

# Compiler construction

Martin Steffen

February 20, 2017

## Contents

<b>1 Abstract</b>	<b>1</b>
1.1 Parsing . . . . .	1
1.1.1 First and follow sets . . . . .	1
1.1.2 Top-down parsing . . . . .	15
1.1.3 Bottom-up parsing . . . . .	43
<b>2 Reference</b>	<b>78</b>

## 1 Abstract

### Abstract

This is the *handout* version of the slides. It contains basically the same content, only in a way which allows more compact printing. Sometimes, the overlays, which make sense in a presentation, are not fully rendered here.

Besides the material of the slides, the handout versions may also contain additional remarks and background information which may or may not be helpful in getting the bigger picture.

### 1.1 Parsing

31. 1. 2017

#### 1.1.1 First and follow sets

##### Overview

- First and Follow set: general concepts for grammars
  - textbook looks at one parsing technique (top-down) [Louden, 1997, Chap. 4] before studying First/Follow sets
  - we: take First/Follow sets *before* any parsing technique
- two *transformation* techniques for grammars, both *preserving* the accepted language
  1. removal for left-recursion
  2. left factoring

##### First and Follow sets

- general concept for grammars
  - certain types of analyses (e.g. parsing):
    - info needed about possible “forms” of *derivable* words,
1. First-set of  $A$ 
    - which terminal symbols can appear at the start of strings *derived from* a given nonterminal  $A$

2. Follow-set of  $A$  Which terminals can follow  $A$  in some *sentential form*.

### 3. Remarks

- sentential form: word *derived from* grammar's starting symbol
- later: different algos for First and Follow sets, for all non-terminals of a given grammar
- mostly straightforward
- one complication: *nullable* symbols (non-terminals)
- Note: those sets depend on grammar, not the language

## First sets

**Definition 1** (First set). Given a grammar  $G$  and a non-terminal  $A$ . The *First-set* of  $A$ , written  $First_G(A)$  is defined as

$$First_G(A) = \{a \mid A \Rightarrow_G^* a\alpha, \quad a \in \Sigma_T\} + \{\epsilon \mid A \Rightarrow_G^* \epsilon\}. \quad (1)$$

**Definition 2** (Nullable). Given a grammar  $G$ . A non-terminal  $A \in \Sigma_N$  is *nullable*, if  $A \Rightarrow^* \epsilon$ .

### 1. Nullable

The definition here of being nullable refers to a non-terminal symbol. When concentrating on context-free grammars, as we do for parsing, that's basically the only interesting case. In principle, one can define the notion of being nullable analogously for arbitrary words from the whole alphabet  $\Sigma = \Sigma_T + \Sigma_N$ . The form of productions in CFGs makes it obvious, that the only words which actually may be nullable are words containing only non-terminals. Once a terminal is derived, it can never be "erased". It's equally easy to see, that a word  $\alpha \in \Sigma_N^*$  is nullable iff all its non-terminal symbols are nullable. The same remarks apply to context-sensitive (but not general) grammars.

For level-0 grammars in the Chomsky-hierarchy, also words containing terminal symbols may be nullable, and nullability of a word, like most other properties in that setting, becomes undecidable.

### 2. First and follow set

One point worth noting is that the first and the follow sets, while seemingly quite similar, *differ* in one aspect. The First set is about words derivable from a given non-terminal  $A$ . The follow set is about words derivable from the starting symbol! As a consequence, non-terminals  $A$  which are not *reachable* from the grammar's starting symbol have, by definition, an empty follow set. In contrast, non-terminals unreachable from the start symbol may well have a non-empty first-set. In practice a grammar containing unreachable non-terminals is ill-designed, so that distinguishing feature in the definition of the first and the follow set for a non-terminal may not matter so much. Nonetheless, when *implementing* the algo's for those sets, those subtle points do matter! In general, to avoid all those fine points, one works with grammars satisfying a number of common-sense restrictions. One are so called *reduced grammars*, where, informally, all symbols "play a role" (all are reachable, all can derive into a word of terminals).

## Examples

- Cf. the Tiny grammar
- in Tiny, as in most languages

$$First(if-stmt) = \{"if"\}$$

- in many languages:

$$First(assign-stmt) = \{\text{identifier}, "("\}$$

- typical *Follow* (see later) for statements:

$$Follow(stmt) = \{";", "\text{end}", "\text{else}", "\text{until"}\}$$

## Remarks

- note: special treatment of the empty word  $\epsilon$
  - in the following: if grammar  $G$  clear from the context
    - $\Rightarrow^*$  for  $\Rightarrow_G^*$
    - $First$  for  $First_G$
    - ...
  - definition so far: “top-level” for start-symbol, only
  - next: a more general definition
    - definition of First set of arbitrary symbols (and even words)
    - and also: definition of First for a symbol *in terms of* First for “other symbols” (connected by *productions*)
- $\Rightarrow$  recursive definition

## A more algorithmic/recursive definition

- grammar symbol  $X$ : terminal or non-terminal or  $\epsilon$

**Definition 3** (First set of a symbol). Given a grammar  $G$  and grammar symbol  $X$ . The *First-set* of  $X$ , written  $First(X)$ , is defined as follows:

1. If  $X \in \Sigma_T + \{\epsilon\}$ , then  $First(X) = \{X\}$ .
2. If  $X \in \Sigma_N$ : For each production

$$X \rightarrow X_1 X_2 \dots X_n$$

- (a)  $First(X)$  contains  $First(X_1) \setminus \{\epsilon\}$
- (b) If, for some  $i < n$ , all  $First(X_1), \dots, First(X_i)$  contain  $\epsilon$ , then  $First(X)$  contains  $First(X_{i+1}) \setminus \{\epsilon\}$ .
- (c) If all  $First(X_1), \dots, First(X_n)$  contain  $\epsilon$ , then  $First(X)$  contains  $\{\epsilon\}$ .

1. Recursive definition of *First*?

The following discussion may be ignored if wished. Even if details and theory behind it is beyond the scope of this lecture, it is worth considering above definition more closely. One may even consider if it is a definition at all (resp. in which way it is a definition).

One naive first impression may be: it's a kind of a “functional definition”, i.e., the above Definition 3 gives a recursive definition of the function *First*. As discussed later, everything gets rather simpler if we would not have to deal with nullable words and  $\epsilon$ -productions. For the point being explained here, let's assume that there are no such productions and get rid of the special cases, cluttering up Definition 3. Removing the clutter gives the following simplified definition:

**Definition 4** (First set of a symbol (no  $\epsilon$ -productions)). Given a grammar  $G$  and grammar symbol  $X$ . The *First-set* of  $X \neq \epsilon$ , written  $First(X)$  is defined as follows:

- (a) If  $X \in \Sigma_T$ , then  $First(X) \supseteq \{X\}$ .
- (b) If  $X \in \Sigma_N$ : For each production

$$X \rightarrow X_1 X_2 \dots X_n ,$$

$$First(X) \supseteq First(X_1).$$

Compared to the previous condition, I did the following 2 minor adaptations (apart from cleaning up the  $\epsilon$ 's): In case (1b), I replaced the English word “contains” with the superset relation symbol  $\supseteq$ . In case (1a), I replaced the equality symbol  $=$  with the superset symbol  $\supseteq$ , basically for consistency with the other case.

Now, with Definition 4 as a simplified version of the original definition being made slightly more explicit and internally consistent: in which way is that a definition at all?

For being a definition for  $\text{First}(X)$ , it seems awfully lax. Already in (1a), it “defines” that  $\text{First}(X)$  should “at least contain  $X$ ”. A similar remark applies to case (1b) for non-terminals. Those two requirements are as such well-defined, but *they don't define  $\text{First}(X)$  in a unique manner!* Definition 4 defines what the set  $\text{First}(X)$  should *at least* contain!

So, in a nutshell, one should not consider Definition 4 a “recursive definition of  $\text{First}(X)$ ” but rather

“a definition of recursive conditions on  $\text{First}(X)$ , which, when satisfied, ensures that  $\text{First}(X)$  contains *at least* all non-terminals we are after”.

What we are *really* after is the *smallest*  $\text{First}(X)$  which satisfies those conditions of the definitions.

Now one may think: the problem is that definition is just “sloppy”. Why does it use the word “contain” resp. the  $\supseteq$ -relation, instead of requiring equality, i.e.,  $=$ ? While plausible at first sight, unfortunately, whether we use  $\supseteq$  or set equality  $=$  in Definition 4 does not change anything (and remember that the original Definition 3 “mixed up” the styles by requiring equality in the case of non-terminals and requiring “contains”, i.e.,  $\supseteq$  for non-terminals).

Anyhow, the core of the matter is not  $=$  vs.  $\supseteq$ . The core of the matter is that “Definition” 4 is *circular*!

Considering that definition of  $\text{First}(X)$  as a plain functional and recursive definition of a procedure missed the fact that grammar can, of course, contain “loops”. Actually, it's almost a characterizing feature of reasonable context-free grammars (or even regular grammars) that they contain “loops” – that's the way they can describe infinite languages.

In that case, obviously, considering Definition 3 with  $=$  instead of  $\supseteq$  as the recursive definition of a function leads immediately to an “infinite regress”, the recursive function won't terminate. So again, that's not helpful.

Technically, such a definition is called to be a *constraint* (or a constraint system, if one considers the whole definition to consist of more than one constraint, namely for different terminals and for different productions).

## For words

**Definition 5** (First set of a word). Given a grammar  $G$  and word  $\alpha$ . The *First-set* of

$$\alpha = X_1 \dots X_n ,$$

written  $\text{First}(\alpha)$  is defined inductively as follows:

1.  $\text{First}(\alpha)$  contains  $\text{First}(X_1) \setminus \{\epsilon\}$
2. for each  $i = 2, \dots, n$ , if  $\text{First}(X_k)$  contains  $\epsilon$  for all  $k = 1, \dots, i-1$ , then  $\text{First}(\alpha)$  contains  $\text{First}(X_i) \setminus \{\epsilon\}$
3. If all  $\text{First}(X_1), \dots, \text{First}(X_n)$  contain  $\epsilon$ , then  $\text{First}(\alpha)$  contains  $\{\epsilon\}$ .

1. Concerning the definition of First

The definition here is of course very close to the definition of inductive case of the previous definition, i.e., the first set of a non-terminal. Whereas the previous definition was a recursive, this one is not.

Note that the word  $\alpha$  may be empty, i.e.,  $n = 0$ . In that case, the definition gives  $\text{First}(\epsilon) = \{\epsilon\}$  (due to the 3rd condition in the above definition). In the definitions, the empty word  $\epsilon$  plays a specific, mostly technical role. The original, non-algorithmic version of Definition 1, makes it already clear, that the first set *not precisely* corresponds to the set of terminal symbols that can appear at the beginning of a derivable word. The correct intuition is that it corresponds to that set of terminal symbols *together* with  $\epsilon$  as a special case, namely when the initial symbol is nullable.

That may raise two questions. 1) Why does the definition makes that as special case, as opposed to just using the more “straightforward” definition without taking care of the nullable situation? 2) What role does  $\epsilon$  play here?

The second question has no “real” answer, it’s a choice which is being made which could be made differently. What the definition from equation (1) in fact says is: “give the set of terminal symbols in the derivable word **and** indicate whether or not the start symbol is *nullable*. The information might as well be interpreted as a *pair* consisting of a set of terminals *and* a boolean (indicating nullability). The fact that the definition of *First* as presented here uses  $\epsilon$  to indicate that additional information is a particular choice of representation (probably due to historical reasons: “they always did it like that …”). For instance, the influential “Dragon book” [Aho et al., 2007, Section 4.4.2] uses the  $\epsilon$ -based definition. The textbooks [Appel, 1998a] (and its variants) don’t use  $\epsilon$  as indication for nullability.

In order that this definition works, it is important, obviously, that  $\epsilon$  is *not* a terminal symbol, i.e.,  $\epsilon \notin \Sigma_T$  (which is generally assumed).

Having clarified 2), namely that using  $\epsilon$  is a matter of conventional choice, remains question 1), why bother to include nullability-information in the definition of the first-set *at all*, why bother with the “extra information” of nullability? For that, there is a real technical reason: For the *recursive* definitions to work, we need the information whether or not a symbol or word is *nullable*, therefore it’s given back as information.

As a further point concerning the first sets: The slides give 2 definitions, Definition 1 and Definition 3. Of course they are intended to mean the same. The second version is a more recursive or algorithmic version, i.e., closer to a recursive algorithm. If one takes the first one as the “real” definition of that set, in principle we would be obliged to prove that both versions actually describe the same same (resp. that the recursive definition *implements* the original definition). The same remark applies also to the non-recursive/iterative code that is shown next.

### Pseudo code

```

for all non-terminals A do
    First[A] := {}
end
while there are changes to any First[A] do
    for each production A → X1...Xn do
        k := 1;
        continue := true
        while continue = true and k ≤ n do
            First[A] := First[A] ∪ First[Xk] \ {ε}
            if ε ∉ First[Xk] then continue := false
            k := k + 1
        end;
        if continue = true
            First[A] := First[A] ∪ {ε}
        end;
    end

```

If only we could do away with special cases for the empty words ... for grammar without  $\epsilon$ -productions.<sup>1</sup>

```

for all non-terminals A do
    First[A] := {} // counts as change
end
while there are changes to any First[A] do
    for each production A → X1...Xn do
        First[A] := First[A] ∪ First[X1]
    end;
end

```

### Example expression grammar (from before)

$$\begin{array}{lcl}
 \text{exp} & \rightarrow & \text{exp addop term} \mid \text{term} \\
 \text{addop} & \rightarrow & + \mid - \\
 \text{term} & \rightarrow & \text{term mulop factor} \mid \text{factor} \\
 \text{mulop} & \rightarrow & * \\
 \text{factor} & \rightarrow & (\text{exp}) \mid \text{number}
 \end{array} \tag{2}$$

<sup>1</sup>A production of the form  $A \rightarrow \epsilon$ .

### Example expression grammar (expanded)

$exp \rightarrow exp\ addop\ term$	(3)
$exp \rightarrow term$	
$addop \rightarrow +$	
$addop \rightarrow -$	
$term \rightarrow term\ mulop\ factor$	
$term \rightarrow factor$	
$mulop \rightarrow *$	
$factor \rightarrow ( exp )$	
$factor \rightarrow n$	

### “Run” of the algo

nr	pass 1	pass 2	pass 3
1	$exp \rightarrow exp\ addop\ term$		
2	$exp \rightarrow term$		
3	$addop \rightarrow +$		
4	$addop \rightarrow -$		
5	$term \rightarrow term\ mulop\ factor$		
6	$term \rightarrow factor$		
7	$mulop \rightarrow *$		
8	$factor \rightarrow ( exp )$		
9	$factor \rightarrow n$		

#### 1. How the algo works

The first thing to observe: the grammar does not contain  $\epsilon$ -productions. That, very fortunately, simplifies matters considerably! It should also be noted that the table from above is a schematic illustration of a particular *execution strategy* of the pseudo-code. The pseudo-code itself leaves out details of the evaluation, notably *the order* in which non-deterministic choices are done by the code. The main body of the pseudo-code is given by two nested loops. Even if details (of data structures) are not given, one possible way of interpreting the code is as follows: the outer while-loop figures out which of the entries in the `First`-array have “recently” been changed, remembers that in a “collection” of non-terminals  $A$ ’s, and that collection is then worked off (i.e. iterated over) on the inner loop. Doing it like that leads to the “passes” shown in the table. In other words, the two dimensions of the table represent the fact that there are 2 nested loops.

Having said that: it’s not the only way to “traverse the productions of the grammar”. One could arrange a version with only one loop and a collection data structure, which contains all productions  $A \rightarrow X_1 \dots X_n$  such that `First[A]` has “recently been changed”. That data structure therefore contains all the productions that “still need to be treated”. Such a collection data structure containing “all the work still to be done” is known as *work-list*, even if it needs not technically be a list. It can be a queue, i.e., following a FIFO strategy, it can be a stack (realizing LIFO), or some other strategy or heuristic. Possible is also a randomized, i.e., non-deterministic strategy (which is sometimes known as chaotic iteration).

#### 2. Recursion vs. iteration

### “Run” of the algo

Grammar rule	Pass 1	Pass 2	Pass 3
$exp \rightarrow exp$ $\quad addop \ term$			
$exp \rightarrow term$			$\text{First}(exp) = \{(, \text{number}\}$
$addop \rightarrow +$	$\text{First}(addop) = \{+\}$		
$addop \rightarrow -$	$\text{First}(addop) = \{+, -\}$		
$term \rightarrow term$ $\quad mulop \ factor$			
$term \rightarrow factor$		$\text{First}(term) = \{(, \text{number}\}$	
$mulop \rightarrow *$	$\text{First}(mulop) = \{*\}$		
$factor \rightarrow ( \ exp \ )$	$\text{First}(factor) = \{( \ )\}$		
$factor \rightarrow \text{number}$	$\text{First}(factor) = \{(, \text{number}\}$		

### Collapsing the rows & final result

- results per pass:

	1	2	3
$exp$			$\{(, \text{n}\}$
$addop$	$\{+, -\}$		
$term$		$\{(, \text{n}\}$	
$mulop$	$\{*\}$		
$factor$	$\{(, \text{n}\}$		

- final results (at the end of pass 3):

	$First[ ]$
$exp$	$\{(, \text{n}\}$
$addop$	$\{+, -\}$
$term$	$\{(, \text{n}\}$
$mulop$	$\{*\}$
$factor$	$\{(, \text{n}\}$

### Work-list formulation

```

for all non-terminals A do
    First[A] := {}
    WL      := P    // all productions
end
while WL ≠ ∅ do
    remove one (A → X1...Xn) from WL
    if   First[A] ≠ First[A] ∪ First[X1]
    then First[A] := First[A] ∪ First[X1]
        add all productions (A → X'1'...X'm) to WL
    else skip
end

```

- worklist here: “collection” of productions
- alternatively, with slight reformulation: “collection” of non-terminals instead also possible

## Follow sets

**Definition 6** (Follow set (ignoring  $\$$ )). Given a grammar  $G$  with start symbol  $S$ , and a non-terminal  $A$ . The *Follow-set* of  $A$ , written  $Follow_G(A)$ , is

$$Follow_G(A) = \{a \mid S \Rightarrow_G^* \alpha_1 A a \alpha_2, \quad a \in \Sigma_T\} . \quad (4)$$

- More generally:  $\$$  as special end-marker

$$S \$ \Rightarrow_G^* \alpha_1 A a \alpha_2, \quad a \in \Sigma_T + \{ \$ \} .$$

- typically: start symbol *not* on the right-hand side of a production

### 1. Special symbol $\$$

The symbol  $\$$  can be interpreted as “end-of-file” (EOF) token. It’s standard to assume that the start symbol  $S$  does not occur on the right-hand side of any production. In that case, the follow set of  $S$  contains  $\$$  as *only* element. Note that the follow set of other non-terminals may well contain  $\$$ .

As said, it’s common to assume that  $S$  does not appear on the right-hand side on any production. For a start,  $S$  won’t occur “naturally” there anyhow in practical programming language grammars. Furthermore, with  $S$  occurring only on the left-hand side, the grammar has a slightly nicer shape insofar as it makes its algorithmic treatment slightly nicer. It’s basically the same reason why one sometimes assumes that for instance, control-flow graphs has one “isolated” entry node (and/or an isolated exit node), where being isolated means, that no edge in the graph goes (back) into the entry node; for exits nodes, the condition means, no edge goes out. In other words, while the graph can of course contain loops or cycles, the entry node is not part of any such loop. That is done likewise to (slightly) simplify the treatment of such graphs. Slightly more generally and also connected to control-flow graphs: similar conditions about the shape of loops (not just for the entry and exit nodes) have been worked out, which play a role in loop optimization and intermediate representations of a compiler, such as static single assignment forms.

Coming back to the condition here concerning  $\$$ : even if a grammar would not immediately adhere to that condition, it’s trivial to transform it into that form by adding another symbol and make that the new start symbol, replacing the old.

### 2. Special symbol $\$$

It seems that [Appel, 1998b] does not use the special symbol in his treatment of the follow set, but the dragon book uses it. It is used to represent the symbol (not otherwise used) “right of the start symbol”, resp., the symbol right of a non-terminal which is at the right end of a derived word.

## Follow sets, recursively

**Definition 7** (Follow set of a non-terminal). Given a grammar  $G$  and nonterminal  $A$ . The *Follow-set* of  $A$ , written  $Follow(A)$  is defined as follows:

1. If  $A$  is the start symbol, then  $Follow(A)$  contains  $\$$ .
  2. If there is a production  $B \rightarrow \alpha A \beta$ , then  $Follow(A)$  contains  $First(\beta) \setminus \{\epsilon\}$ .
  3. If there is a production  $B \rightarrow \alpha A \beta$  such that  $\epsilon \in First(\beta)$ , then  $Follow(A)$  contains  $Follow(B)$ .
- $\$$ : “end marker” special symbol, only to be contained in the follow set

## More imperative representation in pseudo code

```

Follow [  $S$  ] := { $\$$ }
for all non-terminals  $A \neq S$  do
    Follow [  $A$  ] := {}
end
while there are changes to any Follow-set do
    for each production  $A \rightarrow X_1 \dots X_n$  do
        for each  $X_i$  which is a non-terminal do
            Follow [  $X_i$  ] := Follow [  $X_i$  ]  $\cup$  First(  $X_{i+1} \dots X_n$  )  $\setminus \{\epsilon\}$ 
            if  $\epsilon \in$  First(  $X_{i+1} X_{i+2} \dots X_n$  )
            then Follow [  $X_i$  ] := Follow [  $X_i$  ]  $\cup$  Follow [  $A$  ]
        end
    end
end

```

Note!  $\text{First}() = \{\epsilon\}$

## Expression grammar once more

### “Run” of the algo

nr	pass 1	pass 2
1	$exp \rightarrow exp \text{ addop } term$	
2	$exp \rightarrow term$	
5	$term \rightarrow term \text{ mulop } factor$	
6	$term \rightarrow factor$	
8	$factor \rightarrow ( exp )$	

normalsize

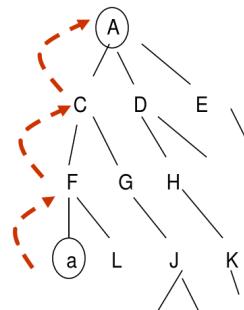
1. Recursion vs. iteration

### “Run” of the algo

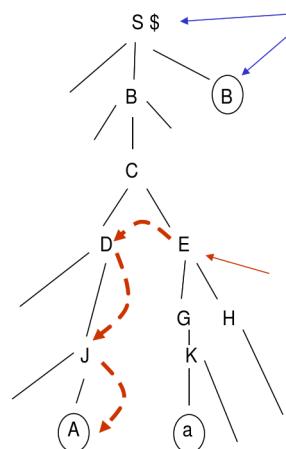
Grammar rule	Pass 1	Pass 2
$exp \rightarrow exp \text{ addop}$ $term$	$\text{Follow}(exp) = \{\$, +, -\}$ $\text{Follow}(addop) = \{(\text{, number})\}$ $\text{Follow}(term) = \{\$, +, -\}$	$\text{Follow}(term) = \{\$, +, -, *, )\}$
$exp \rightarrow term$		
$term \rightarrow term \text{ mulop}$ $factor$	$\text{Follow}(term) = \{\$, +, -, *\}$ $\text{Follow}(mulop) = \{(\text{, number})\}$ $\text{Follow}(factor) = \{\$, +, -, *\}$	$\text{Follow}(factor) = \{\$, +, -, *, )\}$
$term \rightarrow factor$		
$factor \rightarrow (\text{ exp })$	$\text{Follow}(exp) = \{\$, +, -, )\}$	

### Illustration of first/follow sets

$a \in First(A)$



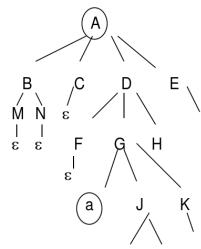
$a \in Follow(A)$



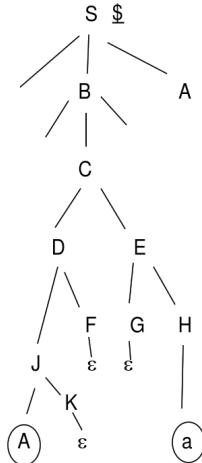
- red arrows: illustration of information flow in the algs
- run of *Follow*:
  - relies on *First*
  - in particular  $a \in First(E)$  (right tree)
- $\$ \in Follow(B)$

### More complex situation (nullability)

$a \in First(A)$



$a \in Follow(A)$



### 1. Remarks

In the tree on the left,  $B, M, N, C$ , and  $F$  are *nullable*. That is marked in that the resulting first sets contain  $\epsilon$ . There will also be exercises about that.

### Some forms of grammars are less desirable than others

- **left-recursive** production:

$$A \rightarrow A\alpha$$

more precisely: example of *immediate* left-recursion

- 2 productions with **common “left factor”**:

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \quad \text{where } \alpha \neq \epsilon$$

#### 1. Left-recursive and unfactored grammars

At the current point in the presentation, the importance of those conditions might not yet be clear. In general, it's that certain kind of parsing techniques require absence of left-recursion and of common left-factors. Note also that a left-linear production is a special case of a production with immediate left recursion. In particular, recursive descent parsers would not work with left-recursion, as they would require “unbounded look-ahead”. For that kind of parsers, left-recursion needs to be avoided.

Why common left-factors are undesirable should at least intuitively be clear: we see this also on the next slide (the two forms of conditionals). It's intuitively clear, that a parser, when encountering an **if** (and the following boolean condition and perhaps the **then** clause) cannot decide immediately which rule applies. It should also be intuitively clear that that's what a parser does: inputting a stream of tokens and trying to figure out which sequence of rules are responsible for that stream (or else reject the input). The amount of additional information, at each point of the parsing process, to *determine* which rule is responsible next is called the *look-ahead*. Of course, if the grammar is ambiguous, no unique decision may be possible (no matter the look-ahead). Ambiguous grammars are unwelcome as specification for parsers.

On a very high-level, the situation can be compared with the situation for regular languages/automata. Non-deterministic automata may be ok for specifying the *language* (they can more easily be connected to regular expressions), but they are not so useful for specifying a scanner *program*. There, deterministic automata are necessary. Here, grammars with left-recursion, grammars with common factors, or even ambiguous grammars may be ok for specifying a context-free language. For instance, ambiguity may be caused by unspecified precedences or non-associativity. Nonetheless, how to obtain a grammar representation more suitable to be more or less directly translated to a parser is an issue less clear cut compared to regular languages. Already the question whether or not a given grammar is ambiguous or not is undecidable. If ambiguous, there'd be no point in turning it into a practical parser. Also the question, what's an acceptable form of grammar depends on what class of parsers one is after (like a top-down parser or a bottom-up parser).

## Some simple examples for both

- left-recursion

$$exp \rightarrow exp + term$$

- classical example for common left factor: rules for conditionals

$$\begin{array}{lcl} if-stmt & \rightarrow & \text{if } ( exp ) stmt \text{ end} \\ & | & \text{if } ( exp ) stmt \text{ else } stmt \text{ end} \end{array}$$

## Transforming the expression grammar

$$\begin{array}{lcl} exp & \rightarrow & exp \text{ addop } term \mid term \\ addop & \rightarrow & + \mid - \\ term & \rightarrow & term \text{ mulop } factor \mid factor \\ mulop & \rightarrow & * \\ factor & \rightarrow & ( exp ) \mid \text{number} \end{array}$$

- obviously left-recursive
- remember: this variant used for proper **associativity**!

## After removing left recursion

$$\begin{array}{lcl} exp & \rightarrow & term \text{ exp}' \\ exp' & \rightarrow & addop \text{ term } exp' \mid \epsilon \\ addop & \rightarrow & + \mid - \\ term & \rightarrow & factor \text{ term}' \\ term' & \rightarrow & mulop \text{ factor } term' \mid \epsilon \\ mulop & \rightarrow & * \\ factor & \rightarrow & ( exp ) \mid \text{n} \end{array}$$

- still *unambiguous*
- unfortunate: *associativity* now different!
- note also:  $\epsilon$ -productions & nullability

## Left-recursion removal

1. Left-recursion removal A transformation process to turn a CFG into one without left recursion
2. Explanation

- price:  $\epsilon$ -productions
- 3 cases to consider
  - immediate (or direct) recursion
    - \* simple
    - \* general
  - *indirect* (or mutual) recursion

## Left-recursion removal: simplest case

1. Before

$$A \rightarrow A\alpha \mid \beta$$

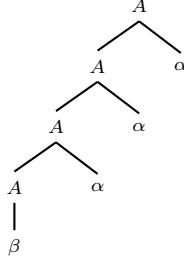
2. space

3. After

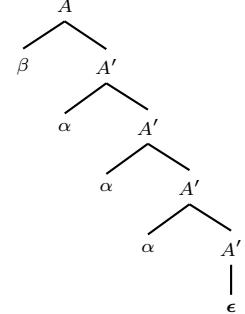
$$\begin{array}{lcl} A & \rightarrow & \beta A' \\ A' & \rightarrow & \alpha A' \mid \epsilon \end{array}$$

### Schematic representation

$$A \rightarrow A\alpha \mid \beta$$



$$\begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \epsilon \end{array}$$



### Remarks

- both grammars generate the same (context-free) language (= set of words over terminals)
- in EBNF:

$$A \rightarrow \beta\{\alpha\}$$

- two *negative* aspects of the transformation
  1. generated language unchanged, but: change in resulting structure (parse-tree), i.a.w. change in **associativity**, which may result in change of *meaning*
  2. introduction of  $\epsilon$ -productions
- more concrete example for such a production: grammar for expressions

### Left-recursion removal: immediate recursion (multiple)

1. Before

$$\begin{array}{l} A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_n \\ \mid \beta_1 \mid \dots \mid \beta_m \end{array}$$

2. space

3. After

$$\begin{array}{l} A \rightarrow \beta_1 A' \mid \dots \mid \beta_m A' \\ A' \rightarrow \alpha_1 A' \mid \dots \mid \alpha_n A' \\ \mid \epsilon \end{array}$$

4. EBNF

Note: can be written in EBNF as:

$$A \rightarrow (\beta_1 \mid \dots \mid \beta_m)(\alpha_1 \mid \dots \mid \alpha_n)^*$$

**Removal of: general left recursion** Assume non-terminals  $A_1, \dots, A_m$

```

for i := 1 to m do
  for j := 1 to i-1 do
    replace each grammar rule of the form  $A_i \rightarrow A_j \beta$  by //  $i < j$ 
    rule  $A_i \rightarrow \alpha_1 \beta \mid \alpha_2 \beta \mid \dots \mid \alpha_k \beta$ 
    where  $A_j \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_k$ 
    is the current rule(s) for  $A_j$  // current, i.e., at the current stage of algo
  end
  { corresponds to  $i = j$  }
  remove, if necessary, immediate left recursion for  $A_i$ 
end
  
```

“current” = rule in the current stage of algo

**Example (for the general case)** let  $A = A_1, B = A_2$

$$\begin{array}{lcl} A & \rightarrow & Ba \mid Aa \mid c \\ B & \rightarrow & Bb \mid Ab \mid d \end{array}$$

$$\begin{array}{lcl} A & \rightarrow & BaA' \mid cA' \\ A' & \rightarrow & aA' \mid \epsilon \\ B & \rightarrow & Bb \mid Ab \mid d \end{array}$$

$$\begin{array}{lcl} A & \rightarrow & BaA' \mid cA' \\ A' & \rightarrow & aA' \mid \epsilon \\ B & \rightarrow & Bb \mid BaA'b \mid cA'b \mid d \end{array}$$

$$\begin{array}{lcl} A & \rightarrow & BaA' \mid cA' \\ A' & \rightarrow & aA' \mid \epsilon \\ B & \rightarrow & cA'bB' \mid dB' \\ B' & \rightarrow & bB' \mid aA'bB' \mid \epsilon \end{array}$$

### Left factor removal

- CFG: not just describe a context-free languages
- also: intended (indirect) description of a **parser** for that language
- ⇒ common left factor undesirable
- cf.: *determinization* of automata for the lexer

#### 1. Simple situation

(a) before

$$A \rightarrow \alpha\beta \mid \alpha\gamma \mid \dots$$

(b) after

$$\begin{array}{lcl} A & \rightarrow & \alpha A' \mid \dots \\ A' & \rightarrow & \beta \mid \gamma \end{array}$$

### Example: sequence of statements

#### 1. sequences of statements

(a) Before

$$\begin{array}{lcl} stmt-seq & \rightarrow & stmt ; stmt-seq \\ & | & stmt \end{array}$$

(b) After

$$\begin{array}{lcl} stmt-seq & \rightarrow & stmt \; stmt-seq' \\ stmt-seq' & \rightarrow & ; \; stmt-seq \mid \epsilon \end{array}$$

### Example: conditionals

#### 1. Before

$$\begin{array}{lcl} if-stmt & \rightarrow & if( exp ) stmt-seq end \\ & | & if( exp ) stmt-seq else stmt-seq end \end{array}$$

#### 2. After

$$\begin{array}{lcl} if-stmt & \rightarrow & if( exp ) stmt-seq else-or-end \\ else-or-end & \rightarrow & else stmt-seq end \mid end \end{array}$$

### Example: conditionals (without else)

#### 1. Before

$$\begin{array}{lcl} if-stmt & \rightarrow & if( exp ) stmt-seq \\ & | & if( exp ) stmt-seq else stmt-seq \end{array}$$

#### 2. After

$$\begin{array}{lcl} if-stmt & \rightarrow & if( exp ) stmt-seq else-or-empty \\ else-or-empty & \rightarrow & else stmt-seq \mid \epsilon \end{array}$$

## Not all factorization doable in “one step”

1. Starting point

$$A \rightarrow abcB \mid abC \mid aE$$

2. After 1 step

$$\begin{array}{l} A \rightarrow abA' \mid aE \\ A' \rightarrow cB \mid C \end{array}$$

3. After 2 steps

$$\begin{array}{l} A \rightarrow aA'' \\ A'' \rightarrow bA' \mid E \\ A' \rightarrow cB \mid C \end{array}$$

4. longest left factor

- note: we choose the *longest* common prefix (= longest left factor) in the first step

## Left factorization

```

while there are changes to the grammar do
  for each nonterminal A do
    let  $\alpha$  be a prefix of max. length that is shared
        by two or more productions for A
    if  $\alpha \neq \epsilon$ 
    then
      let  $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$  be all
          prod. for A and suppose that  $\alpha_1, \dots, \alpha_k$  share  $\alpha$ 
          so that  $A \rightarrow \alpha\beta_1 \mid \dots \mid \alpha\beta_k \mid \alpha_{k+1} \mid \dots \mid \alpha_n$ ,
          that the  $\beta_j$ 's share no common prefix, and
          that the  $\alpha_{k+1}, \dots, \alpha_n$  do not share  $\alpha$ .
      replace rule  $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$  by the rules
       $A \rightarrow \alpha A' \mid \alpha_{k+1} \mid \dots \mid \alpha_n$ 
       $A' \rightarrow \beta_1 \mid \dots \mid \beta_k$ 
    end
  end
end

```

### 1.1.2 Top-down parsing

#### What's a parser generally doing

1. task of parser = syntax analysis

- input: stream of **tokens** from lexer
- output:
  - **abstract syntax tree**
  - or meaningful diagnosis of source of *syntax error*

- 2.

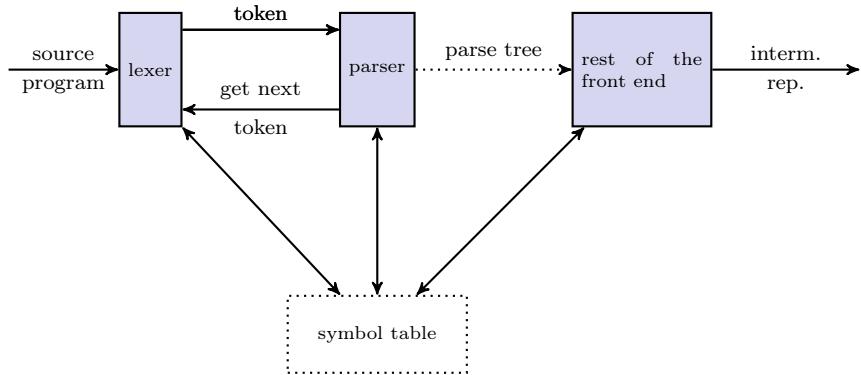
- the full “power” (i.e., expressiveness) of CFGs not used
- thus:
  - consider *restrictions* of CFGs, i.e., a specific subclass, and/or
  - *represented* in specific ways (no left-recursion, left-factored ...)

3. Syntax errors (and other errors) Since almost by definition, the *syntax* of a language are those aspects covered by a context-free grammar, a *syntax error* thereby is a violation of the grammar, something the parser has to detect. Given a CFG, typically given in BNF resp. implemented by a tool supporting a BNF variant, the parser (in combination with the lexer) must generate an AST *exactly* for those programs that adhere to the grammar and must *reject* all others. One says, the parser *recognizes* the given grammar. An important practical part when rejecting a program is to generate a meaningful error message, giving hints about potential location of the error and potential reasons. In the most minimal way, the parser should inform the programmer where

the parser tripped, i.e., telling how far, from left to right, it was able to proceed and informing where it stumbled: “parser error in line xxx/at character position yyy”). One typically has higher expectations for a real parser than just the line number, but that’s the basics.

It may be noted that also the subsequent phase, the *semantic analysis*, which takes the abstract syntax tree as input, may report errors, which are then no longer syntax errors but more complex kind of errors. One typical kind of error in the semantic phase is a *type error*. Also there, the minimal requirement is to indicate the probable location(s) where the error occurs. To do so, in basically all compilers, the nodes in an abstract syntax tree will contain information concerning the position in the original file, the resp. node corresponds to (like line-numbers, character positions). If the parser would not add that information into the AST, the semantic analysis would have no way to relate potential errors it finds to the original, concrete code in the input. Remember: the compiler goes in *phases*, and once the parsing phase is over, there’s no going back to scan the file *again*.

### Lexer, parser, and the rest



### Top-down vs. bottom-up

- all parsers (together with lexers): *left-to-right*
- remember: parsers operate with *trees*
  - parse tree (concrete syntax tree): representing grammatical derivation
  - abstract syntax tree: data structure
- 2 fundamental classes
- while parser eats through the token stream, it grows, i.e., builds up (at least conceptually) the parse tree:
  1. Bottom-up Parse tree is being grown from the leaves to the root.
  2. Top-down Parse tree is being grown from the root to the leaves.
  3. AST
  - while parse tree mostly conceptual: parsing build up the concrete data structure of AST bottom-up vs. top-down.

### Parsing restricted classes of CFGs

- parser: better be “efficient”
- full complexity of CFLs: not really needed in practice<sup>2</sup>
- classification of CF languages vs. CF grammars, e.g.:

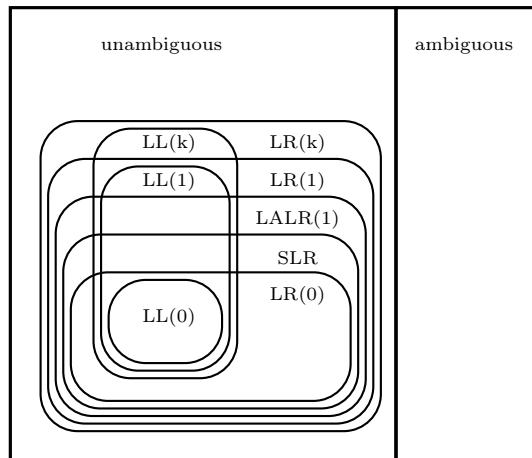
<sup>2</sup>Perhaps: if a parser has trouble to figure out if a program has a syntax error or not (perhaps using back-tracking), probably humans will have similar problems. So better keep it simple. And time in a compiler may be better spent elsewhere (optimization, semantical analysis).

- left-recursion-freedom: condition on a grammar
- ambiguous language vs. ambiguous grammar
- classification of grammars  $\Rightarrow$  classification of *language*
  - a CF language is (inherently) ambiguous, if there's no unambiguous grammar for it
  - a CF language is top-down parseable, if there exists a grammar that allows top-down parsing
  - ...
- in practice: classification of parser generating tools:
  - based on accepted notation for grammars: (BNF or some form of EBNF etc.)

### Classes of CFG grammars/languages

- *maaaaany* have been proposed & studied, including their relationships
- lecture concentrates on
  - top-down parsing, in particular
    - \* **LL(1)**
    - \* **recursive descent**
  - bottom-up parsing
    - \* **LR(1)**
    - \* **SLR**
    - \* **LALR(1)** (the class covered by yacc-style tools)
- grammars typically written in *pure* BNF

### Relationship of some grammar (not language) classes

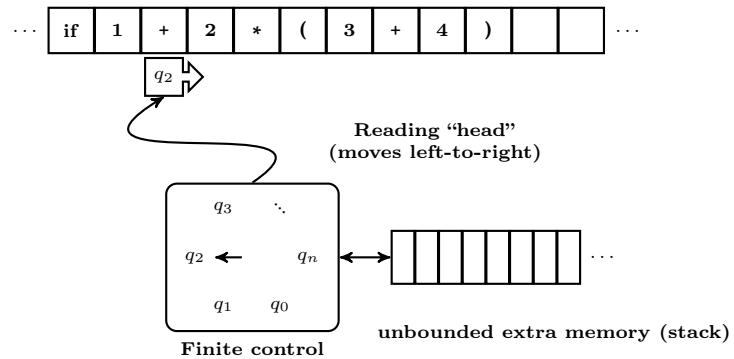


taken from [Appel, 1998c]

### General task (once more)

- Given: a CFG (but appropriately restricted)
- Goal: “systematic method” s.t.
  1. for every given word  $w$ : check syntactic correctness
  2. [build AST/representation of the parse tree as side effect]
  3. [do reasonable error handling]

## Schematic view on “parser machine”



Note: sequence of *tokens* (not characters)

## Derivation of an expression

*exp term exp' factor term' exp' number term' exp' number term' exp' number exp' number exp' number addop term exp' number+ term exp'*

- ## 1. factors and terms

$exp$	$\rightarrow$	$term \ exp'$	(5)
$exp'$	$\rightarrow$	$addop \ term \ exp' \mid \epsilon$	
$addop$	$\rightarrow$	$+ \mid -$	
$term$	$\rightarrow$	$factor \ term'$	
$term'$	$\rightarrow$	$mulop \ factor \ term' \mid \epsilon$	
$mulop$	$\rightarrow$	$*$	
$factor$	$\rightarrow$	$( \ exp ) \mid n$	

## Remarks concerning the derivation Note:

- input = stream of tokens
  - there: 1... stands for token class **number** (for readability/concreteness), in the grammar: just **number**
  - in full detail: pair of token class and token value  $\langle \text{number}, 1 \rangle$

## Notation:

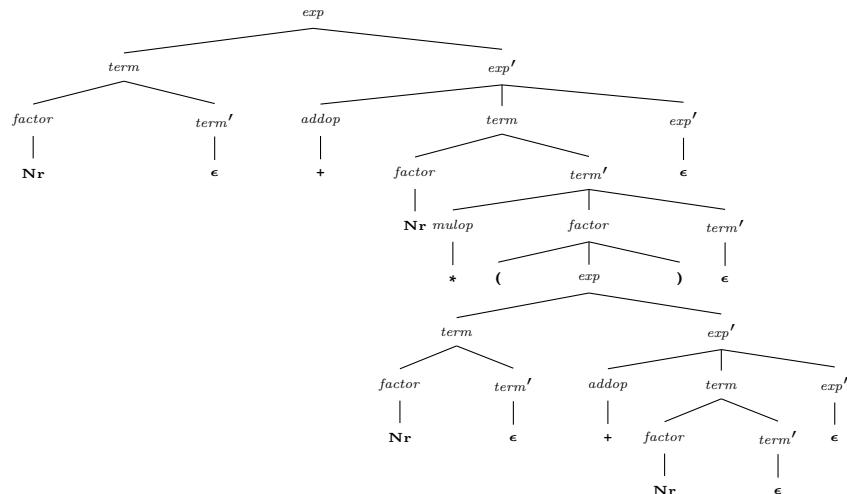
- underline: the *place* (occurrence of *non-terminal* where production is used)
  - ~~crossed out~~:

- *terminal = token* is considered treated
  - parser “moves on”
  - later implemented as **match** or **eat** procedure

Not as a “film” but at a glance: reduction sequence

exp  
term exp'  
factor term' exp'  
number term' exp'  
number term' exp'  
numbere exp'  
number exp'  
number addop term exp'  
number+term exp'  
number +term exp'  
number +factor term' exp'  
number +number term' exp'  
number +number term' exp'  
number + number mulop factor term' exp'  
number + number \*factor term' exp'  
number + number \*( exp ) term' exp'  
number + number \*{ exp } term' exp'  
number + number \*( exp ) term' exp'  
number + number \*( term exp' ) term' exp'  
number + number \*( factor term' exp' ) term' exp'  
number + number \*( number term' exp' ) term' exp'  
number + number \*( number term' exp' ) term' exp'  
number + number \*( numbere exp' ) term' exp'  
number + number \*( number exp' ) term' exp'  
number + number \*( number addop term exp' ) term' exp'  
number + number \*( number+term exp' ) term' exp'  
number + number \*( number +term exp' ) term' exp'  
number + number \*( number +factor term' exp' ) term' exp'  
number + number \*( number +number term' exp' ) term' exp'  
number + number \*( number + numbere exp' ) term' exp'  
number + number \*( number + number exp' ) term' exp'  
number + number \*( number + numbere ) term' exp'  
number + number \*( number + number ) term' exp'  
number + number \*( number + number )  $\epsilon$  exp'  
number + number \*( number + number ) exp'  
number + number \*( number + number )  $\epsilon$   
number + number \*( number + number )

Best viewed as a tree



## Non-determinism?

- not a “free” expansion/reduction/generation of some word, but
    - reduction of start symbol towards the *target word of terminals*

*exp*  $\Rightarrow^*$  1 + 2 \* (3 + 4)

- i.e.: input stream of tokens “guides” the derivation process (at least it fixes the target)
  - but: how much “guidance” does the target word (in general) gives?

Two principle sources of non-determinism here

1. Using production  $A \rightarrow \beta$

$$S \Rightarrow^* \alpha_1 \textcolor{red}{A} \alpha_2 \Rightarrow \alpha_1 \textcolor{red}{\beta} \alpha_2 \Rightarrow^* w$$

2. Conventions

- $\alpha_1, \alpha_2, \beta$ : word of terminals and nonterminals
- $w$ : word of terminals, only
- $A$ : one non-terminal

3. 2 choices to make

- (a) **where**, i.e., on **which occurrence of a non-terminal** in  $\alpha_1 A \alpha_2$  to apply a production<sup>3</sup>
- (b) **which production** to apply (for the chosen non-terminal).

### Left-most derivation

- that's the *easy* part of non-determinism
- taking care of “where-to-reduce” non-determinism: *left-most* derivation
- notation  $\Rightarrow_l$
- the example derivation earlier used that

### Non-determinism vs. ambiguity

- Note: the “where-to-reduce”-non-determinism  $\neq$  ambiguity of a grammar<sup>4</sup>
- in a way (“theoretically”): where to reduce next is *irrelevant*:
  - the order in the sequence of derivations *does not matter*
  - what does matter: the **derivation tree** (aka the **parse tree**)

**Lemma 8** (Left or right, who cares).  $S \Rightarrow_l^* w \iff S \Rightarrow_r^* w \iff S \Rightarrow^* w$ .

- however (“practically”): a (deterministic) parser implementation: must make a *choice*

1. Using production  $A \rightarrow \beta$

$$S \Rightarrow^* \alpha_1 \textcolor{red}{A} \alpha_2 \Rightarrow \alpha_1 \textcolor{red}{\beta} \alpha_2 \Rightarrow^* w$$

$$S \Rightarrow_l^* w_1 \textcolor{red}{A} \alpha_2 \Rightarrow w_1 \textcolor{red}{\beta} \alpha_2 \Rightarrow_l^* w$$

### What about the “which-right-hand side” non-determinism?

$$A \rightarrow \beta \mid \gamma$$

1. Is that the correct choice?

$$S \Rightarrow_l^* w_1 \textcolor{red}{A} \alpha_2 \Rightarrow w_1 \textcolor{red}{\beta} \alpha_2 \Rightarrow_l^* w$$

- reduction with “guidance”: don’t lose sight of the target  $w$ 
  - “past” is fixed:  $w = \textcolor{blue}{w}_1 w_2$
  - “future” is not:

$$\textcolor{red}{A} \alpha_2 \Rightarrow_l \textcolor{red}{\beta} \alpha_2 \Rightarrow_l^* w_2 \quad \text{or else} \quad \textcolor{red}{A} \alpha_2 \Rightarrow_l \textcolor{red}{\gamma} \alpha_2 \Rightarrow_l^* w_2 ?$$

2. Needed (minimal requirement): In such a situation, “future target”  $w_2$  must *determine* which of the rules to take!

---

<sup>3</sup>Note that  $\alpha_1$  and  $\alpha_2$  may contain non-terminals, including further occurrences of  $A$ .

<sup>4</sup>A CFG is ambiguous, if there exists a word (of terminals) with 2 different parse trees.

## Deterministic, yes, but still impractical

$$A\alpha_2 \Rightarrow_l \beta\alpha_2 \Rightarrow_l^* w_2 \quad \text{or else} \quad A\alpha_2 \Rightarrow_l \gamma\alpha_2 \Rightarrow_l^* w_2 ?$$

- the “target”  $w_2$  is of *unbounded length*!
- $\Rightarrow$  impractical, therefore:
1. Look-ahead of length  $k$  resolve the “which-right-hand-side” non-determinism inspecting only fixed-length prefix of  $w_2$  (for *all* situations as above)
  2. LL( $k$ ) grammars  
CF-grammars which *can* be parsed doing that.<sup>5</sup>

## Parsing LL(1) grammars

- *this lecture*: we don’t do LL( $k$ ) with  $k > 1$
  - LL(1): particularly easy to understand and to implement (efficiently)
  - not as expressive than LR(1) (see later), but still kind of decent
1. LL(1) parsing principle Parse from 1) left-to-right (as always anyway), do a 2) **left-most** derivation and resolve the “which-right-hand-side” non-determinism by looking 3) **1 symbol ahead**.
  2. Explanation
    - two flavors for LL(1) parsing here (both are top-down parsers)
      - *recursive descent*<sup>6</sup>
      - *table-based LL(1)* parser
    - *predictive* parsers

## Sample expr grammar again

### 1. factors and terms

$$\begin{array}{lcl}
 \text{exp} & \rightarrow & \text{term exp}' \\
 \text{exp}' & \rightarrow & \text{addop term exp}' \mid \epsilon \\
 \text{addop} & \rightarrow & + \mid - \\
 \text{term} & \rightarrow & \text{factor term}' \\
 \text{term}' & \rightarrow & \text{mulop factor term}' \mid \epsilon \\
 \text{mulop} & \rightarrow & * \\
 \text{factor} & \rightarrow & (\text{exp}) \mid \text{n}
 \end{array} \tag{6}$$

## Look-ahead of 1: straightforward, but not trivial

- look-ahead of 1:
    - not much of a look-ahead, anyhow
    - just the “current token”
- $\Rightarrow$  read the next token, and, based on that, decide
- but: what if there’s *no more symbols*?
- $\Rightarrow$  read the next token if there is, and decide based on the token *or else* the fact that there’s none left<sup>7</sup>
1. Example: 2 productions for non-terminal *factor*

$$\text{factor} \rightarrow (\text{exp}) \mid \text{number}$$

### 2. Remark that situation is *trivial*, but that’s not all to LL(1) ...

---

<sup>5</sup>Of course, one can always write a parser that “just makes some decision” based on looking ahead  $k$  symbols. The question is: will that allow to capture *all* words from the grammar and *only* those.

<sup>6</sup>If one wants to be very precise: it’s recursive descent with one look-ahead and without back-tracking. It’s the single most common case for recursive descent parsers. Longer look-aheads are possible, but less common. Technically, even allowing back-tracking can be done using recursive descent as principle (even if not done in practice).

<sup>7</sup>Sometimes “special terminal”  $\$$  used to mark the end (as mentioned).

## Recursive descent: general set-up

1. global variable, say `tok`, representing the “current token” (or pointer to current token)
2. parser has a way to *advance* that to the next token (if there’s one)
1. Idea For each *non-terminal nonterm*, write one procedure which:
  - succeeds, if starting at the current token position, the “rest” of the token stream starts with a syntactically correct word of terminals representing *nonterm*
  - fail otherwise
2. Rest
  - ignored (for right now): when doing the above successfully, build the *AST* for the accepted nonterminal.
3. Concepts (no pension, gratis, more esoteric details)

Feel free to ignore the following remarks.

(a) LL grammar

(b) LL(k) grammar

The definition of LL(K) grammar is inherently connected to the notion of *left-derivation* (or left-most derivation). One “L” in the LL stands for left-derivations (the other “L” means, the parser works from left-to-right, which is standard, anyhow). The relationship to predictive, recursive descent parsers is a bit unclear. The English Wikipedia states that LL-grammars correspond to the ones parsable via predictive recursive-descent parsers. That contradicts the statement that the correspondence is rather with LL(\*), which are grammars with unlimited look-ahead. Even of course the collection of all LL(k) languages imply unlimited look-head considering the whole class of languages, for me it seems still a difference (at least in the case that the look-ahead is not just dependent on the grammar, but (also) on the input word).

The definition of LL(k) is connected to grammars as primary concept, not languages (which is clear, since it’s about derivations). Furthermore, it’s based on being “deterministic” in the sense discussed: it’s not about arbitrary expansions, but expansions leading to a fixed result (which we want to parse). Now, LL(k) parsing is all about “predictive” parsing (top-down). Doing a left-derivation takes care of one source of non-determinism, namely picking the redex. Remains the second source of non-determinism, namely picking the right-hand side. Determinism means that at each point, there is only *one* correct choice (for a successful parse, of course). So one should better say, there are not *two* possible choices for a production.

Let’s have a look at the situation in a left-derivation:

$$S \Rightarrow_l^* wA\beta \Rightarrow_l w\alpha\beta .$$

The last step is the one where the right decision needs to be made. What needs to be decided is  $\alpha$  and the decision need to be determined by the **terminals following**  $w$ . Of course, in the post-state after the derivation step in the run of the parsers, the symbols following  $w$  may not be all terminals. So we need to derive on, until we obtain a string of terminals. Actually, the definition not just continues deriving until one obtains enough terminals, the reduction continues until all non-terminals are removed. That leads to the following definition.

**Definition 9** (LL(k)-grammar). Assume

$$S \Rightarrow_l^* wA\beta \left\{ \begin{array}{l} \Rightarrow_l w\alpha_1\beta \Rightarrow_l^* wv_1 \\ \Rightarrow_l w\alpha_2\beta \Rightarrow_l^* wv_2 \end{array} \right.$$

If  $\text{first}_k(v_1) = \text{first}_k(v_2)$ , then  $\alpha_1 = \alpha_2$ .

The definition of *first* here (used for strings of terminals) is rather intuitive and simply consists in the first  $k$  terminal symbols, if possible. If the string is shorter, then, obviously, the whole word is taken. Note that the definition here, for a string of terminals, is a set, even if that seems unnecessary, since it’s simply a singleton set. Only in the more general setting when

considering non-terminals as well, we need to generalize it to larger sets. Compare that also to our definition of the first sets, which is slightly different.

The differences seem more in the technical details. In the (recursive) definition of the first-sets used in the lecture it was important to include the  $\epsilon$ . The reason for that was basically to be able to formulate straightforwardly a *recursive* definition of the first-set. For that it is important to include the information if a non-terminal is *nullable*. Nonetheless, the inclusion of  $\epsilon$  in the set itself is a bit of an “anomaly”, resp. one could have done it differently. Definition 10 of the first sets covers all words in sentential form (including the empty one). In that way, this definition is more general as the one we used before. Technically, Definition 10 distinguishes between words using terminal only, and words containing at least one non-terminal, both cases are disjoint.

Actually, for  $k = 1$  and for one single non-terminal symbol, the definition here and the earlier one are *almost* the same. The only subtle difference is that case (3(b)ii) here requires that  $\alpha$  can be reduced to a word of *terminals*. In our definition of the lecture, that’s not required. That’s a minor point, and for “decent” grammars, there is *no* difference, because it’s a reasonable requirement that all non-terminal symbols can be reduced to a word of terminals. Under that assumption, both definitions are indeed identical for  $k = 1$ . In particular, the inclusion of  $\epsilon$  in the result for nullable nonterminal symbols is covered!

We should also remember, that the lecture did indeed cover first sets of words already (but not as primary definition). That was done in Definition 5, which was not based on  $\Rightarrow_l^*$ , but rather a recursive definition, generalizing the recursive characterization from Definition 3.

**Definition 10** (First-set of a word). Assume a grammar  $G$ .

- i. Let  $w \in \Sigma_T^*$ . Then  $First_k(w) = \{w\}$  if  $|w| \leq k$ , otherwise  $First_k(w) = \{v \mid w = vv', |v| = k\}$ .
- ii. Let  $\alpha \in (\Sigma_T + \Sigma_N)^* \setminus \Sigma_T^*$ . Then  $First_k(\alpha) = \{v \in \Sigma_T^* \mid \alpha \Rightarrow_G^* w, First_k(w) = v\}$ .

The following is a characterization of LL( $k$ )-grammars.

**Lemma 11** (LL( $k$ )). A context-free grammar is an LL( $k$ )-grammar iff for all non-terminals  $A$  and for all pairs of productions  $A \rightarrow \alpha_1$  and  $A \rightarrow \alpha_2$  with  $\alpha_1 \neq \alpha_2$  and  $S \Rightarrow_l^* wA\beta$ :

$$First_k(\alpha_1\beta) \cap First_k(\alpha_2\beta) = \emptyset .$$

It’s not much different from the original definition, though. One difference is that in the definition, one reduces to words of terminals (and consequently, the simpler form of the first-sets is used, from Definition 10(3(b)i)). In the characterization here, the reduction is to arbitrary words. In both cases, left derivations are required, of course.

Like the first-sets, also the follow sets can be generalized to  $k > 1$ . In the lecture, we actually defined the follow set only for non-terminals. One can generalize that for larger  $k$ ’s and for general words.

**Definition 12** (Follow-set of a word). Assume a grammar  $G$  and a  $\alpha \in (\Sigma_T + \Sigma_N)^*$ . Then

$$Follow_k(\alpha) = \{w \mid S \Rightarrow_l^* \beta_1\alpha\beta_2, w \in First_k(\beta_2)\} .$$

In Definition 6 from the lecture earlier (where we ignored  $\$$  and where  $k = 1$ ), the follow set was done as follows: start from the start-symbol to reach the non-terminal  $A$  in question. All terminal symbols which can occur directly after  $A$ , that’s the follow-set (ignoring, as said,  $\$$ ). The definition here uses words  $\alpha$  instead of non-terminals, which does in itself not change much. Secondly, since we are generalizing  $k$ , the first set consists of words (from  $\Sigma_T^*$ , and of length  $\leq k$ ), not letters. Note that it’s not the case that it consists of words of exactly length  $k$ . Not also: the current definiton implicitly still treats the case that there is *no* symbol following. In the previous definition, this was indicated by a special marker  $\$$ , representing the end-of-parse. Here, at least for  $k > 0$ , the fact that  $\epsilon \in Follow_k(\alpha)$  means that the end of the parse is reachable. Note that the definition here is formulated slightly differently in one aspect, namely relying on  $First_k$  as auxiliary definition. It does not make a technical difference though.

**Lemma 13** (LL(1)). A reduced context-free grammar is an LL(1)-grammar iff for all non-terminals  $A$  and for all pairs of productions  $A \rightarrow \alpha_1$  and  $A \rightarrow \alpha_2$  with  $\alpha_1 \neq \alpha_2$ :

$$First_1(\{\alpha_1\}Follow_1(A)) \cap First_1(\{\alpha_2\}Follow_1(A)) = \emptyset .$$

(c) Reduced CFG

It seems as follows. A reduced grammar is one where one does not have superfluous symbols. A symbol (terminal or non-terminal) is superfluous, if it cannot be reached from the start symbol or that, starting from it, one cannot reach a string of  $\Sigma_T^*$ . The latter condition is non-trivial only for non-terminals, of course. A symbol satisfying those conditions is called *reachable* resp. *terminating*. As a side remark. Considered as reduction, the latter condition would more precisely be called *weakly terminating*.

(d) LL(\*) and unbounded look-ahead

Each LL(k) grammar operates with a fixed maximal look-ahead of  $k$ . There is the notion of LL(\*)-grammars, which is unclear. It's said that it corresponds to unbounded look-ahead. For LL(k) grammars it's clear that the  $k$  is per grammar (or language). Of course the collection of all LL(k) grammars do not have an upper bound on the look-ahead (since the hierarchy is infinite). But still a question is, if this is the same as having *one* grammar with unbounded look-ahead. There is a difference between LL(k) and LL(\*). The latter class is known as *LL-regular grammars*. Basically, the parsing can check if the following token(s) belong to a regular language.

(e) LL(1) grammar

LL(1) grammars are grammars whose *predict sets* are disjoint (for the same non-terminal) is LL(1). It seems that this is not the primary definition (as LL(1) grammars are defined via LL(k) grammars).

(f) Recursive descent and predictive parsing

Actually, recursive descent parsing and predictive parsing are *not* the same in general. Predictive parsing seems to be defined as a subclass of recursive descent parsing (at least according to Wikipedia), namely when *backtracking* is not needed. Since backtracking for parsing is a bad idea, it might explain that one sometimes has the impression that the notions recursive-descent parsing and predictive parsing are used synonymously. Actually, it seems that predictive parsing and recursive descent is not only often discussed as if they were the same, but also as if the prediction is based on just one next (the current) symbol.

To make it explicit (even if perhaps needless to say, perhaps): predictive parsing also does not mean that one need one *one lookahead* to process. In the lecture we deal only with the simplest case of predictive parsing, as illustration for recursive descent parsing, namely for one look-ahead.

(g) Left-derivation and non-determinism.

Grammars, when considered as language generators, "generate" languages, resp. one word of its language, but using sequences of productions, it's a form of rewriting, starting from the grammar's start symbol (even if there are differences, especially in the purpose, between rewriting and expanding grammar productions). This expansion or generation process is basically always non-deterministic (and non-confluent), for obvious reasons. For using productions of a CF grammar, there are exactly 2 sources of non-determinism (and the situation is basically analogous for term rewriting). One may be caused by the form of the grammar, namely there are two rules with the same non-terminal on the left-hand side. When assuming that a grammar contains a *set* of productions (as opposed to a multi-set), two different rules with the same non-terminal left-hand side have two different right-hand sides (which is the natural thing to assume anyhow, it makes no sense to write the "same" rule twice). Under this assumption, of course, applying two different rules with the same left-hand side leads to two different successor states.

The second source of non-determinism is due to the fact that a word in sentential form may contain more than one non-terminal, i.e., more than one place where a production applies. In the context of rewriting, it means there is more than one *redex*. This kind of non-determinism is likewise basically unavoidable in most CFG. It can only be avoided, if the productions are *linear*.

It should be noted that this picture and discussion related to non-determinism is *misleading*. For rewriting, the intermediate forms are interpreted as "states" and non-determinism means that different states are reached (or at least states which are then not reconcilable via *confluence*). In particular, the derivation is more seen as a sequence. Let's call the steps under the rewriting perspective "reductions" and in the perspective of generative grammars "expansions" or "generation".

For grammar derivation, the situation rather different: the process may be performed in a sequence of steps, but conceptually is seen as a *tree* (the *parse tree*). However, again the parse tree is a simplification or rather an *abstraction* of the derivation process, it's not the same. Given a sentential form which contains 2 different non-terminals, then obviously one can apply a rule at two different places (of course assuming that each non-terminal is covered by at least one rule, which is another obvious common-sense well-formedness assumption on grammars). As far as reductions go, the situation corresponds to *parallel reductions* because there's no overlap. That's a general characteristic of context-free grammars. As a consequence, the grammar generation would be confluent (seen as reduction), *where it not for for the non-determinism in the productions itself*. If the productions were *functional* in that each non-terminal had only one right-hand side, CFG as reduction system would be confluent.

Of course, a grammars of that form would not be of much use. as there would mostly generate the empty language (from the perspective of language generation, not sentential reduction). The reason is that any meaningful context-free grammar will contain recursion (otherwise the language would be finite; not even full regular languages can be captured without recursion, namely with left-linear, resp. right-linear grammars). But *with* recursion and without the possibility of having the “alternative” to exit the “loop” the generation process will never be able to terminate.

There is another difference between reduction and generation. Perhaps that one is not so relevant, but in general, rewriting is understood as *term rewriting* not string rewriting. Of course, string writing can be seen as a special case of term rewriting. In term-rewriting, especially when dealing with using rewriting for providing evaluation semantics for programming languages, the issue of determinism (and confluence) resp. the lack thereof, also shows up. In that context, there are strategies like left-most *inner-most* or left-most *outer-most* or similar. Since there, one is dealing with *terms* or *trees*, this way of specifying a reduction strategy makes sense. In a linear setting with strings, of course, things are simpler: one can only distinguish between left-most reduction vs. right-most reduction (as the two obvious deterministic strategies). If we consider strings as special case of trees or terms, then one need to fix associativity of concatenation. The natural choice seems to be right-associativity. In that, given an alphabet  $\Sigma$ , each symbol is treated as a function symbol of arity 1, and one needs an additional constant symbol to denote the end of the string/tree. In that way, the first symbol on the right-hand end of the string is the *root* of the tree. However, in that way, there is no good match of grammar expansion with term rewriting, because in a grammar, the *non-terminals* are replaced. In the string-as-term setting, that would mean that variables could stand for “function symbols”, something which is not supported. In term-rewriting, variables are of arity 0, they cannot have arguments. It seems that it could be covered by higher-order rewriting (but it's not 100% clear).

Finally, and that's perhaps the really crucial difference between reduction in rewriting and grammatical expansion wrt. of the question of determinism is the following. Determinism in the grammar setting is not about if the expansion leads to a definite endpoint (when each steps is “determined” in itself, which will not be the case for standard grammars). It's in contrast that the *end-point* of the expansion is *a priori* fixed. So the question is not if there are more than one end-point or not, it's whether, given a fixed endpoint, there are two different ways to get there.

In principle, with the exception of uninteresting grammars, there are always different “ways to get there”, when interpreted as sequences of reduction or derivation steps. The order of applying productions is *irrelevant* which is also the reason why the core of a derivation it not the sequence of steps, but the abstraction thereof, which is the *parse tree*. It's the uniqueness of the derivation *tree* which matters, not the uniqueness of the derivation sequence (which mostly won't hold anyway). Uniqueness of derivation or parse trees is *unambiguity* of the grammar.

It seems clear that ambiguous grammars cannot be used for predictive top-down parsing. For recursive descent parsing, it's not so clear, basically because recursive descent parsing seems not to be too clearly defined right now and seems to allow back-tracking. If one can do that, one can explore alternatives, so perhaps even ambiguous grammars can be treated by general recursive descent parsers. What is clear is that at least LL( $k$ ) grammars cannot capture ambiguous languages.

- (h) Predictive recursive descent with 1 look-ahead

As said, recursive descent parsing does not imply predictive parsing. Consequently, it's clear that it's more powerful than LL(k) parsing. Also, *predictive* parsing does not mean one has just one look-ahead. Nonetheless, the lecture gives the impression that recursive descent parsers work on a look-ahead of size one. Probably that's because in that case, the implementation is rather simple. In a language with more advanced pattern matching, also a look-ahead of fixed maximal size would be easy to implement, as well.

(i) Higher-order rewriting

Some material can be found in the van der Pol's Ph.D thesis. It seems that higher-order rewriting does more than just introducing variables with arity, namely mechanisms dealing with boundness of variables. The notion of HRS is very much connected to the  $\lambda$ -calculus. However, having variables with arities *is* indeed a feature for higher-order rewrite systems!

(j) Term-rewriting and transduction

(k) Tree transducers

(l) Terminals vs. Non-terminals.

The discussion here is more like speculation.

In the discussions about left-derivations, the issue of non-termination came up. The reason why non-termination was an issue there is that one basic rule for CFGs says that the terminals and non-terminals are disjoint, and to be accepted in the language, the sentential sentence must consist of terminal, only. That's of course understandable. The non-terminals are used for the internal representation of the grammar, in the same way that for FSA, the states (which correspond to non-terminals) are of course strictly separate from the letters, an automaton accepts.

Nonetheless: if one had non-terminals which simply represent a terminal "as well", one would have introduced something like "accepting states". If one restricts attention to *left-linear* grammars, which correspond to automata, then accepting states could be represented by  $\epsilon$ -productions. Of course, that leads then to non-determinism, which is exactly what is needed to find the exit in an infinite recursion.

So, there is still some non-determinism ("continue-or-stop") but it seems less harmful than general forms of non-determinism. So one may wonder, if that form of grammar is useful. It would a form which would be deterministic, except that  $\epsilon$ -productions are allowed. That somehow corresponds to a situation, where some of the non-terminals are "also terminals".

Perhaps that has not been studied much, as one typically want to get rid of  $\epsilon$ -productions. It can be shown that all CFG languages can be written *without*  $\epsilon$ -productions with one exception: if the language contains  $\epsilon$ , then that obviously does not work.

(m) Predict sets

*Predict-sets* are often mentioned in connection with recursive-descent parsing.

(n) LL(k) parsing and recursive descent

For a discussion, see stackoverflow or the link here. There are also recursive ascent parsers which correspond to LALR grammars. There are also special forms of grammars, tailor-made for top-down parsers, so called PEG (parsing expression grammars). For instance they cannot be ambiguous. Recursive-descent parsers are more expressive than LL(k) parsers, they correspond to LL(\*)-parsers (unlimited look-ahead). However, it can mean that what is discussed there as recursive descent is something different from the way it's covered in the lecture, because here it seems more like 1-look-ahead.

It is also said that recursive descent is an early implementation for LL(1) grammars.

### Recursive descent method factor for nonterminal factor

```
final int LPAREN=1,RPAREN=2,NUMBER=3,
PLUS=4,MINUS=5,TIMES=6;
```

```
void factor () {
    switch (tok) {
        case LPAREN: eat(LPAREN); expr (); eat(RPAREN);
        case NUMBER: eat(NUMBER);
    }
}
```

## Recursive descent

```
type token = LPAREN | RPAREN | NUMBER
| PLUS | MINUS | TIMES
```

```
let factor () = (* function for factors *)
match !tok with
| LPAREN -> eat(LPAREN); expr(); eat(RPAREN)
| NUMBER -> eat(NUMBER)
```

## Slightly more complex

- previous 2 rules for *factor*: situation not always as immediate as that
- 1. LL(1) principle (again) given a non-terminal, the next *token* must determine the choice of right-hand side<sup>8</sup>
- 2. First  
⇒ definition of the *First set*

**Lemma 14** (LL(1) (without nullable symbols)). *A reduced context-free grammar without nullable non-terminals is an LL(1)-grammar iff for all non-terminals A and for all pairs of productions  $A \rightarrow \alpha_1$  and  $A \rightarrow \alpha_2$  with  $\alpha_1 \neq \alpha_2$ :*

$$First_1(\alpha_1) \cap First_1(\alpha_2) = \emptyset.$$

## Common problematic situation

- often: common *left factors* problematic

$$\begin{array}{lcl} if-stmt & \rightarrow & \text{if ( exp ) stmt} \\ & | & \text{if ( exp ) stmt else stmt} \end{array}$$

- requires a look-ahead of (at least) 2
  - ⇒ try to rearrange the grammar
1. *Extended BNF* ([Louden, 1997] suggests that)

$$if-stmt \rightarrow \text{if ( exp ) stmt [else stmt]}$$

1. *left-factorizing*:

$$\begin{array}{lcl} if-stmt & \rightarrow & \text{if ( exp ) stmt else-part} \\ else-part & \rightarrow & \epsilon \mid \text{else stmt} \end{array}$$

## Recursive descent for left-factored *if-stmt*

```
procedure ifstmt()
begin
  match ("if");
  match "(";
  exp();
  match ")";
  stmt();
  if token = "else"
  then match ("else");
    stmt()
  end
end;
```

---

<sup>8</sup>It must be the next token/terminal in the sense of *First*, but it need not be a token *directly* mentioned on the right-hand sides of the corresponding rules.

## Left recursion is a no-go

### 1. factors and terms

$$\begin{aligned}
 \text{exp} &\rightarrow \text{exp addop term} \mid \text{term} \\
 \text{addop} &\rightarrow + \mid - \\
 \text{term} &\rightarrow \text{term mulop factor} \mid \text{factor} \\
 \text{mulop} &\rightarrow * \\
 \text{factor} &\rightarrow (\text{exp}) \mid \text{number}
 \end{aligned} \tag{7}$$

### 2. Left recursion explanation

- consider treatment of  $\text{exp}$ :  $\text{First}(\text{exp})$ ?
  - whatever is in  $\text{First}(\text{term})$ , is in  $\text{First}(\text{exp})$ <sup>9</sup>
  - even if only *one* (left-recursive) production  $\Rightarrow$  *infinite* recursion.

### 3. Left-recursion Left-recursive grammar *never* works for recursive descent.

## Removing left recursion may help

### 1. Pseudo code

$$\begin{aligned}
 \text{exp} &\rightarrow \text{term exp}' \\
 \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\
 \text{addop} &\rightarrow + \mid - \\
 \text{term} &\rightarrow \text{factor term}' \\
 \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\
 \text{mulop} &\rightarrow * \\
 \text{factor} &\rightarrow (\text{exp}) \mid \text{n}
 \end{aligned}$$

```

procedure exp()
begin
    term();
    exp'()
end

```

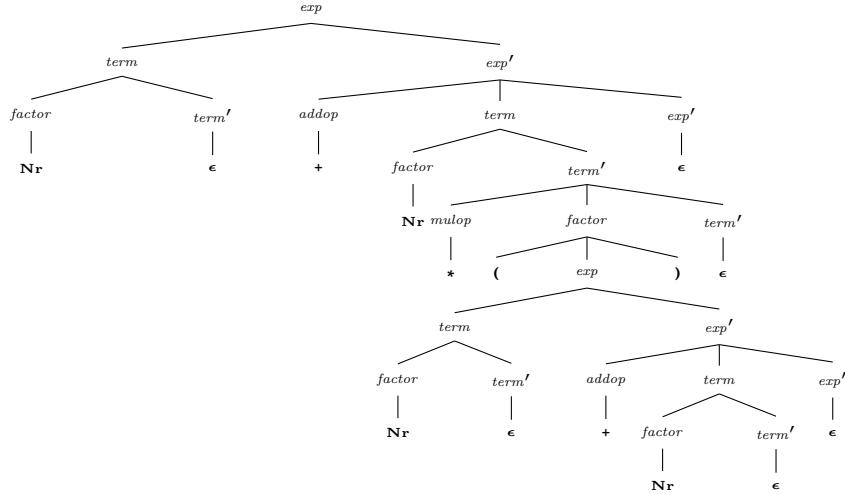
```

procedure exp'()
begin
  case token of
    "+": match("+");
    term();
    exp'()
    "-": match("-");
    term();
    exp'()
  end
end

```

## Recursive descent works, alright, but ...

<sup>9</sup> And it would not help to *look-ahead* more than 1 token either.



... who wants this form of trees?

## The two expression grammars again

### 1. factors and terms

#### (a) Precedence & assoc.

$$\begin{array}{lcl}
 \text{exp} & \rightarrow & \text{exp addop term} \mid \text{term} \\
 \text{addop} & \rightarrow & + \mid - \\
 \text{term} & \rightarrow & \text{term mulop factor} \mid \text{factor} \\
 \text{mulop} & \rightarrow & * \\
 \text{factor} & \rightarrow & (\text{exp}) \mid \text{number}
 \end{array}$$

#### (b) explanation

- clean and straightforward rules
- left-recursive

### 2. no left recursion

#### (a) no left-rec.

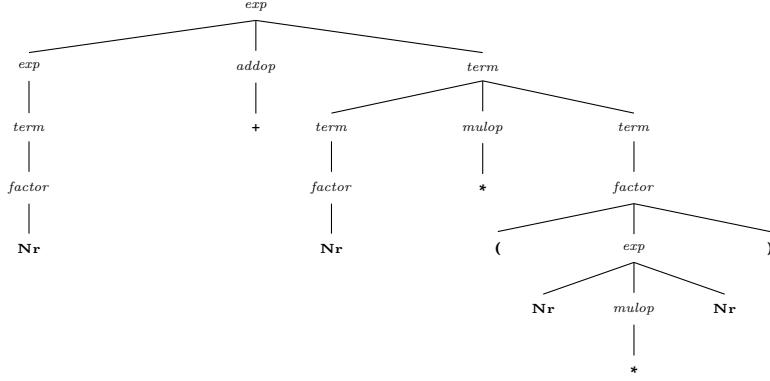
$$\begin{array}{lcl}
 \text{exp} & \rightarrow & \text{term exp}' \\
 \text{exp}' & \rightarrow & \text{addop term exp}' \mid \epsilon \\
 \text{addop} & \rightarrow & + \mid - \\
 \text{term} & \rightarrow & \text{factor term}' \\
 \text{term}' & \rightarrow & \text{mulop factor term}' \mid \epsilon \\
 \text{mulop} & \rightarrow & * \\
 \text{factor} & \rightarrow & (\text{exp}) \mid \text{n}
 \end{array}$$

#### (b) Explanation

- no left-recursion
- assoc. / precedence ok
- rec. descent parsing ok
- but: just “unnatural”
- non-straightforward parse-trees

## Left-recursive grammar with nicer parse trees

$$1 + 2 * (3 + 4)$$



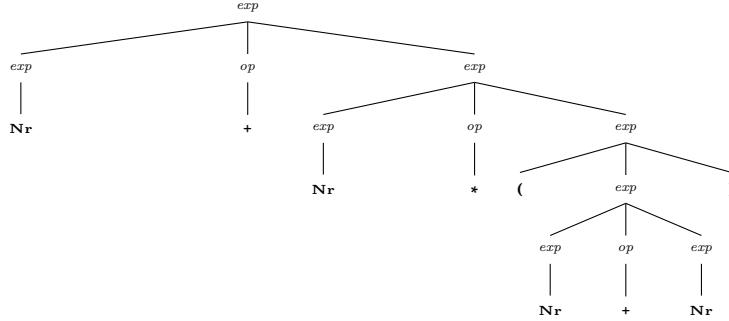
### The simple “original” expression grammar (even nicer)

#### 1. Flat expression grammar

$$\begin{aligned} \text{exp} &\rightarrow \text{exp } \text{op } \text{exp} \mid (\text{exp}) \mid \text{number} \\ \text{op} &\rightarrow + \mid - \mid * \end{aligned}$$

#### 2. Nice tree

$$1 + 2 * (3 + 4)$$



### Associativity problematic

#### 1. Precedence & assoc.

$$\begin{aligned} \text{exp} &\rightarrow \text{exp } \text{addop } \text{term} \mid \text{term} \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{term } \text{mulop } \text{factor} \mid \text{factor} \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow (\text{exp}) \mid \text{number} \end{aligned}$$

#### 2. Example plus and minus

##### (a) Formula

$$3 + 4 + 5$$

parsed “as”

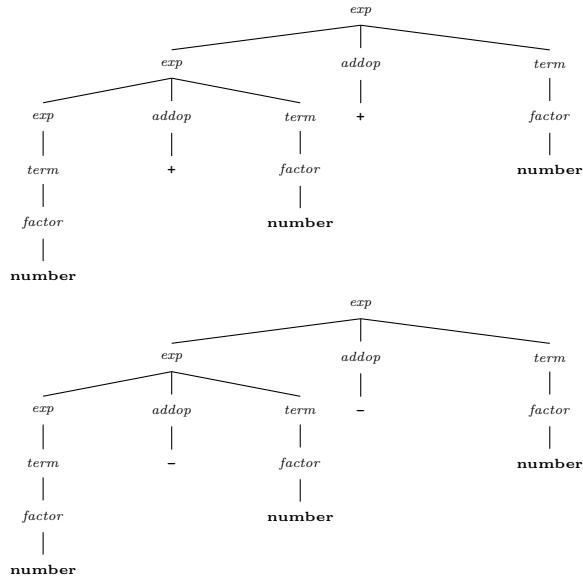
$$(3 + 4) + 5$$

$$3 - 4 - 5$$

parsed “as”

$$(3 - 4) - 5$$

(b) Tree



Now use the grammar without left-rec (but right-rec instead)

1. No left-rec.

$$\begin{aligned}
 \text{exp} &\rightarrow \text{term exp}' \\
 \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\
 \text{addop} &\rightarrow + \mid - \\
 \text{term} &\rightarrow \text{factor term}' \\
 \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\
 \text{mulop} &\rightarrow * \\
 \text{factor} &\rightarrow (\text{exp}) \mid \text{n}
 \end{aligned}$$

2. Example minus

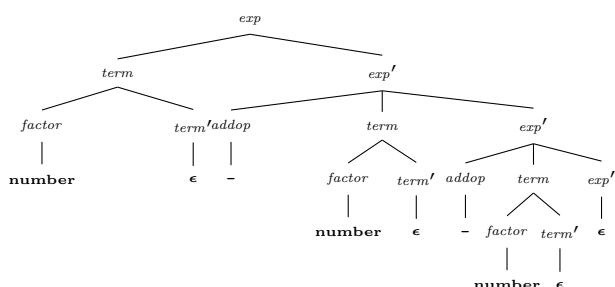
(a) Formula

$$3 - 4 - 5$$

parsed “as”

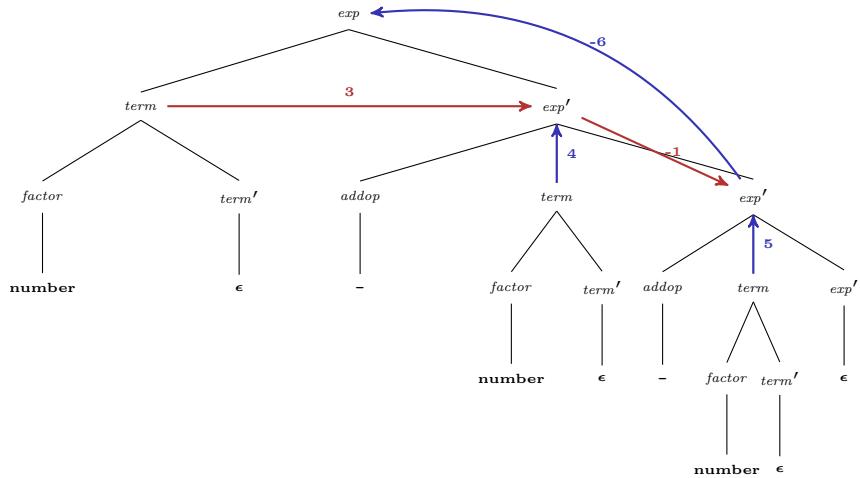
$$3 - (4 - 5)$$

(b) Tree



But if we need a “left-associative” AST?

- we want  $(3 - 4) - 5$ , not  $3 - (4 - 5)$



Code to “evaluate” ill-associated such trees correctly

```

function exp' (valsofar: int): int;
begin
  if token = '+' or token = '-'
  then
    case token of
      '+': match ('+');
              valsofar := valsofar + term;
      '-': match ('-');
              valsofar := valsofar - term;
    end case;
    return exp'(valsofar);
  else return valsofar
end;
  
```

- extra “accumulator” argument `valsofar`
- instead of evaluating the expression, one could build the AST with the appropriate associativity instead:
- instead of `valueSoFar`, one had `rootOfTreeSoFar`

“Designing” the syntax, its parsing, & its AST

- trade offs:
  1. starting from: design of the language, how much of the syntax is left “implicit”<sup>10</sup>
  2. which language class? Is LL(1) good enough, or something stronger wanted?
  3. how to parse? (top-down, bottom-up, etc.)
  4. parse-tree/concrete syntax trees vs. ASTs

### AST vs. CST

- once steps 1.–3. are fixed: *parse-trees* fixed!
- parse-trees = *essence* of grammatical derivation process
- often: parse trees only “conceptually” present in a parser
- AST:
  - *abstractions* of the parse trees

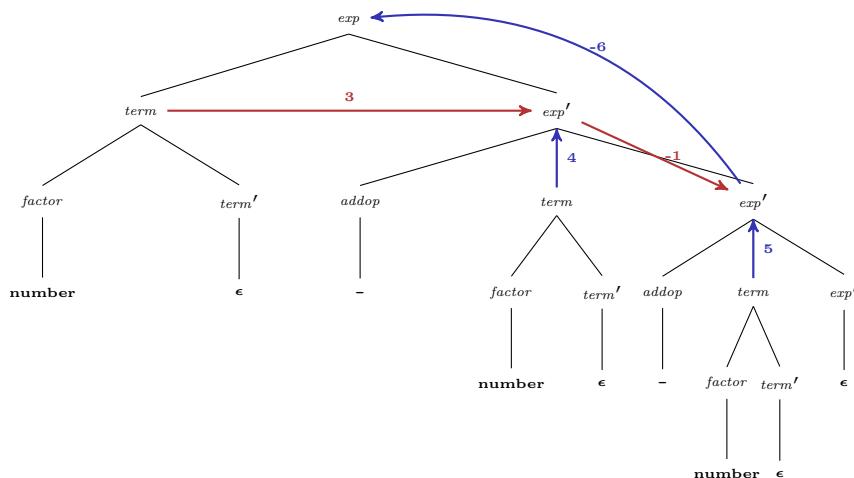
<sup>10</sup>Lisp is famous/notorious in that its surface syntax is more or less an explicit notation for the ASTs. Not that it was originally planned like this ...

- essence of the parse tree
- actual tree data structure, as output of the parser
- typically on-the fly: AST built while the parser parses, i.e. while it executes a derivation in the grammar

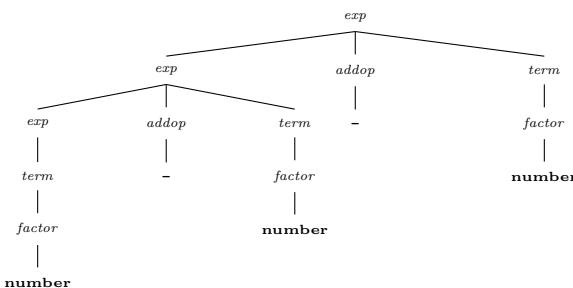
1. AST vs. CST/parse tree Parser "builds" the AST data structure while "doing" the parse tree

### AST: How “far away” from the CST?

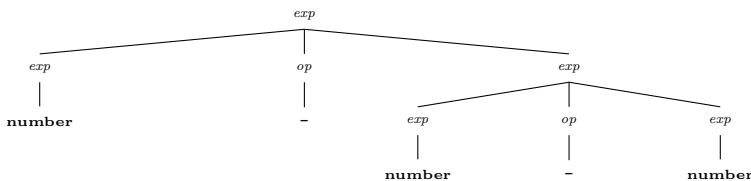
- AST: only thing relevant for later phases  $\Rightarrow$  better be *clean* ...
- AST “=” CST?
  - building AST becomes straightforward
  - possible choice, if the grammar is not designed “weirdly”,



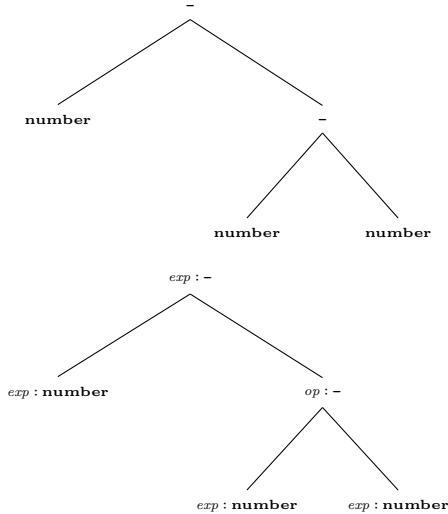
parse-trees like that better be cleaned up as AST



slightly more reasonable looking as AST (but underlying grammar not directly useful for recursive descent)



That parse tree looks reasonable clear and intuitive



### Wouldn't that be the best AST here?

Certainly minimal amount of nodes, which is nice as such. However, what is missing (which might be interesting) is the fact that the 2 nodes labelled “-” are *expressions*!

### This is how it's done (a recipe)

1. Assume, one has a “non-weird” grammar  $\begin{array}{l} exp \rightarrow exp \ op \ exp \mid (exp) \mid number \\ op \rightarrow + \mid - \mid * \end{array}$
2. Explanation
  - typically that means: assoc. and precedences etc. are fixed *outside* the non-weird grammar
    - by massaging it to an equivalent one (no left recursion etc.)
    - or (better): use parser-generator that allows to *specify* assoc ... like “ “\*” binds stronger than “+”, it *associates* to the left ...”
  - „ without cluttering the grammar.
  - if grammar for *parsing* is not as clear: do a second one describing the ASTs
3. Remember (independent from parsing) BNF describe **trees**

### This is how it's done (recipe for OO data structures)

1. Recipe
  - turn each **non-terminal** to an **abstract class**
  - turn each **right-hand** side of a given non-terminal as (non-abstract) **subclass** of the class for considered non-terminal
  - chose fields & constructors of concrete classes appropriately
  - **terminal**: concrete class as well, field/constructor for token's *value*

### Example in Java

```

$$\begin{array}{l} exp \rightarrow exp \ op \ exp \mid (exp) \mid number \\ op \rightarrow + \mid - \mid * \end{array}$$

```

```
abstract public class Exp { }
```

```
public class BinExp extends Exp { // exp -> exp op exp
    public Exp left, right;
    public Op op;
    public BinExp(Exp l, Op o, Exp r) {
        left=l; op=o; right=r;}
}
```

```

public class ParentheticalExp extends Exp { // exp -> ( op )
    public Exp exp;
    public ParentheticalExp(Exp e) {exp = e;}
}

```

```

public class NumberExp extends Exp { // exp -> NUMBER
    public number; // token value
    public Number(int i) {number = i;}
}

```

```

abstract public class Op { // non-terminal = abstract
}

```

```

public class Plus extends Op { // op -> "+"
}

```

```

public class Minus extends Op { // op -> "-"
}

```

```

public class Times extends Op { // op -> "*"
}

```

3 – (4 – 5)

```

Exp e = new BinExp(
    new NumberExp(3),
    new Minus(),
    new BinExp(new ParentheticalExpr(
        new NumberExp(4),
        new Minus(),
        new NumberExp(5))))
}

```

### Pragmatic deviations from the recipe

- it's nice to have a guiding principle, but no need to carry it too far ...
- To the very least: the `ParentheticalExpr` is completely without purpose: grouping is captured by the tree structure  
⇒ that class is *not* needed
- some might prefer an implementation of

$$op \rightarrow + \mid - \mid *$$

as simply integers, for instance arranged like

```

public class BinExp extends Exp { // exp -> exp op exp
    public Exp left, right;
    public int op;
    public BinExp(Exp l, int o, Exp r) {pos=p; left=l; oper=o; right=r;}
    public final static int PLUS=0, MINUS=1, TIMES=2;
}

```

and used as `BinExpr.PLUS` etc.

### Recipe for ASTs, final words:

- space considerations for AST representations are irrelevant in most cases
- clarity and cleanliness trumps “quick hacks” and “squeezing bits”
- some deviation from the recipe or not, the advice still holds:

1. Do it systematically A clean grammar is **the** specification of the syntax of the language and thus the parser. It is also a means of **communicating** with humans (at least with pros who (of course) can read BNF) what the syntax is. A clean grammar is a very systematic and structured thing which consequently *can* and *should* be **systematically** and **cleanly** represented in an AST, including judicious and systematic choice of names and conventions (nonterminal *exp* represented by class **Exp**, non-terminal *stmt* by class **Stmt** etc)

2. Louden

- a word on [Louden, 1997]: His C-based representation of the AST is a bit on the “bit-squeezing” side of things ...

### Extended BNF may help alleviate the pain

1. BNF

$$\begin{array}{l} \textit{exp} \rightarrow \textit{exp addop term} \mid \textit{term} \\ \textit{term} \rightarrow \textit{term mulop factor} \mid \textit{factor} \end{array}$$

2. EBNF

$$\begin{array}{l} \textit{exp} \rightarrow \textit{term}\{ \textit{addop term} \} \\ \textit{term} \rightarrow \textit{factor}\{ \textit{mulop factor} \} \end{array}$$

3. Explanation but remember:

- EBNF just a notation, just because we do not see (left or right) recursion in { ... }, does not mean there is no recursion.
- not all parser generators support EBNF
- however: often easy to translate into loops- <sup>11</sup>
- does not offer a *general* solution if associativity etc. is problematic

### Pseudo-code representing the EBNF productions

```
procedure exp;
begin
  term;           { recursive call }
  while token = "+" or token = "-"
  do
    match(token);
    term;          // recursive call
  end
end
```

```
procedure term;
begin
  factor;         { recursive call }
  while token = "*"
  do
    match(token);
    factor;        // recursive call
  end
end
```

### How to produce “something” during RD parsing?

1. Recursive descent So far: RD = top-down (parse-)tree traversal via recursive procedure.<sup>12</sup> Possible outcome: termination or failure.
2. Rest

<sup>11</sup>That results in a parser which is somehow not “pure recursive descent”. It’s “recusive descent, but sometimes, let’s use a while-loop, if more convenient concerning, for instance, associativity”

<sup>12</sup>Modulo the fact that the tree being traversed is “conceptual” and not the input of the traversal procedure; instead, the traversal is “steered” by stream of tokens.

- Now: instead of returning “nothing” (return type `void` or similar), return some meaningful, and build that up during traversal
- for illustration: procedure for expressions:
  - return type `int`,
  - while traversing: *evaluate* the expression

### Evaluating an *exp* during RD parsing

```
function exp() : int;
var temp: int
begin
  temp := term();           { recursive call }
  while token = "+" or token = "-"
    case token of
      "+": match ("+");
      temp := temp + term();
      "-": match ("-")
      temp := temp - term();
    end
  end
  return temp;
end
```

### Building an AST: expression

```
function exp() : syntaxTree;
var temp, newtemp: syntaxTree
begin
  temp := term();           { recursive call }
  while token = "+" or token = "-"
    case token of
      "+": match ("+");
      newtemp := makeOpNode("+");
      leftChild(newtemp) := temp;
      rightChild(newtemp) := term();
      temp := newtemp;
      "-": match ("-")
      newtemp := makeOpNode("-");
      leftChild(newtemp) := temp;
      rightChild(newtemp) := term();
      temp := newtemp;
    end
  end
  return temp;
end
```

- note: the use of `temp` and the `while` loop

### Building an AST: factor

$$factor \rightarrow ( exp ) \mid number$$

```
function factor() : syntaxTree;
var fact: syntaxTree
begin
  case token of
    "(": match "(");
    fact := exp();
    match ")";
    number:
    match (number)
    fact := makeNumberNode(number);
    else : error ... // fall through
  end
  return fact;
end
```

### Building an AST: conditionals

$$if-stmt \rightarrow if ( exp ) stmt [else stmt]$$

```

function ifStmt() : syntaxTree;
var temp: syntaxTree
begin
  match ("if");
  match "(");
  temp := makeStmtNode("if")
  testChild(temp) := exp();
  match ")";
  thenChild(temp) := stmt();
  if token = "else"
  then match "else";
    elseChild(temp) := stmt();
  else elseChild(temp) := nil;
  end
  return temp;
end

```

### Building an AST: remarks and “invariant”

- LL(1) requirement: each procedure/function/method (covering one specific non-terminal) decides on alternatives, looking only at the current token
- call of function A for non-terminal A:
  - upon entry: first terminal symbol for A in `token`
  - upon exit: first terminal symbol *after* the unit derived from A in `token`
- `match("a")` : checks for "a" in `token` and eats the token (if matched).

### LL(1) parsing

- remember LL(1) grammars & LL(1) parsing principle:
  1. LL(1) parsing principle 1 look-ahead enough to resolve “which-right-hand-side” non-determinism.
  2. Further remarks
    - instead of recursion (as in RD): *explicit stack*
    - decision making: collated into the **LL(1) parsing table**
    - LL(1) parsing table:
      - finite data structure  $M$  (for instance 2 dimensional array)<sup>13</sup>
      - $M : \Sigma_N \times \Sigma_T \rightarrow ((\Sigma_N \times \Sigma^*) + \text{error})$
      - $M[A, a] = w$
    - we assume: pure BNF

### Construction of the parsing table

1. Table recipe
  - If  $A \rightarrow \alpha \in P$  and  $\alpha \Rightarrow^* \mathbf{a}\beta$ , then add  $A \rightarrow \alpha$  to table entry  $M[A, \mathbf{a}]$
  - If  $A \rightarrow \alpha \in P$  and  $\alpha \Rightarrow^* \epsilon$  and  $S \$ \Rightarrow^* \beta A \mathbf{a} \gamma$  (where  $\mathbf{a}$  is a token (=non-terminal) or  $\$$ ), then add  $A \rightarrow \alpha$  to table entry  $M[A, \mathbf{a}]$
2. Table recipe (again, now using our old friends *First* and *Follow*) Assume  $A \rightarrow \alpha \in P$ .
  - If  $\mathbf{a} \in First(\alpha)$ , then add  $A \rightarrow \alpha$  to  $M[A, \mathbf{a}]$ .
  - If  $\alpha$  is nullable and  $\mathbf{a} \in Follow(A)$ , then add  $A \rightarrow \alpha$  to  $M[A, \mathbf{a}]$ .

<sup>13</sup>Often, the entry in the parse table does not contain a full rule as here, needed is only the *right-hand-side*. In that case the table is of type  $\Sigma_N \times \Sigma_T \rightarrow (\Sigma^* + \text{error})$ . We follow the convention of this book.

### Example: if-statements

- grammars is left-factored and not left recursive

$$\begin{array}{lcl}
 \text{stmt} & \rightarrow & \text{if-stmt} \mid \text{other} \\
 \text{if-stmt} & \rightarrow & \text{if} (\text{exp}) \text{stmt} \text{else-part} \\
 \text{else-part} & \rightarrow & \text{else stmt} \mid \epsilon \\
 \text{exp} & \rightarrow & 0 \mid 1
 \end{array}$$

	<i>First</i>	<i>Follow</i>
<i>stmt</i>	<b>other, if</b>	\$, else
<i>if-stmt</i>	<b>if</b>	\$, else
<i>else-part</i>	<b>else, <math>\epsilon</math></b>	\$, else
<i>exp</i>	<b>0, 1</b>	)

### Example: if statement: “LL(1) parse table”

$M[N, T]$	<b>if</b>	<b>other</b>	<b>else</b>	0	1	\$
<i>statement</i>	<i>statement</i> $\rightarrow \text{if-stmt}$	<i>statement</i> $\rightarrow \text{other}$				
<i>if-stmt</i>	<i>if-stmt</i> $\rightarrow$ <b>if</b> ( <i>exp</i> ) <i>statement</i> <i>else-part</i>					
<i>else-part</i>			<b>else-part</b> $\rightarrow$ <b>else</b> <i>statement</i> <b>else-part</b> $\rightarrow \epsilon$			<i>else-part</i> $\rightarrow \epsilon$
<i>exp</i>				<i>exp</i> $\rightarrow 0$	<i>exp</i> $\rightarrow 1$	

- 2 productions in the “red table entry”
- thus: it’s technically *not* an LL(1) table (and it’s not an LL(1) grammar)
- note: removing left-recursion and left-factoring did not help!

### LL(1) table based algo

```

while the top of the parsing stack ≠ $
  if the top of the parsing stack is terminal a
    and the next input token = a
  then
    pop the parsing stack;
    advance the input; // ‘‘match’’
  else if the top of the parsing stack is non-terminal A
    and the next input token is a terminal or $
    and parsing table  $M[A, a]$  contains
      production  $A \rightarrow X_1 X_2 \dots X_n$ 
    then (* generate *)
      pop the parsing stack
      for i := n to 1 do
        push  $X_i$  onto the stack;
  else error
if the top of the stack = $
  
```

```

then accept
end

```

### LL(1): illustration of run of the algo

Table 4.3 LL(1) parsing actions for if-statements using the most closely nested disambiguating rule	Parsing stack	Input	Action
\$ S	i(0)i(1)o e o \$	$S \rightarrow I$	
\$ I	i(0)i(1)o e o \$	$I \rightarrow i(E) S L$	
\$ L S ) E ( i	i(0)i(1)o e o \$	match	
\$ L S ) E (	(0)i(1)o e o \$	match	
\$ L S ) E	0)i(1)o e o \$	$E \rightarrow 0$	
\$ L S ) 0	0)i(1)o e o \$	match	
\$ L S )	)i(1)o e o \$	match	
\$ L S	i(1)o e o \$	$S \rightarrow I$	
\$ L I	i(1)o e o \$	$I \rightarrow i(E) S L$	
\$ L L S ) E ( i	i(1)o e o \$	match	
\$ L L S ) E (	(1)o e o \$	match	
\$ L L S ) E	1)o e o \$	$E \rightarrow 1$	
\$ L L S ) 1	1)o e o \$	match	
\$ L L S )	)o e o \$	match	
\$ L L S	o e o \$	$S \rightarrow o$	
\$ L L o	o e o \$	match	
\$ L L	e o \$	$L \rightarrow e S$	
\$ L S e	e o \$	match	
\$ L S	o \$	$S \rightarrow o$	
\$ L o	o \$	match	
\$ L	\$	$L \rightarrow \varepsilon$	
\$	\$	accept	

\*

**Remark** The most interesting steps are of course those dealing with the dangling else, namely those with the non-terminal *else-part* at the top of the stack. That's where the LL(1) table is ambiguous. In principle, with *else-part* on top of the stack (in the picture it's just L), the parser table allows always to make the decision that the "current statement" resp "current conditional" is done.

### Expressions Original grammar

$$\begin{aligned}
 \text{exp} &\rightarrow \text{exp addop term} \mid \text{term} \\
 \text{addop} &\rightarrow + \mid - \\
 \text{term} &\rightarrow \text{term mulop factor} \mid \text{factor} \\
 \text{mulop} &\rightarrow * \\
 \text{factor} &\rightarrow (\text{exp}) \mid \text{number}
 \end{aligned}$$

left-recursive  $\Rightarrow$  not LL(k)

### Left-rec removed

$$\begin{aligned}
 \text{exp} &\rightarrow \text{term exp}' \\
 \text{exp}' &\rightarrow \text{addop term exp}' \mid \epsilon \\
 \text{addop} &\rightarrow + \mid - \\
 \text{term} &\rightarrow \text{factor term}' \\
 \text{term}' &\rightarrow \text{mulop factor term}' \mid \epsilon \\
 \text{mulop} &\rightarrow * \\
 \text{factor} &\rightarrow (\text{exp}) \mid \text{n}
 \end{aligned}$$

	First	Follow
exp	(, number	\$, )
exp'	+, -, ε	\$, )
addop	+, -	(, number
term	(, number	\$, ), +, -
term'	* , ε	\$, ), +, -
mulop	*	(, number
factor	(, number	\$, ), +, -, *

### Expressions: LL(1) parse table

$M[N, T]$	(	<b>number</b>	)	+	-	*	\$
$exp$	$exp \rightarrow term\ exp'$	$exp \rightarrow term\ exp'$					
$exp'$			$exp' \rightarrow \epsilon$	$exp' \rightarrow addop$ $term\ exp'$	$exp' \rightarrow addop$ $term\ exp'$		$exp' \rightarrow \epsilon$
$addop$				$addop \rightarrow +$	$addop \rightarrow -$		
$term$	$term \rightarrow factor$ $term'$	$term \rightarrow factor$ $term'$					
$term'$			$term' \rightarrow \epsilon$	$term' \rightarrow \epsilon$	$term' \rightarrow \epsilon$	$term' \rightarrow mulop$ $factor$ $term'$	$term' \rightarrow \epsilon$
$mulop$			.			$mulop \rightarrow *$	
$factor$	$factor \rightarrow ( exp )$	$factor \rightarrow number$					

### Error handling

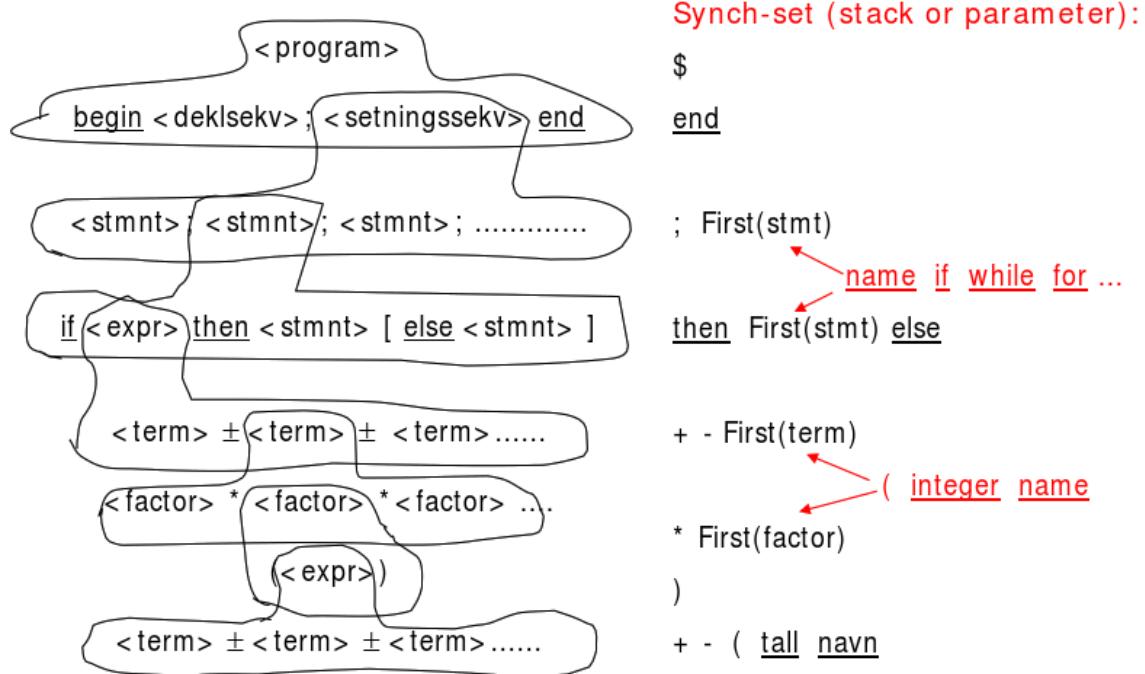
- at the least: do an understandable error message
- give indication of line / character or region responsible for the error in the source file
- potentially *stop* the parsing
- some compilers do *error recovery*
  - give an understandable error message (as minimum)
  - continue reading, until it's plausible to resume parsing  $\Rightarrow$  find more errors
  - however: when finding at least 1 error: no code generation
  - observation: resuming after syntax error is not easy

### Error messages

- important:
  - try to avoid error messages that only occur because of an already reported error!
  - report error as early as possible, if possible at the first point where the program cannot be extended to a correct program.
  - make sure that, after an error, one doesn't end up in a infinite loop without reading any input symbols.
- What's a good error message?
  - assume: that the method `factor()` chooses the alternative (`exp`) but that it, when control returns from method `exp()`, does not find a `)`
  - one could report : `left parenthesis missing`
  - But this may often be confusing, e.g. if what the program text is: `( a + b c )`
  - here the `exp()` method will terminate after `( a + b`, as `c` cannot extend the expression). You should therefore rather give the message `error in expression or left parenthesis missing`.

## Handling of syntax errors using recursive descent

### Method: «Panic mode» with use of «Synchronizing set»



## Syntax errors with sync stack

From the sketch at the previous page we can easily find:

- Which call should continue the execution?
- What input symbol should this method search for before resuming?
- We assume that \$ is added to the synch. stack only by the outermost method (for the start symbol)
- The union of everything on the stack is called the "synch. set", SS

The algorithm for this goes is as follows:

For each coming input symbol, test if it is a member of SS

If so:

- Look through the SS stack from newest to oldest, and find the newest method
  - that are willing to resume at one of these symbol
- This method will itself know how to resume after the actual input symbol

What is *not* easy is to program this without destroying the nice program structure occurring from pure recursive descent.

2

## Procedures for expression with "error recovery"

```

procedure exp ( synchset ) ;
begin
  checkinput ( { (, number }, synchset ) ;
  if not ( token in synchset ) then
    term ( synchset ) ;
    while token = + or token = - do
      match ( token ) ;
      term ( synchset ) ;
    end while ;
    checkinput ( synchset, { (, number } ) ;
  end if ;
end exp ;

```

```

procedure factor ( synchset ) ;
begin
  checkinput ( { (, number }, synchset ) ;
  if not ( token in synchset ) then
    case token of
      ( : match( ( ) ;
      exp ( { } ) ) ; ← Why not the full "synchset"?
      match( ) ;
      number :
        match(number) ;
      else error ;
    end case ;
    checkinput ( synchset, { (, number } ) ;
  end if ;
end factor ;

```

### Main philosophy

The method "checkinput" is called twice: First to check that the construction starts correctly, and secondly to check that the symbol after the construction is legal.

if token in { (, number } then ...

### Uses parameters, not a stack

The procedures must themselves resume execution at the right place inside themselves when they get the control back, or it must terminate immediately if it cannot resume execution on the current symbol.

```
procedure scanto ( synchset ) ;
```

```
begin
  while not ( token in synchset ∪ { $ } ) do
    getToken ;
  end scanto ;
```

```
procedure checkinput ( firstset, followset ) ;
```

```
begin
  if not ( token in firstset ) then
    error ;
    scanto ( firstset ∪ followset ) ;
  end if ;
end;
```

27

### 1.1.3 Bottom-up parsing

#### Bottom-up parsing: intro

"R" stands for *right-most* derivation.

**LR(0)**

- only for very simple grammars

- approx. 300 states for standard programming languages
- only as intro to SLR(1) and LALR(1)

**SLR(1)**

- expressive enough for most grammars for standard PLs

- same number of states as LR(0)
- main focus here

**LALR(1)**

- slightly more expressive than SLR(1)

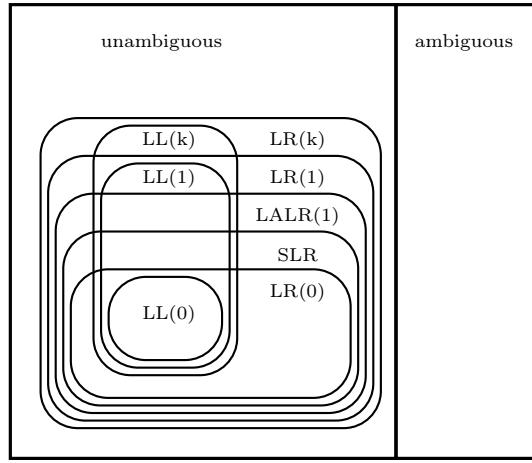
- same number of states as LR(0)
- we look at ideas behind that method as well

**LR(1)** covers all grammars, which can in principle be parsed by looking at the next token

#### 1. Remarks

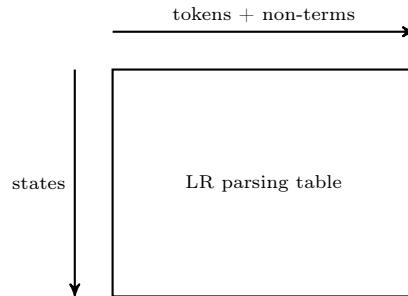
There seems to be a contradiction in the explanation of LR(0): if LR(0) is so weak that it works only for unreasonably simple language, how can one speak about that standard languages have 300 states? The answer is, the other more expressive parsers (SLR(1) and LALR(1)) use the *same* construction of states, so that's why one can estimate the number of states, even if standard languages don't have a LR(0) parser; they may have an LALR(1)-parser, which has, it its core, LR(0)-states.

#### Grammar classes overview (again)



### LR-parsing and its subclasses

- *right-most* derivation (but left-to-right parsing)
- in general: bottom-up parsing more powerful than top-down
- typically: tool-supported (unlike recursive descent, which may well be hand-coded)
- based on *parsing tables* + explicit *stack*
- thankfully: *left-recursion* no longer problematic
- typical tools: yacc and its descendants (like bison, CUP, etc)
- another name: *shift-reduce* parser



### Example grammar

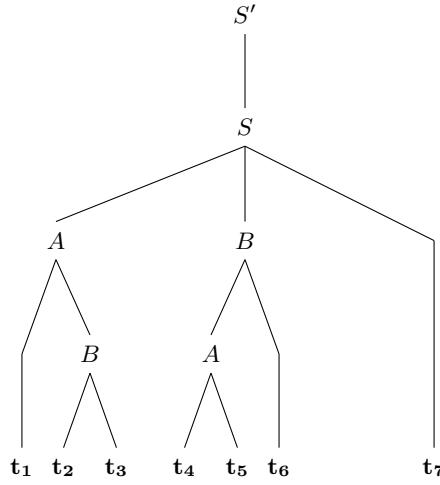
$$\begin{aligned}
 S' &\rightarrow S \\
 S &\rightarrow ABt_7 \mid \dots \\
 A &\rightarrow t_4t_5 \mid t_1B \mid \dots \\
 B &\rightarrow t_2t_3 \mid At_6 \mid \dots
 \end{aligned}$$

- assume: grammar unambiguous
- assume word of terminals  $t_1t_2\dots t_7$  and its (unique) parse-tree
- general agreement for bottom-up parsing:
  - start symbol *never* on the right-hand side or a production
  - **routinely add another “extra” start-symbol** (here  $S'$ )<sup>14</sup>

### Parse tree for $t_1\dots t_7$

---

<sup>14</sup>That will later be relied upon when constructing a DFA for “scanning” the stack, to control the reactions of the stack machine. This restriction leads to a unique, well-defined initial state.



Remember: parse tree independent from left- or right-most-derivation

### LR: left-to right scan, right-most derivation?

- Potentially puzzling question at first sight:

How does the parser *right*-most derivation, when parsing *left*-to-right?

- Discussion

- short answer: parser builds the parse tree **bottom-up**
- derivation:
  - replacement of nonterminals by right-hand sides
  - derivation*: builds (implicitly) a parse-tree *top-down*

- sentential form: word from  $\Sigma^*$  derivable from start-symbol

- Right-sentential form: right-most derivation

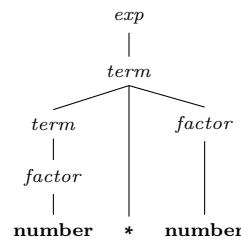
$$S \Rightarrow_r^* \alpha$$

- Slightly longer answer

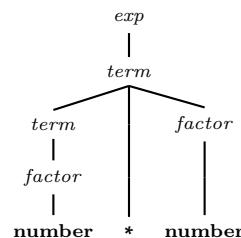
LR parser parses from left-to-right and builds the parse tree bottom-up. When doing the parse, the parser (implicitly) builds a *right-most* derivation **in reverse** (because of bottom-up).

### Example expression grammar (from before)

$$\begin{array}{lcl}
 \textit{exp} & \rightarrow & \textit{exp} \textit{ addop } \textit{ term } \mid \textit{ term } \\
 \textit{addop} & \rightarrow & + \mid - \\
 \textit{term} & \rightarrow & \textit{term} \textit{ mulop } \textit{ factor } \mid \textit{ factor } \\
 \textit{mulop} & \rightarrow & * \\
 \textit{factor} & \rightarrow & (\textit{exp}) \mid \textit{ number }
 \end{array} \tag{8}$$



### Bottom-up parse: Growing the parse tree



$$\begin{array}{lcl}
 \underline{\text{number}} * \text{number} & \hookrightarrow & \underline{\text{factor}} * \text{number} \\
 & \hookrightarrow & \underline{\text{term}} * \underline{\text{number}} \\
 & \hookrightarrow & \underline{\text{term}} * \underline{\text{factor}} \\
 & \hookrightarrow & \underline{\text{term}} \\
 & \hookrightarrow & \text{exp}
 \end{array}$$

**Reduction in reverse = right derivation**

1. Reduction

$$\begin{array}{lcl}
 \underline{\text{n}} * \text{n} & \hookrightarrow & \underline{\text{factor}} * \text{n} \\
 & \hookrightarrow & \underline{\text{term}} * \underline{\text{n}} \\
 & \hookrightarrow & \underline{\text{term}} * \underline{\text{factor}} \\
 & \hookrightarrow & \underline{\text{term}} \\
 & \hookrightarrow & \text{exp}
 \end{array}$$

2. Right derivation

$$\begin{array}{lcl}
 \text{n} * \text{n} & \Leftarrow_r & \underline{\text{factor}} * \text{n} \\
 \Leftarrow_r & & \underline{\text{term}} * \underline{\text{n}} \\
 \Leftarrow_r & & \underline{\text{term}} * \underline{\text{factor}} \\
 \Leftarrow_r & & \underline{\text{term}} \\
 \Leftarrow_r & & \underline{\text{exp}}
 \end{array}$$

3. Underlined entity

- underlined part:
  - *different* in reduction vs. derivation
  - represents the “part being replaced”
    - \* for derivation: right-most non-terminal
    - \* for reduction: indicates the so-called **handle** (or part of it)
- consequently: all intermediate words are *right-sentential forms*

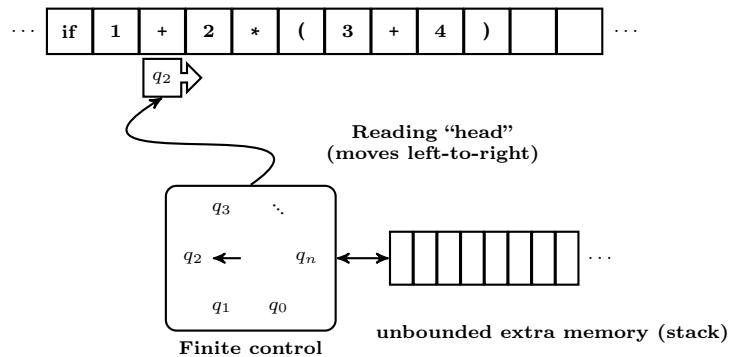
## Handle

**Definition 15** (Handle). Assume  $S \Rightarrow_r^* \alpha \textcolor{red}{A} w \Rightarrow_r \alpha \beta w$ . A production  $A \rightarrow \beta$  at position  $k$  following  $\alpha$  is a *handle* of  $\alpha \beta w$ . We write  $\langle A \rightarrow \beta, k \rangle$  for such a handle.

Note:

- $w$  (right of a handle) contains only terminals
- $w$ : corresponds to the future input still to be parsed!
- $\alpha \beta$  will correspond to the stack content ( $\beta$  the part touched by reduction step).
- the  $\Rightarrow_r$ -derivation-step *in reverse*:
  - one **reduce**-step in the LR-parser-machine
  - adding (implicitly in the LR-machine) a new parent to children  $\beta$  (= **bottom-up**!)
- “handle”-part  $\beta$  can be *empty* ( $= \epsilon$ )

## Schematic picture of parser machine (again)



## General LR “parser machine” configuration

- *Stack*:
  - contains: terminals + non-terminals (+ \$)
  - containing: what has been read already but not yet “processed”
- *position* on the “tape” (= token stream)
  - represented here as word of terminals *not yet read*
  - end of “rest of token stream”: \$, as usual
- *state* of the machine
  - in the following schematic illustrations: *not yet* part of the discussion
  - *later*: part of the parser table, currently we explain *without* referring to the state of the parser-engine
  - currently we assume: tree and rest of the input given
  - the trick ultimately will be: how do achieve the same *without that tree already given* (just parsing left-to-right)

## Schematic run (reduction: from top to bottom)

\$	t <sub>1</sub> t <sub>2</sub> t <sub>3</sub> t <sub>4</sub> t <sub>5</sub> t <sub>6</sub> t <sub>7</sub> \$
\$t <sub>1</sub>	t <sub>2</sub> t <sub>3</sub> t <sub>4</sub> t <sub>5</sub> t <sub>6</sub> t <sub>7</sub> \$
\$t <sub>1</sub> t <sub>2</sub>	t <sub>3</sub> t <sub>4</sub> t <sub>5</sub> t <sub>6</sub> t <sub>7</sub> \$
\$t <sub>1</sub> t <sub>2</sub> t <sub>3</sub>	t <sub>4</sub> t <sub>5</sub> t <sub>6</sub> t <sub>7</sub> \$
\$t <sub>1</sub> B	t <sub>4</sub> t <sub>5</sub> t <sub>6</sub> t <sub>7</sub> \$
\$A	t <sub>4</sub> t <sub>5</sub> t <sub>6</sub> t <sub>7</sub> \$
\$At <sub>4</sub>	t <sub>5</sub> t <sub>6</sub> t <sub>7</sub> \$
\$At <sub>4</sub> t <sub>5</sub>	t <sub>6</sub> t <sub>7</sub> \$
\$AA	t <sub>6</sub> t <sub>7</sub> \$
\$AAt <sub>6</sub>	t <sub>7</sub> \$
\$AB	t <sub>7</sub> \$
\$ABt <sub>7</sub>	\$
\$S	\$
\$S'	\$

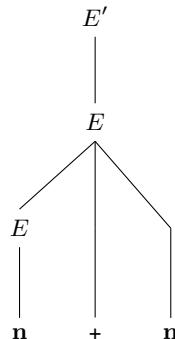
## 2 basic steps: shift and reduce

- parsers reads input and uses stack as intermediate storage
  - so far: no mention of look-ahead (i.e., action depending on the value of the next token(s)), but that may play a role, as well
1. Shift Move the next input symbol (terminal) over to the top of the stack (“push”)
  2. Reduce Remove the symbols of the *right-most* subtree from the stack and replace it by the non-terminal at the root of the subtree (replace = “pop + push”).
  3. Remarks
    - **easy to do if one has the parse tree already!**
    - *reduce* step: popped resp. pushed part = right- resp. left-hand side of handle

## Example: LR parsing for addition (given the tree)

$$\begin{array}{lcl} E' & \rightarrow & E \\ E & \rightarrow & E + n \mid n \end{array}$$

1. CST



2. Run

	parse stack	input	action
1	\$	n + n \$	shift
2	\$ n	+ n \$	red.: $E \rightarrow n$
3	\$ E	+ n \$	shift
4	\$ E +	n \$	shift
5	\$ E + n	\$	reduce $E \rightarrow E + n$
6	\$ E	\$	red.: $E' \rightarrow E$
7	\$ E'	\$	accept

*note:* line 3 vs line 6!; both contain  $E$  on top of stack

3. (right) derivation: reduce-steps “in reverse”

$$\underline{E'} \Rightarrow \underline{E} \Rightarrow \underline{E + n} \Rightarrow \underline{n + n}$$

**Example with  $\epsilon$ -transitions: parentheses**

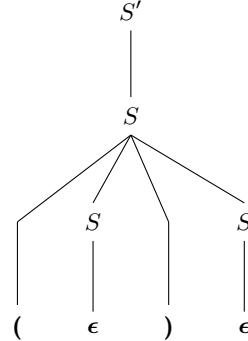
$$\begin{array}{lcl} S' & \rightarrow & S \\ S & \rightarrow & (S)S \mid \epsilon \end{array}$$

side remark: unlike previous grammar, here:

- production with *two* non-terminals in the right
- ⇒ difference between left-most and right-most derivations (and mixed ones)

**Parentheses: tree, run, and right-most derivation**

1. CST



2. Run

	parse stack	input	action
1	\$	( ) \$	shift
2	\$ (	) \$	reduce $S \rightarrow \epsilon$
3	\$ ( S )	) \$	shift
4	\$ ( S )	\$	reduce $S \rightarrow \epsilon$
5	\$ ( S ) S	\$	reduce $S \rightarrow (S)S$
6	\$ S	\$	reduce $S' \rightarrow S$
7	\$ S'	\$	accept

Note: the 2 reduction steps for the  $\epsilon$  productions

3. Right-most derivation and right-sentential forms

$$\underline{S'} \Rightarrow_r \underline{S} \Rightarrow_r ( \underline{S} ) \underline{S} \Rightarrow_r ( \underline{S} ) \Rightarrow_r ( )$$

**Right-sentential forms & the stack** - sentential form: word from  $\Sigma^*$  derivable from start-symbol

1. Right-sentential form: right-most derivation

$$S \Rightarrow_r^* \alpha$$

2. Explanation

- right-sentential forms:

3. Run
  - part of the “run”
  - but: **split** between *stack* and *input*

	parse stack	input	action
1	\$	n + n \$	shift
2	\$ n	+ n \$	red.: $E \rightarrow n$
3	\$ E	+ n \$	shift
4	\$ E +	n \$	shift
5	\$ E + n	\$	reduce $E \rightarrow E + n$
6	\$ E	\$	red.: $E' \rightarrow E$
7	\$ E'	\$	accept

4. Derivation and split

$$\underline{E'} \Rightarrow_r \underline{E} \Rightarrow_r \underline{E+n} \Rightarrow_r \underline{n+n}$$

$$\underline{n+n} \hookrightarrow \underline{E+n} \hookrightarrow \underline{E} \hookrightarrow E'$$

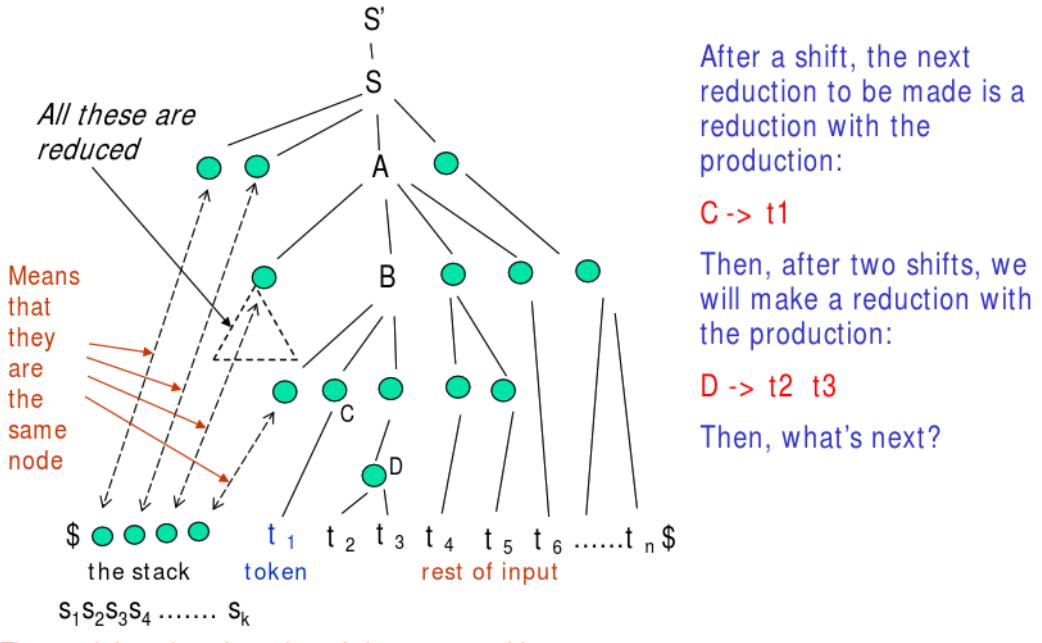
5. Rest

$$\underline{E'} \Rightarrow_r \underline{E} \Rightarrow_r \underline{E+n} \parallel \sim \underline{E+} \parallel \underline{n} \sim \underline{E} \parallel \underline{+n} \Rightarrow_r \underline{n} \parallel \underline{+n} \sim \parallel \underline{n+n}$$

### Viable prefixes of right-sentential forms and handles

- right-sentential form:  $E + n$
- **viable prefixes** of RSF
  - prefixes of that RSF *on the stack*
  - here: 3 viable prefixes of that RSF:  $E$ ,  $E+$ ,  $E+n$
- *handle*: remember the definition earlier
- here: for instance in the sentential form  $n+n$ 
  - handle is production  $E \rightarrow n$  on the *left* occurrence of  $n$  in  $n+n$  (let's write  $n_1 + n_2$  for now)
  - note: in the stack machine:
    - \* the left  $n_1$  on the stack
    - \* rest  $+n_2$  on the input (unread, because of LR(0))
- if the parser engine detects handle  $n_1$  on the stack, it does a *reduce-step*
- However (later): reaction depends on current *state* of the parser engine

### A typical situation during LR-parsing



### General design for an LR-engine

- some ingredients clarified up-to now:
  - bottom-up tree building as reverse right-most derivation,
  - stack vs. input,
  - shift and reduce steps
- however: 1 ingredient missing: next step of the engine may depend on
  - top of the stack (“handle”)
  - look ahead on the input (but not for LL(0))
  - and: current state of the machine (same stack-content, but different reactions at different stages of the parse)

### But what are the states of an LR-parser?

#### 1. General idea:

Construct an NFA (and ultimately DFA) which works on the **stack** (not the input). The alphabet consists of terminals and non-terminals  $\Sigma_T \cup \Sigma_N$ . The language

$$Stacks(G) = \{\alpha \mid \alpha \text{ may occur on the stack during LR-parsing of a sentence in } \mathcal{L}(G)\}$$

is **regular!**

### LR(0) parsing as easy pre-stage

- LR(0): in practice *too simple*, but easy conceptual step towards LR(1), SLR(1) etc.
- LR(1): in practice good enough, LR(k) not used for  $k > 1$
- 1. LR(0) item production with specific “parser position” . in its right-hand side
- 2. Rest
  - . is, of course, a “meta-symbol” (not part of the production)
  - For instance: production  $A \rightarrow \alpha$ , where  $\alpha = \beta\gamma$ , then
- 3. LR(0) item
 
$$A \rightarrow \beta.\gamma$$
- 4. complete and initial items
  - item with dot at the beginning: *initial* item
  - item with dot at the end: *complete* item

### Example: items of LR-grammar

1. Grammar for parentheses: 3 productions

$$\begin{array}{l} S' \rightarrow S \\ S \rightarrow (S)S \mid \epsilon \end{array}$$

2. 8 items

$$\begin{array}{l} S' \rightarrow .S \\ S' \rightarrow S. \\ S \rightarrow .(S)S \\ S \rightarrow (.S)S \\ S \rightarrow (S.)S \\ S \rightarrow (S).S \\ S \rightarrow (S)S. \\ S \rightarrow . \end{array}$$

3. Remarks

- note:  $S \rightarrow \epsilon$  gives  $S \rightarrow .$  as item (not  $S \rightarrow \epsilon.$  and  $S \rightarrow .\epsilon$ )
- side remark: see later, it will turn out: grammar *not*  $LR(0)$

### Another example: items for addition grammar

1. Grammar for addition: 3 productions

$$\begin{array}{l} E' \rightarrow E \\ E \rightarrow E + n \mid n \end{array}$$

2. (coincidentally also:) 8 items

$$\begin{array}{l} E' \rightarrow .E \\ E' \rightarrow E. \\ E \rightarrow .E + n \\ E \rightarrow E.+n \\ E \rightarrow E+.n \\ E \rightarrow E+n. \\ E \rightarrow .n \\ E \rightarrow n. \end{array}$$

3. Remarks: no  $LR(0)$

- also here: it will turn out: *not*  $LR(0)$  grammar

### Finite automata of items

- general set-up: *items* as **states in an automaton**
- automaton: “operates” *not* on the input, **but the stack**
- automaton either
  - first NFA, afterwards made deterministic (subset construction), or
  - directly DFA

1. States formed of sets of items In a state marked by/containing item

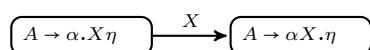
$$A \rightarrow \beta.\gamma$$

- $\beta$  on the *stack*
- $\gamma$ : to be treated next (terminals on the input, but can contain also non-terminals)

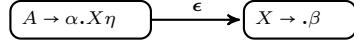
### State transitions of the NFA

- $X \in \Sigma$
- two kind of transitions

1. Terminal or non-terminal



2. Epsilon ( $X$ : non-terminal here)



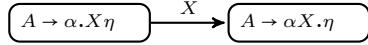
### 3. Explanation

- In case  $X = \text{terminal}$  (i.e. token) =
  - the left step corresponds to a **shift** step<sup>15</sup>
- for non-terminals (see next slide):
  - interpretation more complex: non-terminals are officially never on the input
  - note: in that case, item  $A \rightarrow \alpha.X\eta$  has two (kinds of) outgoing transitions

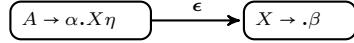
### Transitions for non-terminals and $\epsilon$

- so far: we never pushed a non-terminal from the input to the stack, we **replace** in a **reduce-step** the right-hand side by a left-hand side
- however: the replacement in a **reduce** steps can be seen as
  1. pop right-hand side off the stack,
  2. instead, “assume” corresponding non-terminal on input &
  3. eat the non-terminal and push it on the stack.
- two kind of transitions
  1. the  $\epsilon$ -transition correspond to the “pop” half
  2. that  $X$  transition (for non-terminals) corresponds to that “eat-and-push” part
- assume production  $X \rightarrow \beta$  and *initial* item  $X \rightarrow .\beta$

#### 1. Terminal or non-terminal



#### 2. Epsilon ( $X$ : non-terminal here) Given production $X \rightarrow \beta$ :



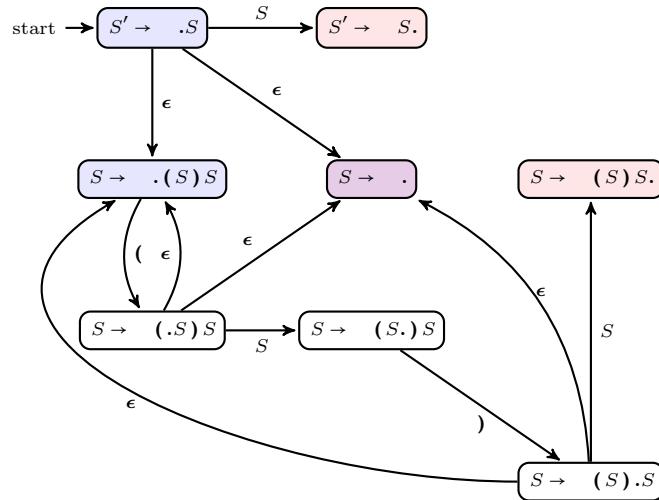
### Initial and final states

1. initial states:
  - we make our lives *easier*
  - we assume (as said): one *extra* start symbol say  $S'$  (augmented grammar)
  - $\Rightarrow$  initial item  $S' \rightarrow .S$  as (only) **initial state**
2. final states:
  - NFA has a specific task, “scanning” the stack, not scanning the input
  - acceptance condition of the *overall* machine: a bit more complex
    - input must be empty
    - stack must be empty except the (new) start symbol
    - NFA has a word to say about acceptance
      - \* but *not* in form of being in an accepting state
      - \* so: no accepting *states*
      - \* but: accepting *action* (see later)

---

<sup>15</sup>We have explained *shift* steps so far as: parser eats one *terminal* (= input token) and pushes it on the stack.

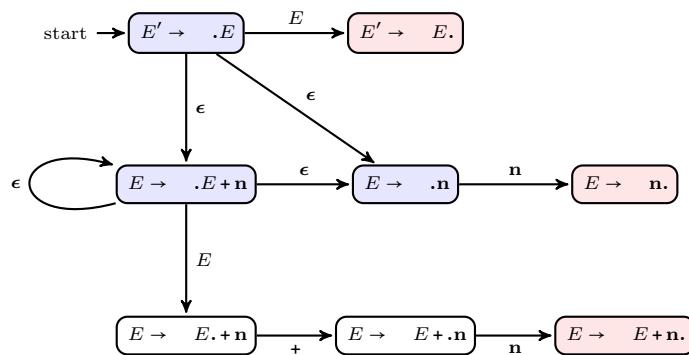
## NFA: parentheses



Remarks on the NFA

- colors for illustration
    - “reddish”: complete items
    - “blueish”: init-item (less important)
    - “violet’tish”: both
  - init-items
    - one per production of the grammar
    - that’s where the  $\epsilon$ -transistions go into, but
    - *with exception* of the initial state (with  $S'$ -production)  
no outgoing edges from the complete items

## NFA: addition

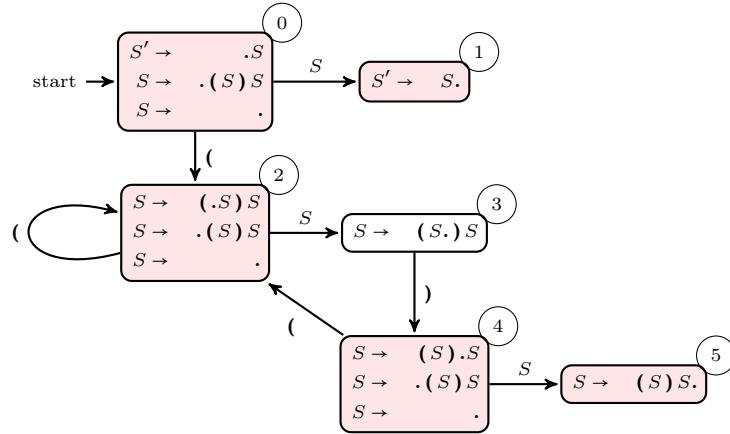


## Determinizing: from NFA to DFA

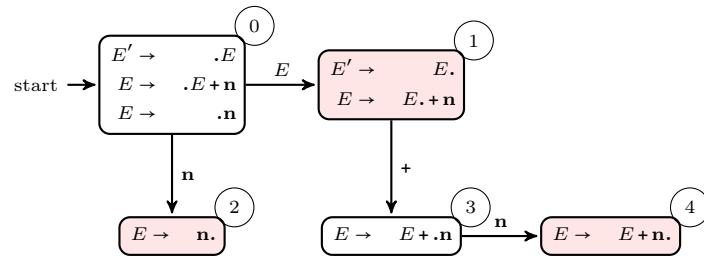
- standard subset-construction<sup>16</sup>
  - states then contains *sets* of items
  - especially important:  $\epsilon$ -closure
  - also: *direct* construction of the DFA possible

<sup>16</sup> Technically, we don't require here a *total* transition function, we leave out any error state.

## DFA: parentheses



## DFA: addition



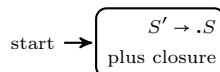
## Direct construction of an LR(0)-DFA

- quite easy: simply build in the closure already

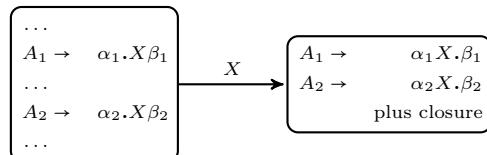
### 1. $\epsilon$ -closure

- if  $A \rightarrow \alpha.B\gamma$  is an item in a state where
- there are productions  $B \rightarrow \beta_1 \mid \beta_2 \dots \Rightarrow$
- add items  $B \rightarrow .\beta_1, B \rightarrow .\beta_2 \dots$  to the state
- continue that process, until saturation

### 2. initial state



## Direct DFA construction: transitions



- $X$ : terminal or non-terminal, both treated uniformly
- All items of the form  $A \rightarrow \alpha.X\beta$  must be included in the post-state
- and all others (indicated by "...") in the pre-state: not included
- re-check the previous examples: outcome is the same

## How does the DFA do the shift/reduce and the rest?

- we have seen: bottom-up parse tree generation
  - we have seen: shift-reduce and the stack vs. input
  - we have seen: the construction of the DFA
1. But: how does it hang together? We need to interpret the “set-of-item-states” in the light of the stack content and figure out the **reaction** in terms of
- transitions in the automaton
  - stack manipulations (shift/reduce)
  - acceptance
  - input (apart from shifting) not relevant when doing LR(0)
2. Determinism and the reaction better be uniquely determined ....

## Stack contents and state of the automaton

- remember: at any given intermediate configuration of stack/input in a run
  1. stack contains words from  $\Sigma^*$
  2. DFA operates deterministically on such words
- the stack contains the “past”: read input (potentially partially reduced)
- when feeding that “past” on the stack into the automaton
  - starting with the oldest symbol (not in a LIFO manner)
  - starting with the DFA’s initial state

$\Rightarrow$  stack content **determines** state of the DFA
- actually: each prefix also determines uniquely a state
- **top state**:
  - state after the complete stack content
  - corresponds to the **current** state of the stack-machine

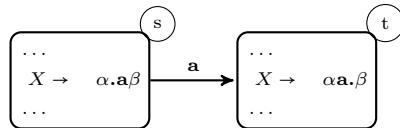
$\Rightarrow$  crucial when determining *reaction*

## State transition allowing a shift

- assume: top-state (= current state) contains item

$$X \rightarrow \alpha.a\beta$$

- construction thus has transition as follows



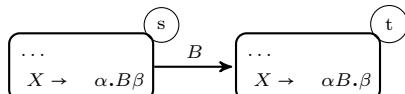
- shift is possible
- if shift is *the* correct operation and **a** is terminal symbol corresponding to the current token: state afterwards = *t*

## State transition: analogous for non-terminals

1. Production

$$X \rightarrow \alpha.B\beta$$

2. Transition

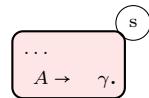


3. Explanation

- same as before, now with non-terminal  $B$
- note: we never read non-term from input
- not officially called a shift
- corresponds to the reaction **followed by** a *reduce* step, it's **not** the reduce step itself
- think of it as follows: reduce and subsequent step
  - not as: *replace* on top of the stack the handle (right-hand side) by non-term  $B$ ,
  - but instead as:
    - (a) pop off the handle from the top of the stack
    - (b) put the non-term  $B$  "back onto the input" (corresponding to the above state  $s$ )
    - (c) eat the  $B$  and *shift* it to the stack
- later: a **goto** reaction in the parse table

### State (not transition) where a reduce is possible

- remember: *complete items* (those with a dot . at the end)
- assume **top state**  $s$  containing complete item  $A \rightarrow \gamma.$



- a complete right-hand side ("handle")  $\gamma$  on the stack and thus done
  - may be replaced by right-hand side  $A$
- $\Rightarrow$  reduce step
- builds up (implicitly) new parent node  $A$  in the bottom-up procedure
  - **Note:**  $A$  on top of the stack instead of  $\gamma$ :<sup>17</sup>
    - **new top state!**
    - remember the "goto-transition" (shift of a non-terminal)

### Remarks: states, transitions, and reduce steps

- ignoring the  $\epsilon$ -transitions (for the NFA)
- there are 2 "kinds" of transitions in the DFA
  1. terminals: real shifts
  2. non-terminals: "following a reduce step"
- 1. No edges to represent (all of) a reduce step!
  - if a reduce happens, parser engine *changes state!*
  - however: this state change is **not** represented by a transition in the DFA (or NFA for that matter)
  - especially *not* by outgoing errors of completed items
- 2. Rest
  - if the (rhs of the) handle is *removed* from top stack:  $\Rightarrow$ 
    - "go back to the (top) state before that handle had been added": *no edge for that*
  - later: stack notation simply remembers the state as part of its configuration

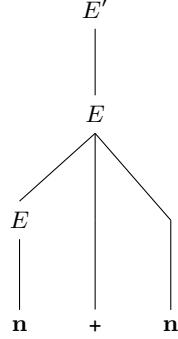
---

<sup>17</sup>Indirectly only: as said, we remove the handle from the stack, and pretend, as if the  $A$  is next on the input, and thus we "shift" it on top of the stack, doing the corresponding  $A$ -transition.

**Example: LR parsing for addition (given the tree)**

$$\begin{array}{lcl} E' & \rightarrow & E \\ E & \rightarrow & E + n \mid n \end{array}$$

1. CST

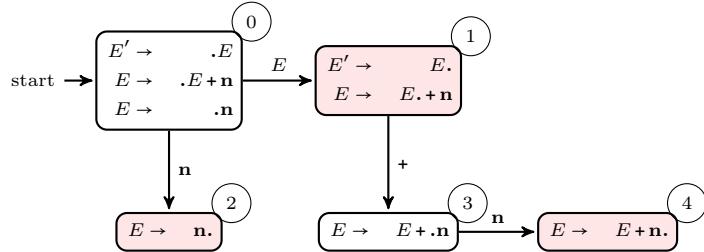


2. Run

	parse stack	input	action
1	\$	n + n \$	shift
2	\$ n	+ n \$	red.: E → n
3	\$ E	+ n \$	shift
4	\$ E +	n \$	shift
5	\$ E + n	\$	reduce E → E + n
6	\$ E	\$	red.: E' → E
7	\$ E'	\$	accept

*note:* line 3 vs line 6!; both contain E on top of stack

**DFA of addition example**



- note line 3 vs. line 6
- both stacks = E ⇒ same (top) state in the DFA (state 1)

**LR(0) grammars**

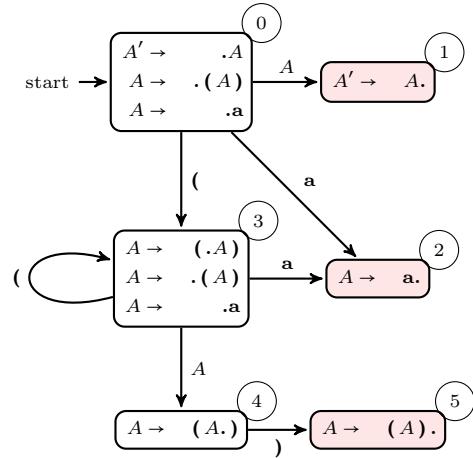
1. LR(0) grammar The top-state alone determines the next step.
2. No LR(0) here

- especially: no shift/reduce conflicts in the form shown
- thus: previous number-grammar is *not* LR(0)

**Simple parentheses**

$$A \rightarrow (A) \mid a$$

1. DFA

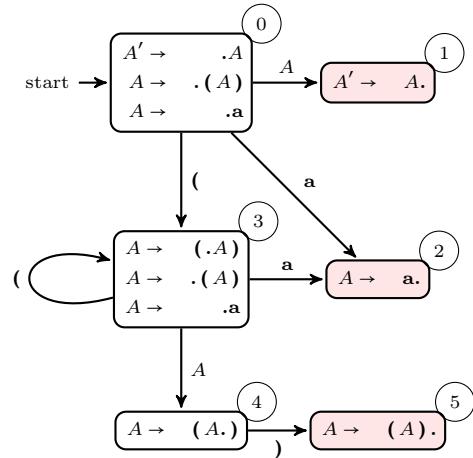


## 2. Remarks

- for *shift*:
  - **many** shift transitions in 1 state allowed
  - shift counts as *one* action (including “shifts” on non-terms)
- but for reduction: also the *production* must be clear

**Simple parentheses is LR(0)**

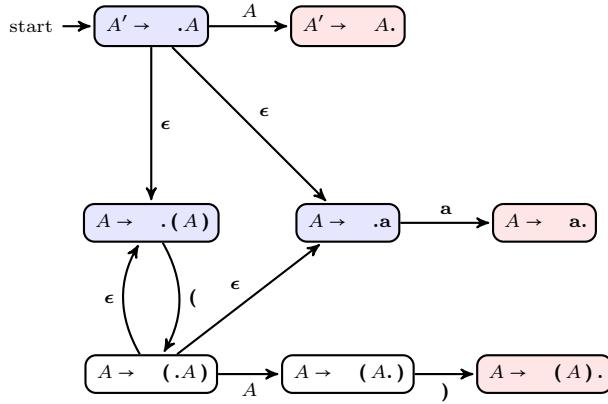
### 1. DFA



## 2. Remarks

state	possible action
0	only shift
1	only red: (with $A' \rightarrow A$ )
2	only red: (with $A \rightarrow a$ )
3	only shift
4	only shift
5	only red (with $A \rightarrow (A)$ )

**NFA for simple parentheses (bonus slide)**



### Parsing table for an LR(0) grammar

- table structure: slightly different for SLR(1), LALR(1), and LR(1) (see later)
- note: the “goto” part: “shift” on non-terminals (only 1 non-terminal  $A$  here)
- corresponding to the  $A$ -labelled transitions
- see the parser run on the next slide

state	action	rule	input	goto
			( a )	$A$
0	shift		3	1
1	reduce	$A' \rightarrow A$		
2	reduce	$A \rightarrow a$		
3	shift		3	2
4	shift			5
5	reduce	$A \rightarrow (A)$		

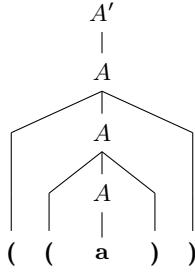
### Parsing of ((a))

stage	parsing stack	input	action
1	$\$_0$	((a))\$	shift
2	$\$_0(3$	(a))\$	shift
3	$\$_0(3(3$	a))\$	shift
4	$\$_0(3(3a_2$	)\$	reduce $A \rightarrow a$
5	$\$_0(3(3A_4$	)\$	shift
6	$\$_0(3(3A_4)_5$	)\$	reduce $A \rightarrow (A)$
7	$\$_0(3A_4$	)\$	shift
8	$\$_0(3A_4)_5$	\$	reduce $A \rightarrow (A)$
9	$\$_0A_1$	\$	accept

- note: stack on the left
  - contains top *state* information
  - in particular: overall **top** state on the right-most end
- note also: **accept** action
  - reduce wrt. to  $A' \rightarrow A$  and
  - *empty stack* (apart from \$,  $A$ , and the state annotation)

⇒ accept

### Parse tree of the parse



- As said:
  - the reduction “contains” the parse-tree
  - reduction: builds it bottom up
  - reduction in reverse: contains a *right-most* derivation (which is “top-down”)
- accept action: corresponds to the parent-child edge  $A' \rightarrow A$  of the tree

### Parsing of erroneous input

- empty slots in the table: “errors”

stage	parsing stack	input	action
1	$\$_0$	$((\text{a})\$$	shift
2	$\$_0($	$(\text{a})\$$	shift
3	$\$_0($ <sub>3</sub> $)$	$\text{a})\$$	shift
4	$\$_0($ <sub>3</sub> $\text{a}_2$	$)\$$	reduce $A \rightarrow \text{a}$
5	$\$_0($ <sub>3</sub> $A_4$	$)\$$	shift
6	$\$_0($ <sub>3</sub> $A_4)_5$	$\$$	reduce $A \rightarrow (\text{A})$
7	$\$_0($ <sub>3</sub> $A_4$	$\$$	????

stage	parsing stack	input	action
1	$\$_0$	$()\$$	shift
2	$\$_0($ <sub>3</sub>	$)\$$	?????

1. Invariant important general invariant for LR-parsing: never shift something “illegal” onto the stack

### LR(0) parsing algo, given DFA

let  $s$  be the current state, on top of the parse stack

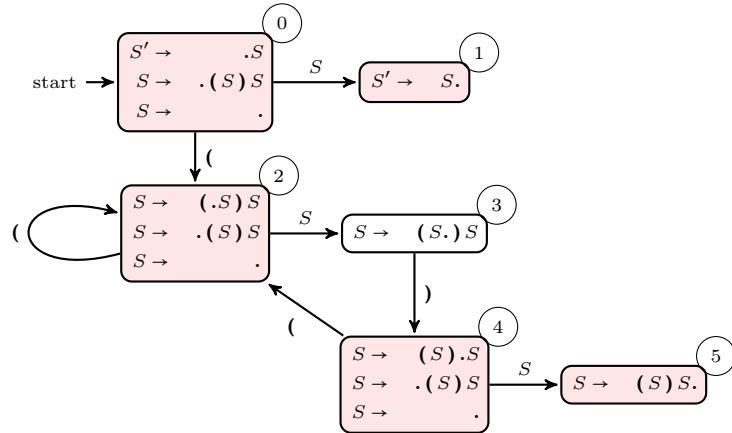
1.  $s$  contains  $A \rightarrow \alpha.X\beta$ , where  $X$  is a *terminal*
  - shift  $X$  from input to top of stack. the new *state* pushed on the stack: state  $t$  where  $s \xrightarrow{X} t$
  - else: if  $s$  does not have such a transition: *error*
2.  $s$  contains a **complete** item (say  $A \rightarrow \gamma.$ ): **reduce** by rule  $A \rightarrow \gamma$ :
  - A reduction by  $S' \rightarrow S$ : **accept**, if input is empty, **error**:
  - else:
    - pop**: remove  $\gamma$  (including “its” states from the stack)
    - back up**: assume to be in state  $u$  which is *now* head state
    - push**: push  $A$  to the stack, new head state  $t$  where  $u \xrightarrow{A} t$  (in the DFA)

### LR(0) parsing algo remarks

- in [Louden, 1997]: slightly differently formulated
- instead of requiring (in the first case):
  - push state  $t$  were  $s \xrightarrow{X} t$  or similar, book formulates
  - push *state containing item*  $A \rightarrow \alpha.X\beta$
- analogous in the second case
- algo (= deterministic) only if LR(0) grammar
  - in particular: cannot have states with *complete item* and item of form  $A\alpha.X\beta$  (otherwise **shift-reduce** conflict)
  - cannot have states with two  $X$ -successors (known as **reduce-reduce** conflict)

DFA parentheses again: LR(0)?

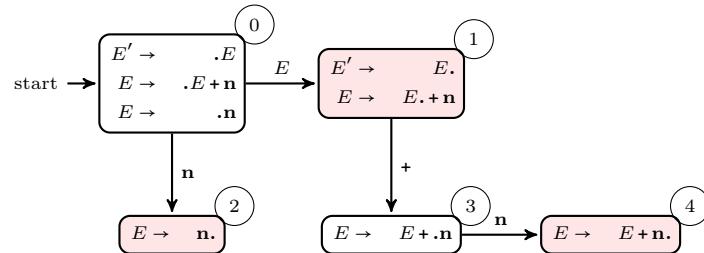
$$\begin{array}{l} S' \rightarrow S \\ S \rightarrow (S)S \mid \epsilon \end{array}$$



Look at states 0, 2, and 4

DFA addition again: LR(0)?

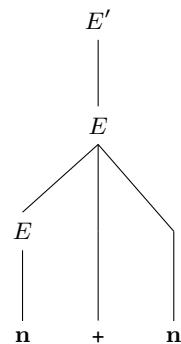
$$\begin{array}{l} E' \rightarrow E \\ E \rightarrow E + n \mid n \end{array}$$



How to make a decision in state 1?

Decision? If only we knew the ultimate tree already ... especially the parts still to come

1. CST



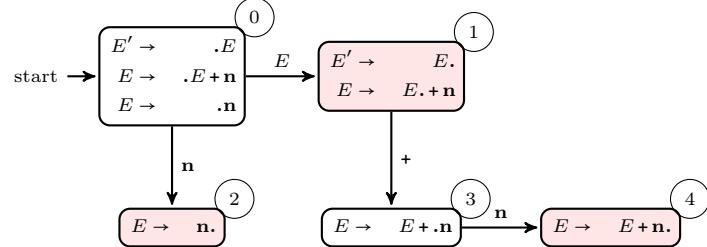
2. Run

	parse stack	input	action
1	\$	$n + n \$$	shift
2	\$ n	$+ n \$$	red.: $E \rightarrow n$
3	\$ E	$+ n \$$	shift
4	\$ E +	$n \$$	shift
5	\$ E + n	$\$$	reduce $E \rightarrow E + n$
6	\$ E	$\$$	red.: $E' \rightarrow E$
7	\$ E'	$\$$	accept

### 3. Explanation

- current stack: represents already known part of the parse tree
- since we don't have the future parts of the tree yet:  
⇒ **look-ahead** on the input (without building the tree as yet)
- LR(1) and its variants: *look-ahead of 1* (= look at the current type of the token)

### Addition grammar (again)



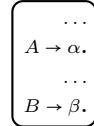
- How to make a decision in state 1? (here: shift vs. reduce)  
⇒ look at the next input symbol (in the token)

### One look-ahead

- LR(0), not useful, too weak
- add look-ahead, here of 1 input symbol (= token)
- different variations of that idea (with slight difference in expressiveness)
- tables slightly changed (compared to LR(0))
- but: still can use the LR(0)-DFAs

### Resolving LR(0) reduce/reduce conflicts

#### 1. LR(0) reduce/reduce conflict:

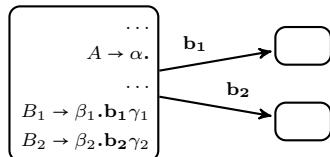


#### 2. SLR(1) solution: use follow sets of non-terms

- If  $Follow(A) \cap Follow(B) = \emptyset$   
⇒ next symbol (in token) decides!
  - if  $\text{token} \in Follow(\alpha)$  then reduce using  $A \rightarrow \alpha$
  - if  $\text{token} \in Follow(\beta)$  then reduce using  $B \rightarrow \beta$
  - ...

### Resolving LR(0) shift/reduce conflicts

#### 1. LR(0) shift/reduce conflict:



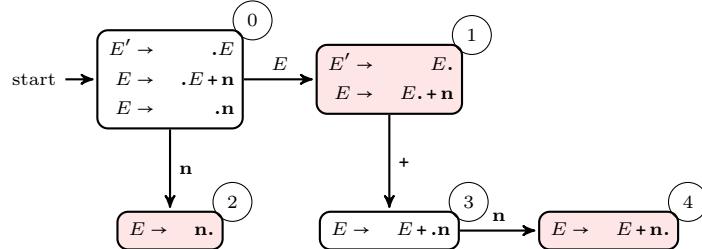
#### 2. SLR(1) solution: again: use follow sets of non-terms

- If  $Follow(A) \cap \{b_1, b_2, \dots\} = \emptyset$   
⇒ next symbol (in token) decides!
  - if  $\text{token} \in Follow(A)$  then reduce using  $A \rightarrow \alpha$ , non-terminal  $A$  determines new top state
  - if  $\text{token} \in \{b_1, b_2, \dots\}$  then shift. Input symbol  $b_i$  determines new top state
  - ...

### SLR(1) requirement on states (as in the book)

- formulated as conditions on the states (of LR(0)-items)
  - given the LR(0)-item DFA as defined
1. SLR(1) condition, on all states  $s$ 
    - (a) For any item  $A \rightarrow \alpha.X\beta$  in  $s$  with  $X$  a terminal, there is no **complete** item  $B \rightarrow \gamma.$  in  $s$  with  $X \in Follow(B).$
    - (b) For any **two complete** items  $A \rightarrow \alpha.$  and  $B \rightarrow \beta.$  in  $s$ ,  $Follow(\alpha) \cap Follow(\beta) = \emptyset$

Revisit addition one more time

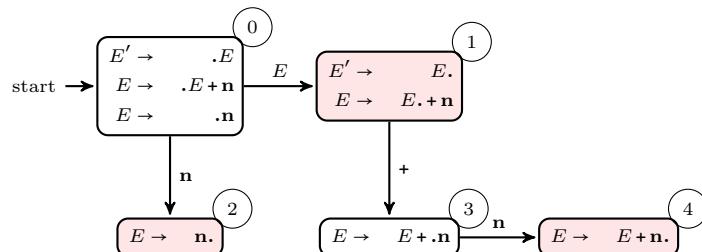


- $Follow(E') = \{\$\}$
- $\Rightarrow$
- shift for +
  - reduce with  $E' \rightarrow E$  for  $\$$  (which corresponds to accept, in case the input is empty)

**SLR(1) algo** let  $s$  be the current state, on top of the parse stack

1.  $s$  contains  $A \rightarrow \alpha.X\beta$ , where  $X$  is a terminal **and  $X$  is the next token on the input**, then
  - shift  $X$  from input to top of stack. the new state pushed on the stack: state  $t$  where  $s \xrightarrow{X} t$ <sup>18</sup>
2.  $s$  contains a *complete* item (say  $A \rightarrow \gamma.$ ) **and the next token in the input is in  $Follow(A)$ :** *reduce by rule  $A \rightarrow \gamma$ :*
  - A reduction by  $S' \rightarrow S$ : *accept*, if input is empty<sup>19</sup>
  - else:
    - pop:** remove  $\gamma$  (including “its” states from the stack)
    - back up:** assume to be in state  $u$  which is *now* head state
    - push:** push  $A$  to the stack, new head state  $t$  where  $u \xrightarrow{A} t$
3. if next token is such that neither 1. or 2. applies: *error*

Parsing table for SLR(1)



state	input			goto
	n	+	\$	
0	$s : 2$			1
1		$s : 3$	accept	
2		$r : (E \rightarrow n)$		
3	$s : 4$			
4		$r : (E \rightarrow E + n)$	$r : (E \rightarrow E + n)$	

for state 2 and 4:  $n \notin Follow(E)$

<sup>18</sup>Cf. to the LR(0) algo: since we checked the existence of the transition before, the else-part is missing now.

<sup>19</sup>Cf. to the LR(0) algo: This happens *now* only if next token is  $\$$ . Note that the follow set of  $S'$  in the *augmented* grammar is always only  $\$$

## Parsing table: remarks

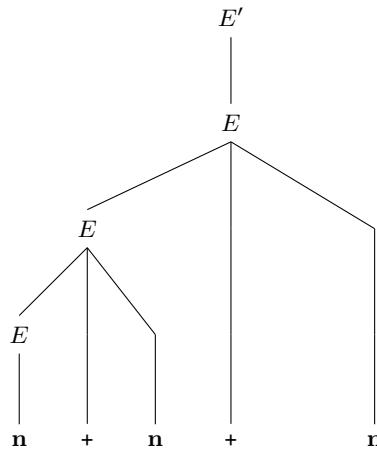
- SLR(1) parsing table: rather similar-looking to the LR(0) one
- differences: reflect the differences in: LR(0)-algo vs. SLR(1)-algo
- same number of rows in the table (= same number of states in the DFA)
- only: columns “arranged differently”
  - LR(0): each state **uniformely**: either shift or else reduce (with given rule)
  - now: non-uniform, **dependent** on the input. But that does not apply to the previous example. We'll see that in the next, then.
- it should be obvious:
  - SLR(1) may resolve LR(0) conflicts
  - but: if the follow-set conditions are not met: SLR(1) *shift-shift* and/or SRL(1) *shift-reduce* conflicts
  - would result in non-unique entries in SRL(1)-table<sup>20</sup>

## SLR(1) parser run (= “reduction”)

state	input			goto
	n	+	\$	
0	s : 2			1
1		s : 3		
2			accept	
3	r : (E → n)			
4	s : 4			
		r : (E → E + n)	r : (E → E + n)	

stage	parsing stack	input	action
1	\$ <sub>0</sub>	n + n + n \$	shift: 2
2	\$ <sub>0</sub> n <sub>2</sub>	+ n + n \$	reduce: E → n
3	\$ <sub>0</sub> E <sub>1</sub>	+ n + n \$	shift: 3
4	\$ <sub>0</sub> E <sub>1</sub> + <sub>3</sub>	n + n \$	shift: 4
5	\$ <sub>0</sub> E <sub>1</sub> + <sub>3</sub> n <sub>4</sub>	+ n \$	reduce: E → E + n
6	\$ <sub>0</sub> E <sub>1</sub>	n \$	shift 3
7	\$ <sub>0</sub> E <sub>1</sub> + <sub>3</sub>	n \$	shift 4
8	\$ <sub>0</sub> E <sub>1</sub> + <sub>3</sub> n <sub>4</sub>	\$	reduce: E → E + n
9	\$ <sub>0</sub> E <sub>1</sub>	\$	accept

## Corresponding parse tree



<sup>20</sup>by which it, strictly speaking, would no longer be an SRL(1)-table :-)

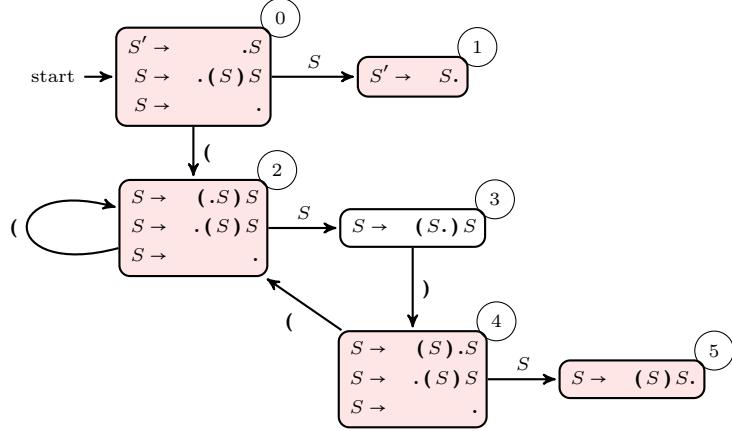
## Revisit the parentheses again: SLR(1)?

1. Grammar: parentheses (from before)

$$\begin{array}{l} S' \rightarrow S \\ S \rightarrow (S)S \mid \epsilon \end{array}$$

2. Follow set  $Follow(S) = \{\}, \$\}$

3. DFA



SLR(1) parse table

state	input			goto
	(	)	\$	S
0	$s : 2$	$r : S \rightarrow \epsilon$	$r : S \rightarrow \epsilon$	1
1			accept	
2	$s : 2$	$r : S \rightarrow \epsilon$	$r : S \rightarrow \epsilon$	3
3		$s : 4$		
4	$s : 2$	$r : S \rightarrow \epsilon$	$r : S \rightarrow \epsilon$	5
5		$r : S \rightarrow (S)S$	$r : S \rightarrow (S)S$	

Parentheses: SLR(1) parser run (= “reduction”)

state	input			goto
	(	)	\$	S
0	$s : 2$	$r : S \rightarrow \epsilon$	$r : S \rightarrow \epsilon$	1
1			accept	
2	$s : 2$	$r : S \rightarrow \epsilon$	$r : S \rightarrow \epsilon$	3
3		$s : 4$		
4	$s : 2$	$r : S \rightarrow \epsilon$	$r : S \rightarrow \epsilon$	5
5		$r : S \rightarrow (S)S$	$r : S \rightarrow (S)S$	

stage	parsing stack	input	action
1	$\$_0$	$()()\$$	shift: 2
2	$\$_0(2$	$)()\$$	reduce: $S \rightarrow \epsilon$
3	$\$_0(2S_3$	$)()\$$	shift: 4
4	$\$_0(2S_3)_4$	$)\$$	shift: 2
5	$\$_0(2S_3)_4(2$	$)\$$	reduce: $S \rightarrow \epsilon$
6	$\$_0(2S_3)_4(2S_3$	$)\$$	shift: 4
7	$\$_0(2S_3)_4(2S_3)_4$	$\$$	reduce: $S \rightarrow \epsilon$
8	$\$_0(2S_3)_4(2S_3)_4S_5$	$\$$	reduce: $S \rightarrow (S)S$
9	$\$_0(2S_3)_4S_5$	$\$$	reduce: $S \rightarrow (S)S$
10	$\$_0S_1$	$\$$	accept

1. Remarks Note how the stack grows, and would continue to grow if the sequence of ( ) would continue. That's characteristic from right-recursive formulation of rules, and may constitute a problem for LR-parsing (stack-overflow).

## SLR(k)

- in principle: straightforward:  $k$  look-ahead, instead of 1
- rarely used in practice, using  $\text{First}_k$  and  $\text{Follow}_k$  instead of the  $k = 1$  versions
- tables grow *exponentially* with  $k$ !

As with other parsing algorithms, the SLR(1) parsing algorithm can be extended to SLR( $k$ ) parsing where parsing actions are based on  $k \geq 1$  symbols of lookahead. Using the sets  $\text{First}_k$  and  $\text{Follow}_k$  as defined in the previous chapter, an SLR( $k$ ) parser uses the following two rules:

1. If state  $s$  contains an item of the form  $A \rightarrow \alpha . X \beta$  ( $X$  a token), and  $Xw \in \text{First}_k(X \beta)$  are the next  $k$  tokens in the input string, then the action is to shift the current input token onto the stack, and the new state to be pushed on the stack is the state containing the item  $A \rightarrow \alpha X . \beta$ .
2. If state  $s$  contains the complete item  $A \rightarrow \alpha .$ , and  $w \in \text{Follow}_k(A)$  are the next  $k$  tokens in the input string, then the action is to reduce by the rule  $A \rightarrow \alpha$ .

SLR( $k$ ) parsing is more powerful than SLR(1) parsing when  $k > 1$ , but at a substantial cost in complexity, since the parsing table grows exponentially in size with  $k$ .

## Ambiguity & LR-parsing

- in principle: LR( $k$ ) (and LL( $k$ )) grammars: *unambiguous*
- definition/construction: free of shift/reduce and reduce/reduce conflict (given the chosen level of lookahead)
- However: ambiguous grammar tolerable, if (remaining) conflicts can be solved “meaningfully” otherwise:
  1. Additional means of disambiguation:
    - (a) by specifying associativity / precedence “outside” the grammar
    - (b) by “living with the fact” that LR parser (commonly) *prioritizes shifts over reduces*
  2. Rest
    - for the second point (“let the parser decide according to its preferences”):
      - use sparingly and cautiously
      - typical example: *dangling-else*
      - even if parsers makes a decision, program may or may not “understand intuitively” the resulting parse tree (and thus AST)
      - grammar with many S/R-conflicts: go back to the drawing board

```

stmt   →  if-stmt | other
if-stmt →  if ( exp ) stmt
          |  if ( exp ) stmt else stmt
exp    →  0 | 1
  
```

**Example of an ambiguous grammar** In the following,  $E$  for  $exp$ , etc.

### Simplified conditionals

1. Simplified “schematic” if-then-else

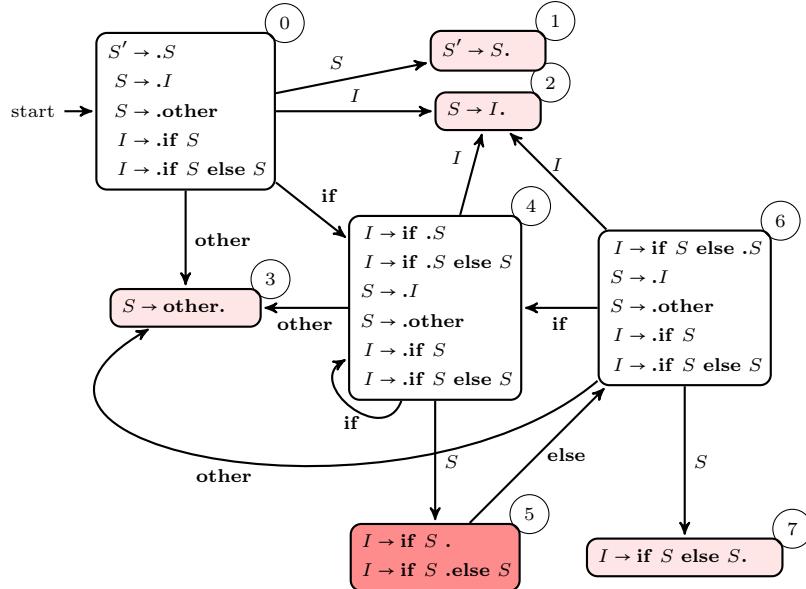
$$\begin{array}{lcl} S & \rightarrow & I \mid \text{other} \\ I & \rightarrow & \text{if } S \mid \text{if } S \text{ else } S \end{array}$$

2. Follow-sets

	Follow
$S'$	$\{\$\}$
$S$	$\{\$\}, \{\text{else}\}$
$I$	$\{\$\}, \{\text{else}\}$

3. Rest

- since ambiguous: at least one conflict must be somewhere



DFA of LR(0) items

### Simple conditionals: parse table

#### 1. Grammar

$$\begin{array}{lcl}
 S & \rightarrow & I \quad (1) \\
 & | & \text{other} \quad (2) \\
 I & \rightarrow & \text{if } S \quad (3) \\
 & | & \text{if } S \text{ else } S \quad (4)
 \end{array}$$

#### 2. SLR(1)-parse-table, conflict resolved

state	input				goto <i>S</i> <i>I</i>
	if	else	other	\$	
0	<i>s</i> : 4		<i>s</i> : 3		1    2
1					
2		<i>r</i> : 1		<i>r</i> : 1	
3		<i>r</i> : 2		<i>r</i> : 2	
4	<i>s</i> : 4		<i>s</i> : 3		5    2
5				<i>r</i> : 3	
6	<i>s</i> : 4		<i>s</i> : 3		7    2
7		<i>r</i> : 4		<i>r</i> : 4	

#### 3. Explanation

- *shift-reduce conflict* in state 5: reduce with rule 3 vs. shift (to state 6)
- conflict there: **resolved** in favor of *shift* to 6
- note: extra start state left out from the table

### Parser run (= reduction)

state	input				goto <i>S</i> <i>I</i>
	if	else	other	\$	
0	<i>s</i> : 4		<i>s</i> : 3		1    2
1					
2		<i>r</i> : 1		<i>r</i> : 1	
3		<i>r</i> : 2		<i>r</i> : 2	
4	<i>s</i> : 4		<i>s</i> : 3		5    2
5				<i>r</i> : 3	
6	<i>s</i> : 4		<i>s</i> : 3		7    2
7		<i>r</i> : 4		<i>r</i> : 4	

stage	parsing stack	input	action
1	\$ <sub>0</sub>	if if other else other \$	shift: 4
2	\$ <sub>0</sub> if <sub>4</sub>	if other else other \$	shift: 4
3	\$ <sub>0</sub> if <sub>4</sub> if <sub>4</sub>	other else other \$	shift: 3
4	\$ <sub>0</sub> if <sub>4</sub> if <sub>4</sub> other <sub>3</sub>	else other \$	reduce: 2
5	\$ <sub>0</sub> if <sub>4</sub> if <sub>4</sub> S <sub>5</sub>	else other \$	shift 6
6	\$ <sub>0</sub> if <sub>4</sub> if <sub>4</sub> S <sub>5</sub> else <sub>6</sub>	other \$	shift: 3
7	\$ <sub>0</sub> if <sub>4</sub> if <sub>4</sub> S <sub>5</sub> else <sub>6</sub> other <sub>3</sub>	\$	reduce: 2
8	\$ <sub>0</sub> if <sub>4</sub> if <sub>4</sub> S <sub>5</sub> else <sub>6</sub> S <sub>7</sub>	\$	reduce: 4
9	\$ <sub>0</sub> if <sub>4</sub> I <sub>2</sub>	\$	reduce: 1
10	\$ <sub>0</sub> S <sub>1</sub>	\$	accept

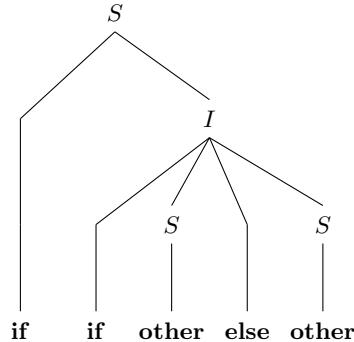
## Parser run, different choice

state	input				goto	
	if	else	other	\$	S	I
0	$s : 4$			$s : 3$		
1					accept	
2			$r : 1$		$r : 1$	
3			$r : 2$		$r : 2$	
4	$s : 4$			$s : 3$		
5			$s : 6$		$r : 3$	
6	$s : 4$			$s : 3$		
7			$r : 4$		$r : 4$	

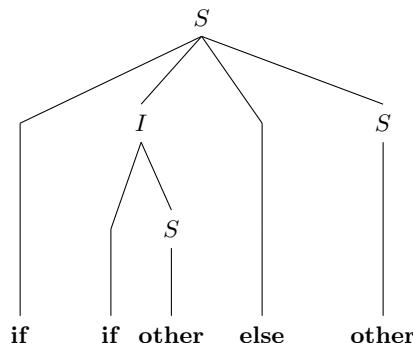
stage	parsing stack	input	action
1	\$ <sub>0</sub>	if if other else other \$	shift: 4
2	\$ <sub>0</sub> if <sub>4</sub>	if other else other \$	shift: 4
3	\$ <sub>0</sub> if <sub>4</sub> if <sub>4</sub>	other else other \$	shift: 3
4	\$ <sub>0</sub> if <sub>4</sub> if <sub>4</sub> other <sub>3</sub>	else other \$	reduce: 2
5	\$ <sub>0</sub> if <sub>4</sub> if <sub>4</sub> S <sub>5</sub>	else other \$	reduce 3
6	\$ <sub>0</sub> if <sub>4</sub> I <sub>2</sub>	else other \$	reduce 1
7	\$ <sub>0</sub> if <sub>4</sub> S <sub>5</sub>	else other \$	shift 6
8	\$ <sub>0</sub> if <sub>4</sub> S <sub>5</sub> else <sub>6</sub>	other \$	shift 3
9	\$ <sub>0</sub> if <sub>4</sub> S <sub>5</sub> else <sub>6</sub> other <sub>3</sub>	\$	reduce 2
10	\$ <sub>0</sub> if <sub>4</sub> S <sub>5</sub> else <sub>6</sub> S <sub>7</sub>	\$	reduce 4
11	\$ <sub>0</sub> S <sub>1</sub>	\$	accept

## Parse trees: simple conditions

- ### 1. shift-precedence: conventional



- ## 2. “wrong” tree



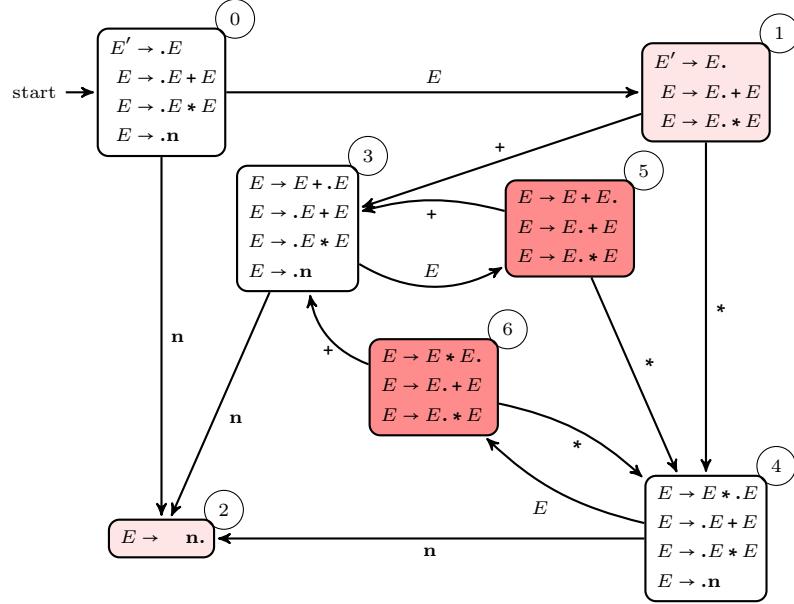
- <sup>3</sup> standard “dangling else” convention “an else belongs to the last previous, still open (= dangling) if-clause”

## Use of ambiguous grammars

- advantage of ambiguous grammars: often simpler
- if ambiguous: grammar guaranteed to have conflicts
- can be (often) resolved by specifying *precedence* and *associativity*
- supported by tools like `yacc` and `CUP` ...

$$\begin{array}{l} E' \rightarrow E \\ E \rightarrow E + E \mid E * E \mid n \end{array}$$

## DFA for + and ×



## States with conflicts

- state 5
  - stack contains ...E+E
  - for input \$: reduce, since shift not allowed from \$
  - for input +; reduce, as + is *left-associative*
  - for input \*: shift, as \* has *precedence* over +
- state 6:
  - stack contains ...E\*E
  - for input \$: reduce, since shift not allowed from \$
  - for input +; reduce, as \* has *precedence* over +
  - for input \*: shift, as \* is *left-associative*
- see also the table on the next slide

## Parse table + and ×

state	input				goto
	n	+	*	\$	
0	s : 2				E
1		s : 3			1
2		r : E → n		r : E → n	
3	s : 2				5
4	s : 2				6
5	r : E → E + E		s : 4	r : E → E + E	
6	r : E → E * E	r : E → E * E	r : E → E * E	r : E → E * E	

1. How about exponentiation (written  $\uparrow$  or  $**$ )? Defined as *right-associative*. See exercise

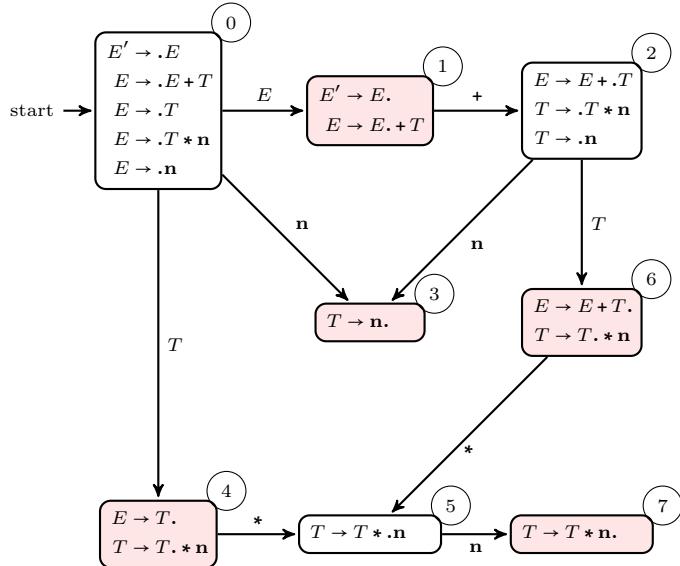
For comparison: unambiguous grammar for + and \*

1. Unambiguous grammar: precedence and left-assoc built in

$$\begin{array}{lcl} E' & \rightarrow & E \\ E & \rightarrow & E+T \mid T \\ T & \rightarrow & T * n \mid n \end{array}$$

	<i>Follow</i>	
$E'$	$\{\$\}$	(as always for start symbol)
$E$	$\{\$, +\}$	
$T$	$\{\$, +, *\}$	

DFA for unambiguous + and ×



DFA remarks

- the DFA now is SLR(1)
  - check states with *complete items*
  - state 1:  $Follow(E') = \{\$\}$
  - state 4:  $Follow(E) = \{\$, +\}$
  - state 6:  $Follow(E) = \{\$, +\}$
  - state 3/7:  $Follow(T) = \{\$, +, *\}$
- in no case there's a shift/reduce conflict (check the outgoing edges vs. the follow set)
- there's not reduce/reduce conflict either

LR(1) parsing

- most general from of LR(1) parsing
  - aka: *canonical* LR(1) parsing
  - usually: considered as unnecessarily “complex” (i.e. LALR(1) or similar is good enough)
  - “stepping stone” towards LALR(1)
1. Basic restriction of SLR(1) Uses *look-ahead*, yes, but only *after* it has built a non-look-ahead DFA (based on **LR(0)-items**)
  2. A help to remember  
SRL(1) “improved” LR(0) parsing LALR(1) is “crippled” LR(1) parsing.

## Limits of SLR(1) grammars

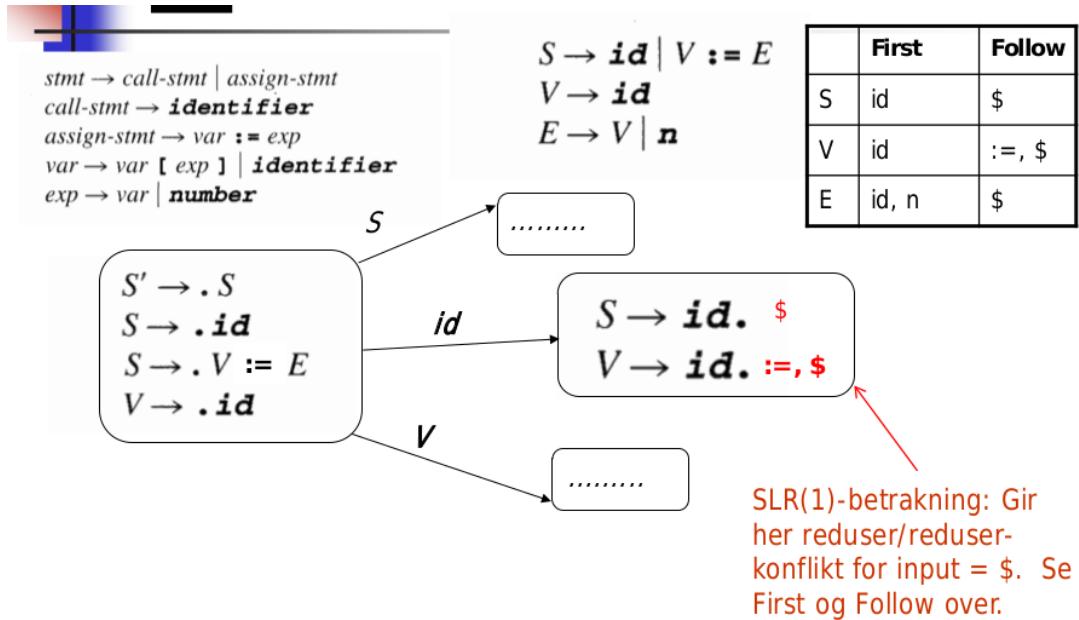
1. Assignment grammar fragment<sup>21</sup>

$$\begin{array}{lcl} \text{stmt} & \rightarrow & \text{call-stmt} \mid \text{assign-stmt} \\ \text{call-stmt} & \rightarrow & \text{identifier} \\ \text{assign-stmt} & \rightarrow & \text{var} := \text{exp} \\ \text{var} & \rightarrow & [\text{exp}] \mid \text{identifier} \\ \text{exp} & \mid & \text{var} \mid \text{n} \end{array}$$

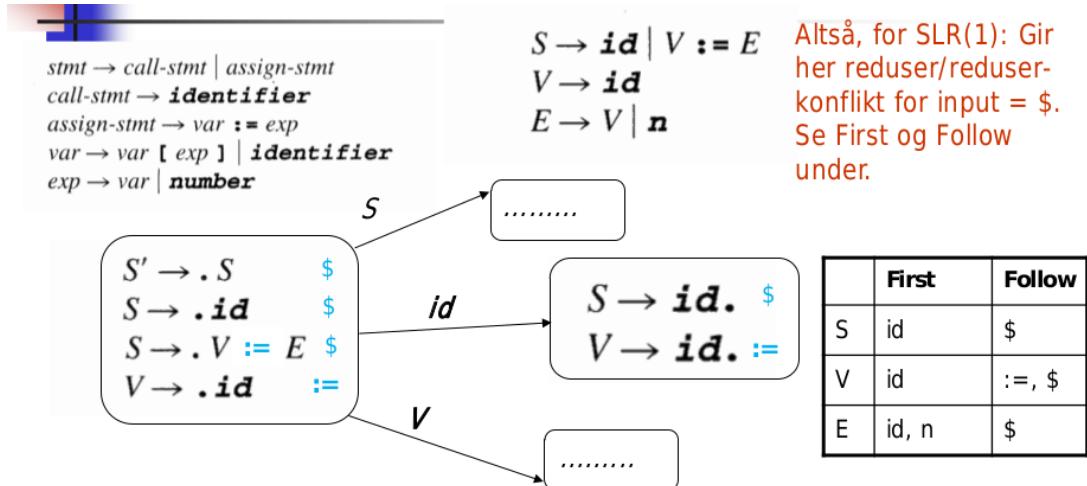
2. Assignment grammar fragment, simplified

$$\begin{array}{lcl} S & \rightarrow & \text{id} \mid V := E \\ V & \rightarrow & \text{id} \\ E & \rightarrow & V \mid n \end{array}$$

## non-SLR(1): Reduce/reduce conflict



## Situation can be saved: more look-ahead



<sup>21</sup>Inspired by Pascal, analogous problems in C ...

### LALR(1) (and LR(1)): Being more precise with the follow-sets

- LR(0)-items: too “indiscriminate” wrt. the follow sets
  - remember the definition of SLR(1) conflicts
  - LR(0)/SLR(1)-states:
    - sets of items<sup>22</sup> due to subset construction
    - the items are LR(0)-items
    - follow-sets as an *after-thought*
1. Add precision in the *states* of the automaton already Instead of using LR(0)-items and, when the LR(0) DFA is done, try to disambiguate with the help of the follow sets for states containing complete items: **make more fine-grained items**:
- **LR(1) items**
  - each *item* with “specific follow information”: look-ahead

### LR(1) items

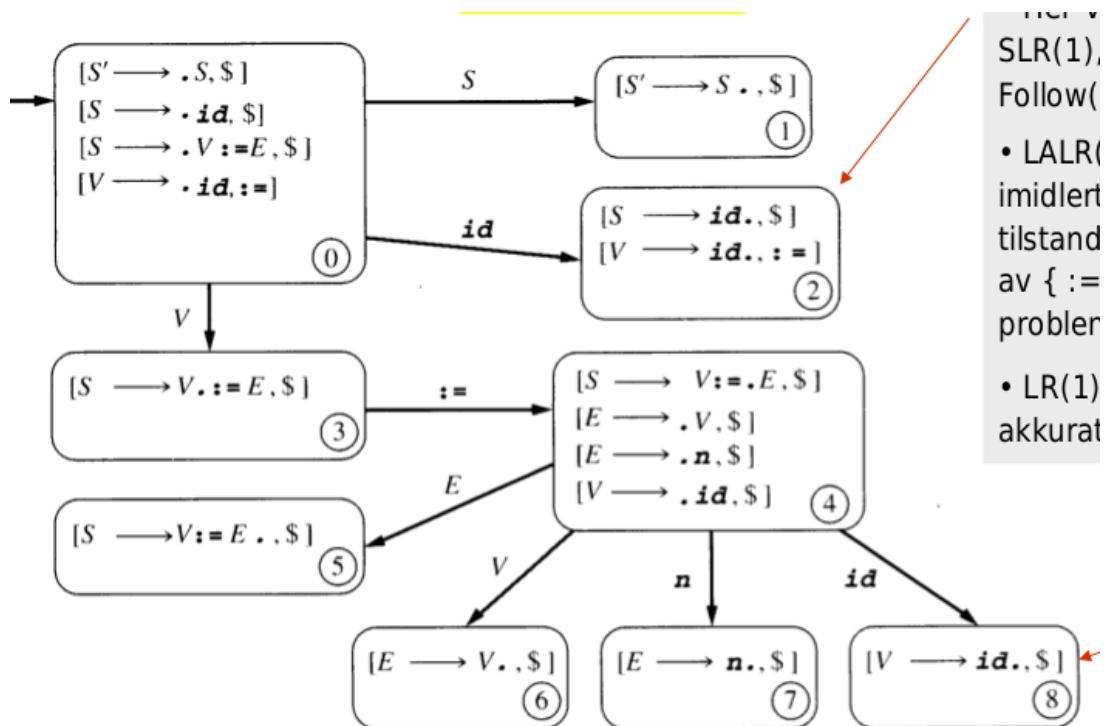
- main idea: simply make the look-ahead part of the item
- obviously: proliferation of states<sup>23</sup>

#### 1. LR(1) items

$$[A \rightarrow \alpha.\beta, a] \quad (9)$$

- $a$ : terminal/token, including  $\$$

### LALR(1)-DFA (or LR(1)-DFA)



### Remarks on the DFA

- Cf. state 2 (seen before)
  - in SLR(1): problematic (reduce/reduce), as  $Follow(V) = \{:=, \$\}$
  - now: disambiguation, by the added information
- LR(1) would give the same DFA

<sup>22</sup>That won't change in principle (but the items get more complex)

<sup>23</sup>Not to mention if we wanted look-ahead of  $k > 1$ , which in practice is not done, though.

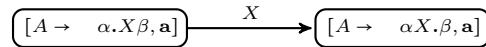
## Full LR(1) parsing

- AKA: **canonical** LR(1) parsing
  - the *best* you can do with 1 look-ahead
  - unfortunately: big tables
  - pre-stage to LALR(1)-parsing
1. SLR(1) LR(0)-item-based parsing, with *afterwards* adding some extra “pre-compiled” info (about follow-sets) to increase expressivity
  2. LALR(1) LR(1)-item-based parsing, but *afterwards* throwing away precision by collapsing states, to save space

## LR(1) transitions: arbitrary symbol

- transitions of the NFA (not DFA)

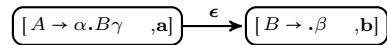
### 1. $X$ -transition



## LR(1) transitions: $\epsilon$

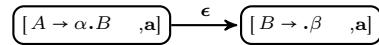
### 1. $\epsilon$ -transition for all

$B \rightarrow \beta_1 \mid \beta_2 \dots$  and all  $b \in First(\gamma a)$



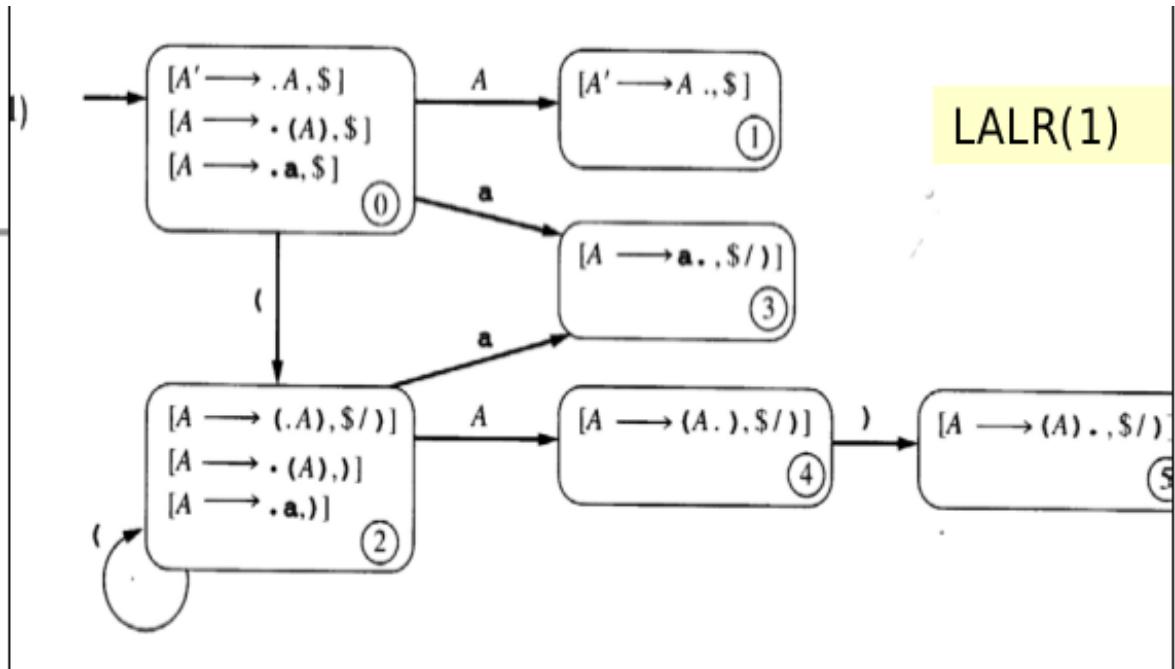
### 2. including special case ( $\gamma = \epsilon$ )

for all  $B \rightarrow \beta_1 \mid \beta_2 \dots$



## LALR(1) vs LR(1)

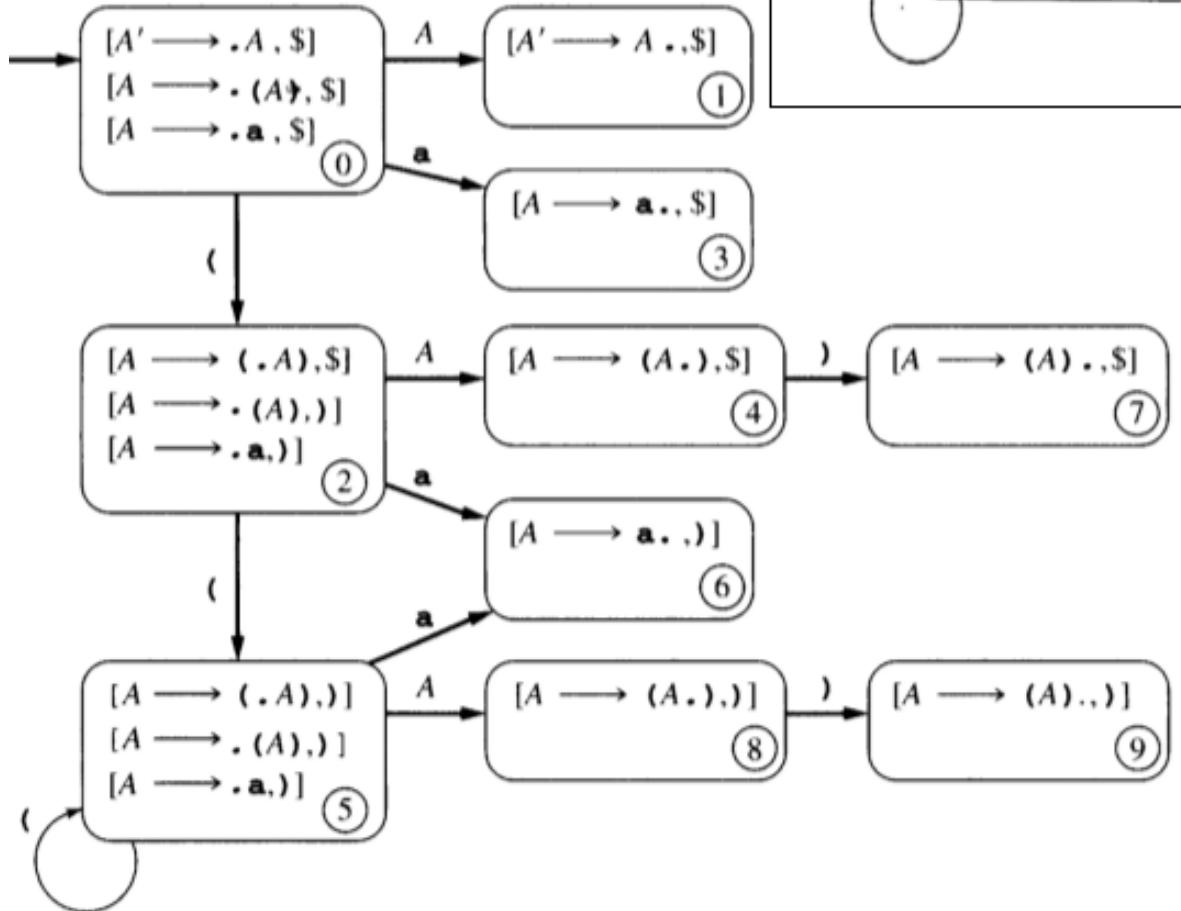
### 1. LALR(1)



### 2. LR(1)

$$A \rightarrow ( A ) \mid a$$

LR(1)



### Core of LR(1)-states

- actually: not done that way in practice
  - main idea: *collapse states with the same core*
1. Core of an LR(1) state = set of LR(0)-items (i.e., ignoring the look-ahead)
  2. Rest
    - observation: core of the LR(1) item = LR(0) item
    - 2 LR(1) states with the same core have same outgoing edges, and those lead to states with the same core

### LALR(1)-DFA by collapse

- collapse all states with the same core
- based on above observations: edges are also consistent
- Result: almost like a LR(0)-DFA but additionally
  - still each individual item has still look ahead attached: the **union** of the “collapsed” items
  - especially for states with *complete* items  $[A \rightarrow \alpha, a, b, \dots]$  is **smaller** than the follow set of  $A$
  - $\Rightarrow$  less unresolved conflicts compared to SLR(1)

## Concluding remarks of LR / bottom up parsing

- all constructions (here) based on BNF (not EBNF)
- *conflicts* (for instance due to ambiguity) can be solved by
  - reformulate the grammar, but generate the same language<sup>24</sup>
  - use *directives* in parser generator tools like `yacc`, `CUP`, `bison` (precedence, assoc.)
  - or (not yet discussed): solve them later via *semantical analysis*
  - NB: *not all* conflicts are solvable, also not in LR(1) (remember ambiguous languages)

## LR/bottom-up parsing overview

	advantages	remarks
LR(0)	defines states <i>also</i> used by SLR and LALR	not really used, many conflicts, very weak
SLR(1)	clear improvement over LR(0) in expressiveness, even if using the same number of states. Table typically with 50K entries	weaker than LALR(1). but often good enough. Ok for hand-made parsers for <i>small</i> grammars
LALR(1)	almost as expressive as LR(1), but number of states as LR(0)!	method of choice for most generated LR-parsers
LR(1)	<i>the</i> method covering <i>all</i> bottom-up, one-look-ahead parseable grammars	large number of states (typically 11M of entries), mostly LALR(1) preferred

Remeber: once the *table* specific for LR(0), ... is set-up, the parsing algorithms all work *the same*

## Error handling

1. Minimal requirement Upon “stumbling over” an error (= deviation from the grammar): give a *reasonable & understandable* error message, indicating also error *location*. Potentially stop parsing
2. Rest
  - for parse error *recovery*
    - one cannot really recover from the fact that the program has an error (an syntax error is a syntax error), but
    - after giving decent error message:
      - \* move on, potentially jump over some subsequent code,
      - \* until parser can *pick up* normal parsing again
      - \* so: meaningful checking code even following a first error
    - avoid: reporting an avalanche of subsequent *spurious* errors (those just “caused” by the first error)
    - “pick up” again after semantic errors: easier than for syntactic errors

## Error messages

- important:
  - avoid error messages that only occur because of an already reported error!
  - report error as early as possible, if possible at the *first point* where the program cannot be extended to a correct program.
  - make sure that, after an error, one doesn’t end up in an *infinite loop* without reading any input symbols.
- What’s a good error message?
  - assume: that the method `factor()` chooses the alternative (`exp`) but that it, when control returns from method `exp()`, does not find a `)`
  - one could report : `left parenthesis missing`
  - But this may often be confusing, e.g. if what the program text is: `( a + b c )`
  - here the `exp()` method will terminate after `( a + b`, as `c` cannot extend the expression). You should therefore rather give the message `error in expression or left parenthesis missing`.

<sup>24</sup>If designing a new language, there’s also the option to massage the language itself. Note also: there are *inherently ambiguous languages* for which there is no *unambiguous* grammar.

## Error recovery in bottom-up parsing

- panic recovery in LR-parsing
    - simple form
    - the only one we shortly look at
  - upon error: recovery ⇒
    - pops parts of the stack
    - ignore parts of the input
  - until “on track again”
  - but: how to do that
  - additional problem: *non-determinism*
    - table: constructed *conflict-free under normal operation*
    - upon error (and clearing parts of the stack + input): no guarantee it's clear how to continue
- ⇒ **heuristic** needed (like panic mode recovery)
1. Panic mode idea
    - try a **fresh start**,
    - promising “fresh start” is: a possible **goto** action
    - thus: back off and take the *next* such goto-opportunity

## Possible error situation

parse stack		input	action
1	$\$_0 a_1 b_2 c_3 ( \underline{d}_4 \underline{d}_5 e_6$	f ) g h ... \$	no entry for f
2	$\$_0 a_1 b_2 c_3 B_v$	g h ... \$	back to normal
3	$\$_0 a_1 b_2 c_3 B_v g_7$	h ... \$	...

state	input				goto			
	...	)	f	g	...	A	B	...
...								
3						<b>u</b>	<b>v</b>	
4			–			–	–	
5			–			–	–	
6		–	–			–	–	
...								
<b>u</b>		–	–	reduce...				
<b>v</b>		–	–	shift: 7				
...								

## Panic mode recovery

1. Algo
  - (a) *Pop* states for the stack *until* a state is found with non-empty **goto** entries
  - (b)
    - If there's legal action on the current input token from one of the goto-states, push token on the stack, *restart* the parse.
    - If there's several such states: *prefer shift* to a reduce
    - Among possible reduce actions: prefer one whose associated non-terminal is least general
  - (c) if no legal action on the current input token from one of the goto-states: *advance input* until there is a legal action (or until end of input is reached)

## Example again

parse stack		input	action
1	$\$_0 a_1 b_2 c_3 ( \underline{d}_4 \underline{d}_5 e_6$	f ) g h ... \$	no entry for f
2	$\$_0 a_1 b_2 c_3 B_v$	g h ... \$	back to normal
3	$\$_0 a_1 b_2 c_3 B_v g_7$	h ... \$	...

- first pop, until in state 3
- then jump over input
  - until next input **g**
  - since **f** and **)** cannot be treated
- choose to goto **v** (shift in that state)

## Panic mode may loop forever

	parse stack	input	action
1	\$ <sub>0</sub>	( n n ) \$	
2	\$ <sub>0</sub> ( <sub>6</sub>	n n ) \$	
3	\$ <sub>0</sub> ( <sub>6</sub> n <sub>5</sub>	n ) \$	
4	\$ <sub>0</sub> ( <sub>6</sub> factor <sub>4</sub>	n ) \$	
6	\$ <sub>0</sub> ( <sub>6</sub> term <sub>3</sub>	n ) \$	
7	\$ <sub>0</sub> ( <sub>6</sub> exp <sub>10</sub>	n ) \$	panic!
8	\$ <sub>0</sub> ( <sub>6</sub> factor <sub>4</sub>	n ) \$	been there before: stage 4!

## Typical yacc parser table

- some variant of the expression grammar again

$$\begin{array}{lcl}
 \text{command} & \rightarrow & \text{exp} \\
 \text{exp} & \rightarrow & \text{term} * \text{factor} \mid \text{factor} \\
 \text{term} & \rightarrow & \text{term} * \text{factor} \mid \text{factor} \\
 \text{factor} & \rightarrow & \mathbf{n} \mid (\text{exp})
 \end{array}$$

State	Input								Goto			
	NUMBER	(	+	-	*	)	\$	command	exp	term	factor	
0	s5	s6							1	2	3	4
1								accept				
2	r1	r1	s7	s8	r1	r1	r1					
3	r4	r4	r4	r4	s9	r4	r4					
4	r6	r6	r6	r6	r6	r6	r6					
5	r7	r7	r7	r7	r7	r7	r7					
6	s5	s6							10	3	4	
7	s5	s6							11	4		
8	s5	s6							12	4		
9	s5	s6*								13		
10			s7	s8			s14					
11	r2	r2	r2	r2	s9	r2	r2					
12	r3	r3	r3	r3	r3	s9	r3	r3				
13	r5	r5	r5	r5	r5	r5	r5	r5				
14	r8	r8	r8	r8	r8	r8	r8	r8				

## Panicking and looping

	parse stack	input	action
1	\$ <sub>0</sub>	( n n ) \$	
2	\$ <sub>0</sub> ( <sub>6</sub>	n n ) \$	
3	\$ <sub>0</sub> ( <sub>6</sub> n <sub>5</sub>	n ) \$	
4	\$ <sub>0</sub> ( <sub>6</sub> factor <sub>4</sub>	n ) \$	
6	\$ <sub>0</sub> ( <sub>6</sub> term <sub>3</sub>	n ) \$	
7	\$ <sub>0</sub> ( <sub>6</sub> exp <sub>10</sub>	n ) \$	panic!
8	\$ <sub>0</sub> ( <sub>6</sub> factor <sub>4</sub>	n ) \$	been there before: stage 4!

- error raised in stage 7, no action possible

- panic:

- pop-off exp<sub>10</sub>
- state 6: 3 goto's

	exp	term	factor
goto to with n next: action there	10	3	4

— reduce r<sub>4</sub>    reduce r<sub>6</sub>

- no shift, so we need to decide between the two reduces

- factor: less general, we take that one

## How to deal with looping panic?

- make sure to detect loop (i.e. previous “configurations”)
- if loop detected: don't repeat but do something special, for instance
  - pop-off more from the stack, and try again

- pop-off and *insist* that a shift is part of the options
1. Left out (from the book and the pensum)
    - more info on error recovery
    - especially: more on `yacc` error recovery
    - it's not pensum, and for the oblig: need to deal with CUP-specifics (not classic `yacc` specifics even if similar) anyhow, and error recovery is not part of the oblig (halfway decent error *handling* is).

## 2 Reference

### References

- [Aho et al., 2007] Aho, A. V., Lam, M. S., Sethi, R., and Ullman, J. D. (2007). *Compilers: Principles, Techniques and Tools*. Pearson/Addison-Wesley, second edition.
- [Appel, 1998a] Appel, A. W. (1998a). *Modern Compiler Implementation in Java*. Cambridge University Press.
- [Appel, 1998b] Appel, A. W. (1998b). *Modern Compiler Implementation in ML*. Cambridge University Press.
- [Appel, 1998c] Appel, A. W. (1998c). *Modern Compiler Implementation in ML/Java/C*. Cambridge University Press.
- [Louden, 1997] Louden, K. (1997). *Compiler Construction, Principles and Practice*. PWS Publishing.

# Index

\$ (end marker symbol), 8  
abstract syntax tree, 17  
ambiguity of a grammar, 22  
associativity, 12, 29, 30  
bottom-up parsing, 43  
complete item, 52  
constraint, 4  
context-free grammar  
    reduced, 25  
CUP, 77  
dangling-else, 73  
determinization, 15  
EBNF, 14, 27, 35, 36  
 $\epsilon$ -production, 12, 13  
First set, 1  
Follow set, 1  
follow set, 8, 47  
grammar  
    ambiguous, 22  
    LL(1), 25  
    LL(K), 24  
    start symbol, 8  
handle, 46  
higher-order rewriting, 26  
initial item, 52  
item  
    complete, 52  
    initial, 52  
LALR(1), 43  
left factor, 12  
left recursion, 12  
left-derivation., 25  
left-factoring, 11, 27, 38  
left-recursion, 11, 12, 28, 29, 38  
    immediate, 11  
linear production, 25  
LL(1), 25, 27  
LL(1) grammars, 37  
LL(1) parse table, 38  
LL(k), 24  
LL(k)-grammar, 24  
LR(0), 43, 52, 66  
LR(1), 43  
non-determinism, 26  
non-terminal symbol, 26  
nullable, 2, 24  
nullable symbols, 1  
parse  
    error, 84  
parse tree, 26  
parser, 17  
    predictive, 27  
    recursive descent, 27  
parsing  
    bottom-up, 43  
precedence, 29  
predict-set, 26  
predictive parser, 27  
prefix  
    viable, 51  
production  
    linear, 25  
recursive descent parser, 27  
reduced context-free grammar, 25  
rewriting, 25  
    higher-order, 26  
sentential form, 1  
shift-reduce parser, 44  
SLR(1), 43, 66  
string rewriting, 26  
syntax error, 17  
term rewriting, 26  
terminal symbol, 26  
transducer  
    tree, 26  
transduction, 26  
tree transducer, 26  
type error, 17  
viable prefix, 51  
worklist, 6, 8, 9  
worklist algorithm, 6, 8, 9  
yacc, 77