on $A_i$ to right recursion using the transformation showed earlier

5. New non-terminals are added at the end of the order & have no left recursion


At the start of the $i^{th}$ outer loop iteration

> *For all k < I, no production that expands $A_k$ contains a non-terminal $A_s$ in its rhs, for s < k*