

INF5110 – Compiler Construction

Parsing

Spring 2016



1. Parsing

Bottom-up parsing

Bibs

1. Parsing

Bottom-up parsing

Bibs

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

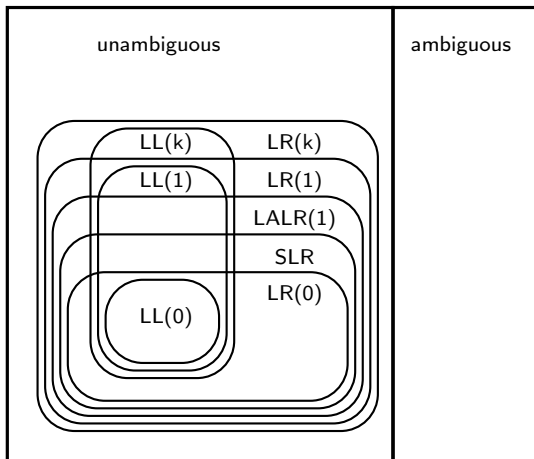
- LR(0)
- only for very simple grammars
 - approx 300 states for standard programming languages
 - only as intro to SLR(1) and LALR(1)

- 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

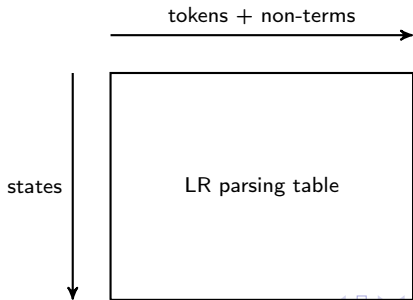
- LR(1) covers all grammars, which can in principle be parsed by looking at the next token

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 it's descendants (like bison, CUP, etc)
- another name: *shift-reduce* parser



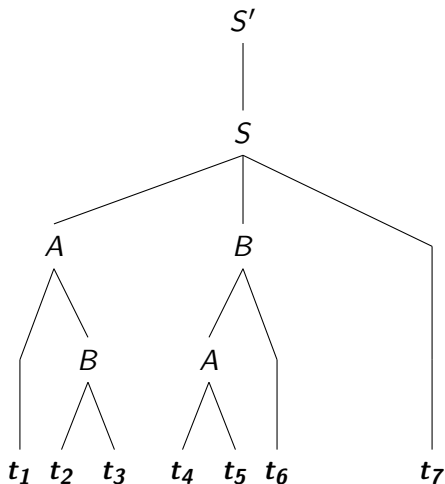
Example grammar

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

- 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 of a production
 - **routinely add another “extra” start-symbol** (here S')¹

¹That will later be used when constructing a DFA for “scanning” the stack, to control the reactions of the stack machine. This restriction leads to a unique, well-define initial state.

Parse tree for $t_1 \dots t_7$



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 can the parser do a *right*-most derivation, when it parses from *left*-to-right?

- 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 Σ^* derivable from start-symbol

Right-sentential form: right-most derivation

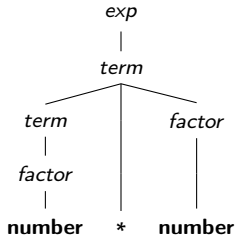
$$S \Rightarrow_r^* \alpha$$

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)

$exp \rightarrow exp\ addop\ term \mid term$ (1)
 $addop \rightarrow + \mid -$
 $term \rightarrow term\ mulop\ term \mid factor$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid \mathbf{number}$



Bottom-up parse: Growing the parse tree

number * number

number * number

Bottom-up parse: Growing the parse tree

$$\begin{array}{c} \textit{factor} \\ | \\ \text{number} \end{array} * \text{number}$$

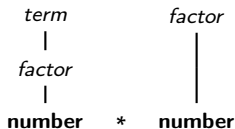
number * number \hookrightarrow factor * number

Bottom-up parse: Growing the parse tree

$$\begin{array}{c} \textit{term} \\ | \\ \textit{factor} \\ | \\ \text{number} \end{array} * \text{number}$$

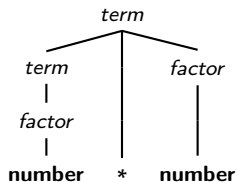
number * number \hookrightarrow factor * number
 \hookrightarrow *term* * number

Bottom-up parse: Growing the parse tree



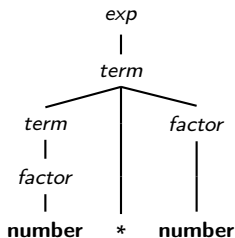
number * number \hookrightarrow factor * number
 \hookrightarrow term * number
 \hookrightarrow term * factor

Bottom-up parse: Growing the parse tree



number * number \hookrightarrow factor * number
 \hookrightarrow term * number
 \hookrightarrow term * factor
 \hookrightarrow term

Bottom-up parse: Growing the parse tree



number * number \hookrightarrow factor * number
 \hookrightarrow term * number
 \hookrightarrow term * factor
 \hookrightarrow term
 \hookrightarrow exp

Reduction in reverse = right derivation

Reduction

$$\begin{aligned}\underline{n} * n &\hookrightarrow \underline{factor} * n \\ &\hookrightarrow \underline{term} * \underline{n} \\ &\hookrightarrow \underline{term} * \underline{factor} \\ &\hookrightarrow \underline{term} \\ &\hookrightarrow \underline{exp}\end{aligned}$$

Right derivation

$$\begin{aligned}n * n &\leftarrow_r \underline{factor} * n \\ &\leftarrow_r \underline{term} * n \\ &\leftarrow_r \underline{term} * \underline{factor} \\ &\leftarrow_r \underline{term} \\ &\leftarrow_r \underline{exp}\end{aligned}$$

- 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**
- note: all intermediate words are *right-sentential forms*

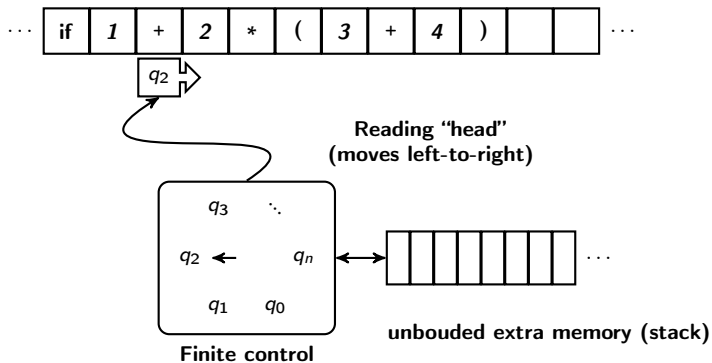
Definition (Handle)

Assume $S \Rightarrow_r^* \alpha A w \Rightarrow_r \alpha \beta w$. A production $A \rightarrow \beta$ at position k following α 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.
- 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 β (= *bottom-up!*)
- “handle” β can be *empty* (= ϵ)

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 will be: how do achieve the same *without that tree already given* (just parsing left-to-right)

Schematic run (reduction: from top to bottom)

\$	$t_1 t_2 t_3 t_4 t_5 t_6 t_7$	\$
\$ t_1	$t_2 t_3 t_4 t_5 t_6 t_7$	\$
\$ $t_1 t_2$	$t_3 t_4 t_5 t_6 t_7$	\$
\$ $t_1 t_2 t_3$	$t_4 t_5 t_6 t_7$	\$
\$ $t_1 B$	$t_4 t_5 t_6 t_7$	\$
\$ A	$t_4 t_5 t_6 t_7$	\$
\$ $A t_4$	$t_5 t_6 t_7$	\$
\$ $A t_4 t_5$	$t_6 t_7$	\$
\$ AA	$t_6 t_7$	\$
\$ $AA t_6$	t_7	\$
\$ AB	t_7	\$
\$ $AB t_7$		\$
\$ 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

Shift

Move the next input symbol (terminal) over to the top of the stack (“push”)

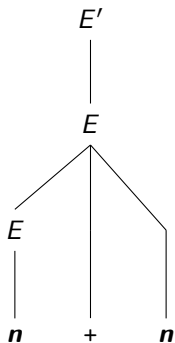
Reduce

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”).

- *easy* to do if one has the parse tree already!

Example: LR parsing for addition (given the tree)

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



	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

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

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

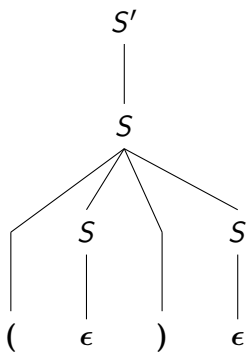
Example with ϵ -transitions: parentheses

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

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



	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	\$\$	\$	reduce $S' \rightarrow S$
7	\$\$S'	\$	accept

Note: the 2 reduction steps for the ϵ productions

Right-most derivation and right-sentential forms

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

Right-sentential forms & the stack

- sentential form: word from Σ^* derivable from start-symbol

Right-sentential form: right-most derivation

$$S \Rightarrow_r^* \alpha$$

- right-sentential forms:
 - 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

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

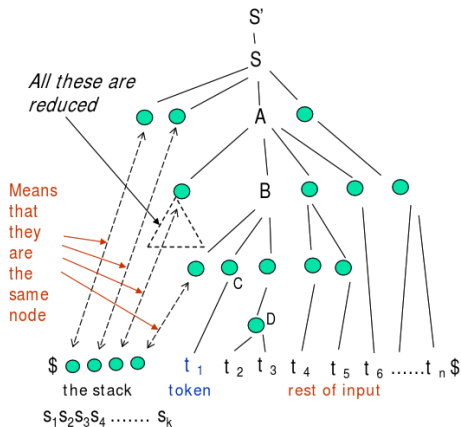
$$\underline{n+n} \hookrightarrow \underline{E+n} \hookrightarrow \underline{E} \hookrightarrow E'$$

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

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): depends on current “state” of the parser engine

A typical situation during LR-parsing



The stack is reduced version of the processed input

After a shift, the next reduction to be made is a reduction with the production:

$C \rightarrow t_1$

Then, after two shifts, we will make a reduction with the production:

$D \rightarrow t_2 t_3$

Then, what's next?

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, one 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)

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

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

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 $k > 1$

LR(0) item

production with specific “parser position” \cdot in its right-hand side

- \cdot is, of course, a “meta-symbol” (not part of the production)
- For instance: production $A \rightarrow \alpha$, where $\alpha = \beta\gamma$, then

LR(0) item

$$A \rightarrow \beta \cdot \gamma$$

- item with dot at the beginning: *initial* item
- item with dot at the end: *complete* item

Example: items of LR-grammar

Grammar for parentheses: 3 productions

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

8 items

$$\begin{aligned}S' &\rightarrow \cdot S \\S' &\rightarrow S \cdot \\S &\rightarrow \cdot (S)S \\S &\rightarrow (\cdot S)S \\S &\rightarrow (S \cdot)S \\S &\rightarrow (S) \cdot S \\S &\rightarrow (S)S \cdot \\S &\rightarrow \cdot\end{aligned}$$

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

Another example: items for addition grammar

Grammar for addition: 3 productions

$$\begin{aligned} E' &\rightarrow E \\ E &\rightarrow E + \text{number} \mid \text{number} \end{aligned}$$

(coincidentally also:) 8 items

$$\begin{aligned} E' &\rightarrow \cdot E \\ E' &\rightarrow E \cdot \\ E &\rightarrow \cdot E + \text{number} \\ E &\rightarrow E \cdot + \text{number} \\ E &\rightarrow E + \cdot \text{number} \\ E &\rightarrow E + \text{number} \cdot \\ E &\rightarrow \cdot \text{number} \\ E &\rightarrow \text{number} \cdot \end{aligned}$$

- 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

States formed of sets of items

In a state marked by/containing item

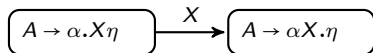
$$A \rightarrow \beta.\gamma$$

- β on the *stack*
- γ : 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

Terminal or non-terminal



Epsilon (X : non-terminal here)



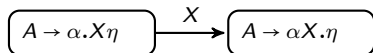
- In case $X = \textit{terminal}$ (i.e. token) =
 - the left step corresponds to a **shift** step²
- for non-terminals (see next slide):
 - interpretation more complex: non-terminals are officially never on the input
 - note: in that case, item $A \rightarrow \alpha.X\eta$ has two (kind of) outgoing transitions

²We have explained **shift** steps so far as: parser eats one **terminal** (= input token) and pushes it on the stack.

Transitions for non-terminals and ϵ

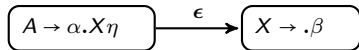
- so far: we never pushed a non-terminal from the input to the stack, we **replace** in a **reduce**-step the right-hand side by a left-hand side
- however: the replacement in a **reduce** steps can be seen as
 1. pop right-hand side off the stack,
 2. instead, “assume” corresponding non-terminal on input &
 3. eat the non-terminal and push it on the stack.
- two kind of transitions
 1. the ϵ -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 $X \rightarrow \cdot\beta$

Terminal or non-terminal



Epsilon (X : non-terminal here)

Given production $X \rightarrow \beta$:



Initial and final states

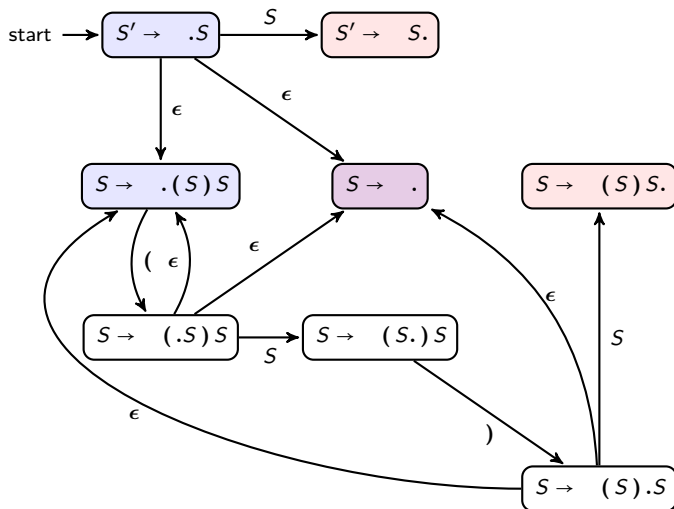
initial states:

- we make our lives *easier*
 - we assume (as said): one *extra* start symbol say S' (augmented grammar)
- ⇒ initial item $S' \rightarrow .S$ as (only) **initial state**

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)

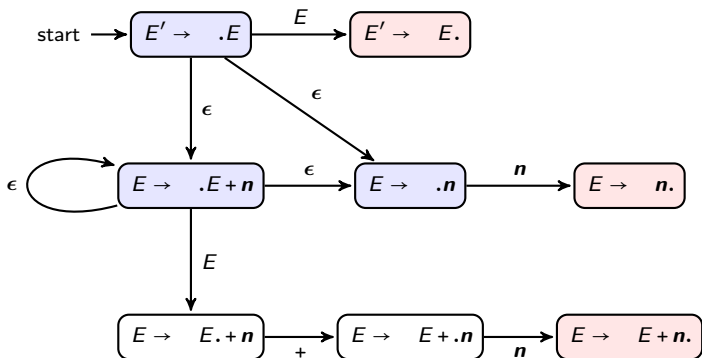
NFA: parentheses



- 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 ϵ -transitions go into, but
 - *with exception* of the initial state (with S' -production)

no outgoing edges from the complete items

NFA: addition

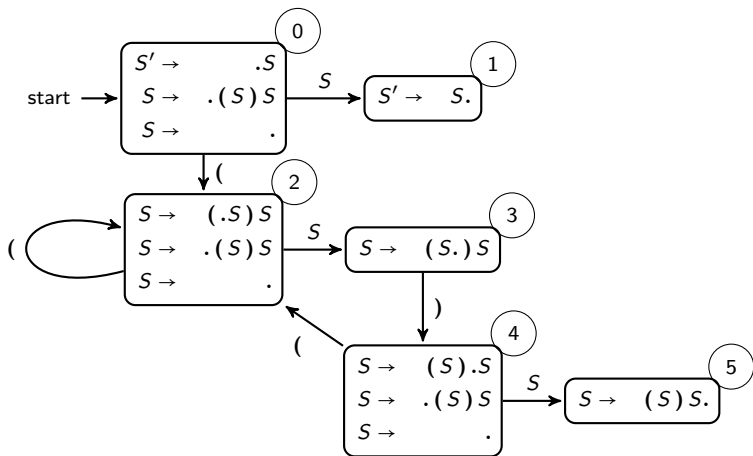


Determinizing: from NFA to DFA

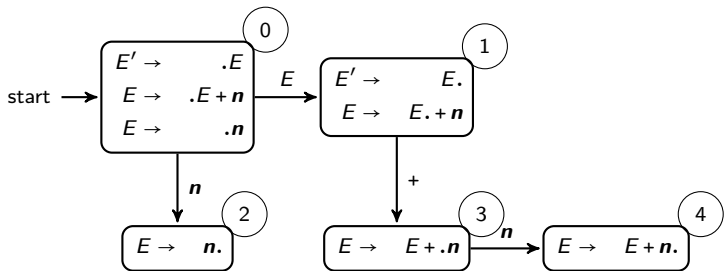
- standard subset-construction³
- states then contains *sets* of items
- especially important: ϵ -closure
- also: *direct* construction of the DFA possible

³technically, we don't require here a *total* transition function, we leave out any error state.

DFA: parentheses



DFA: addition



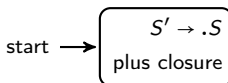
Direct construction of an LR(0)-DFA

- quite easy: simply build in the closure already

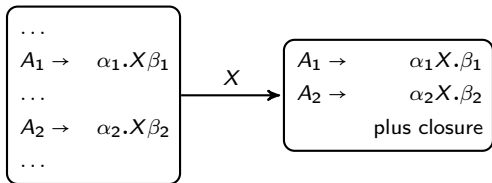
ϵ -closure

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

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

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)

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 Σ^*
 2. DFA operates deterministically on such words
- the stack contains the “past”: read input (and 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

⇒ 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

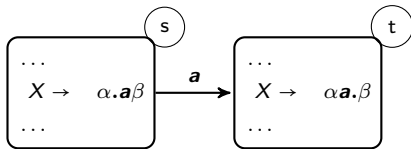
⇒ crucial when determining *reaction*

State transition allowing a shift

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

$$X \rightarrow \alpha \cdot a \beta$$

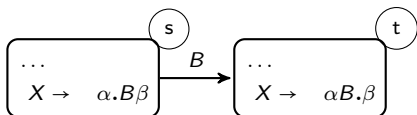
- construction thus has transition as follows



- 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

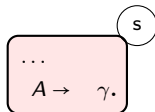
$$X \rightarrow \alpha.B\beta$$



- 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:
 1. pop off the handle from the top of the stack
 2. put the non-term B "back onto the input" (corresponding to the above state s)
 3. 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 \cdot at the end)
- assume **top state** s containing complete item $A \rightarrow \gamma$.



- a complete right-hand side (“handle”) γ on the stack and thus done
 - may be replaced by right-hand side A
- ⇒ reduce step
- builds up (implicitly) new parent node A in the bottom-up procedure
 - **Note:** A on top of the stack instead of γ :⁴
 - **new top state!**
 - remember the “goto-transition” (shift of a non-terminal)

⁴Indirectly 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.

Remarks: states, transitions, and reduce steps

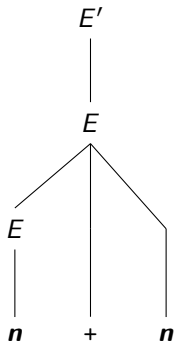
- ignoring the ϵ -transitions (for the NFA)
- there are 2 “kinds” of transitions in the DFA
 1. terminals: reals shifts
 2. non-terminals: “following a reduce step”

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
-
- if 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

Example: LR parsing for addition (given the tree)

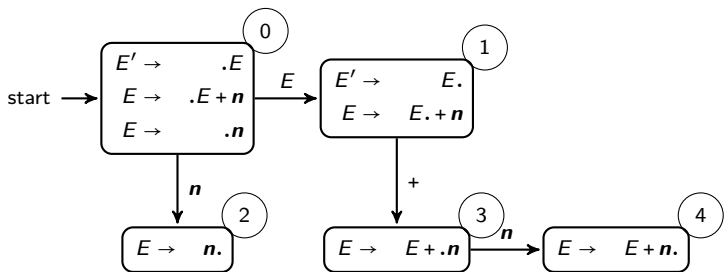
$$\begin{aligned} E' &\rightarrow E \\ E &\rightarrow E+n \mid n \end{aligned}$$



	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

DFA of addition example



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

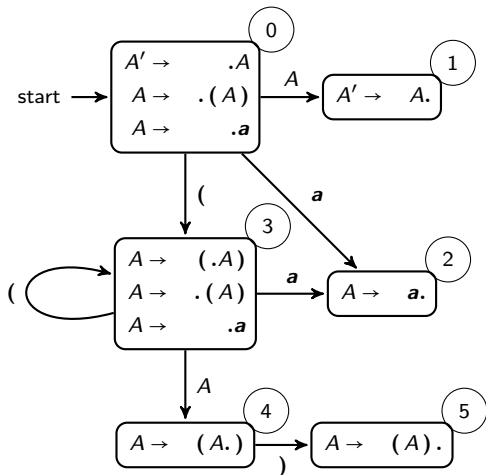
LR(0) grammar

The top-state alone determines the next step.

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

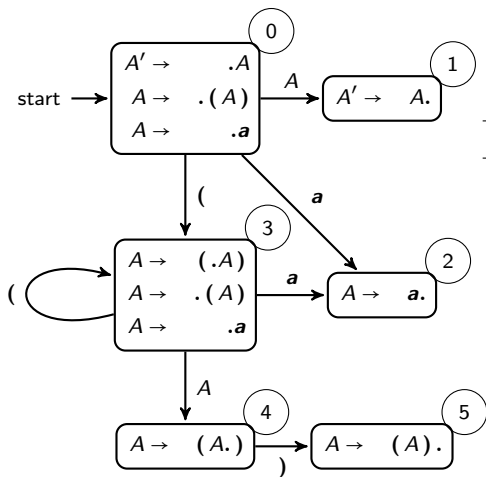
Simple parentheses

$A \rightarrow (A) \mid a$



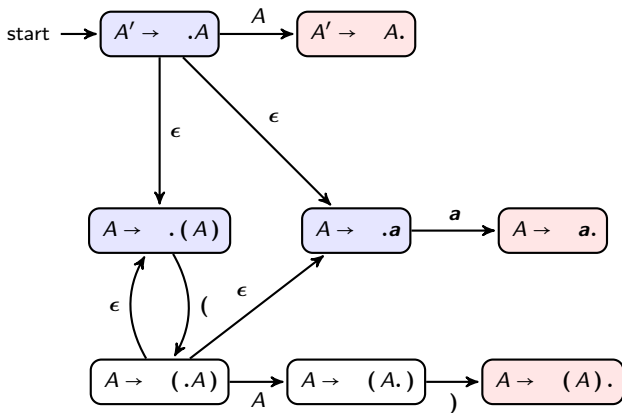
- for *shift*:
 - many shift transitions in state allowed
 - shift counts as *one* action (including “shifts” on non-terms)
- but for reduction: also the *production* must be clear

Simple parentheses is LR(0)



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 one non-terminal here)
- corresponding to the A -labelled transitions
- see the parser run on the next slide

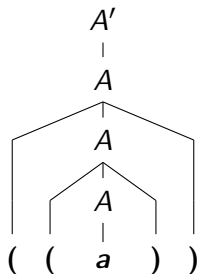
state	action	rule	input	goto
			(a)	A
0	shift		3 2	1
1	reduce	$A' \rightarrow A$		
2	reduce	$A \rightarrow a$		
3	shift		3	4
4	shift			5
5	reduce	$A \rightarrow (A)$		

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)
- \Rightarrow 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	$\$_0$	$((a)\$$	shift
2	$\$_0(\textit{a})$	$(a)\$$	shift
3	$\$_0(\textit{a}(\textit{a}))$	$a)\$$	shift
4	$\$_0(\textit{a}(\textit{a}A_2))$	$)\$$	reduce $A \rightarrow a$
5	$\$_0(\textit{a}(\textit{a}A_4))$	$)\$$	shift
6	$\$_0(\textit{a}(\textit{a}A_4)_5)$	$\$$	reduce $A \rightarrow (A)$
7	$\$_0(\textit{a}A_4)$	$\$$????

stage	parsing stack	input	action
1	$\$_0$	$()\$$	shift
2	$\$_0(\textit{a})$	$)\$$?????

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 γ (including “its” states from the stack)

back up: assume to be in state u which is *now* head state

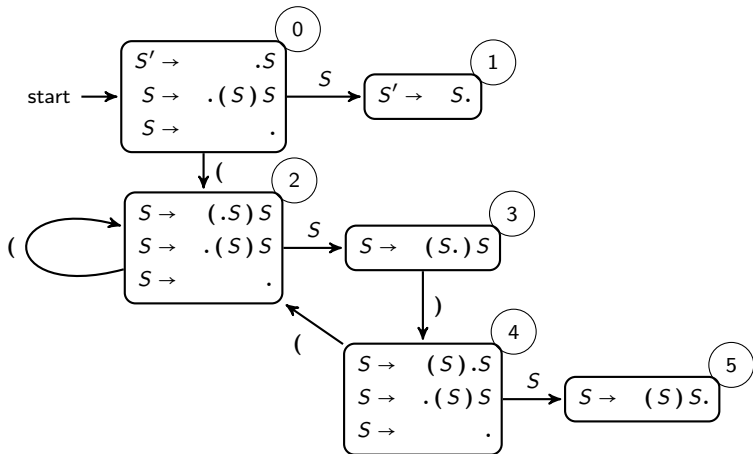
push: push A to the stack, new head state t where

$$u \xrightarrow{A} t$$

- in [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)

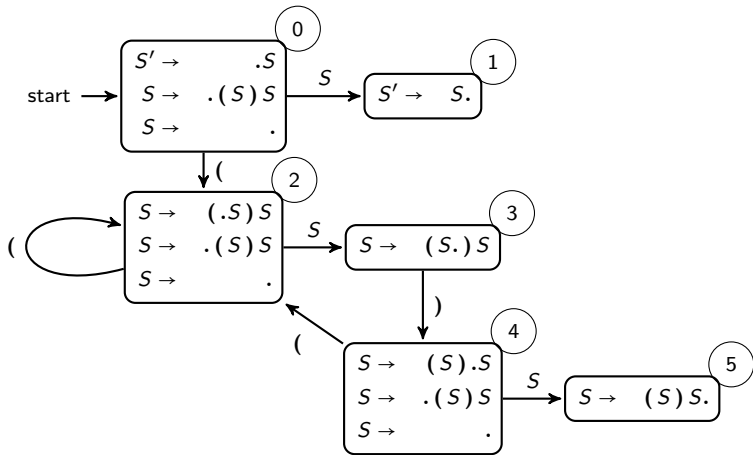
DFA parentheses again: LR(0)?

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



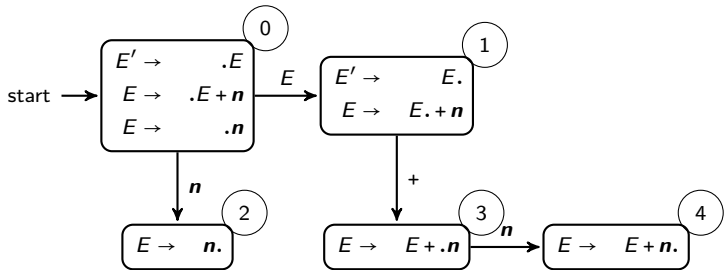
DFA parentheses again: LR(0)?

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

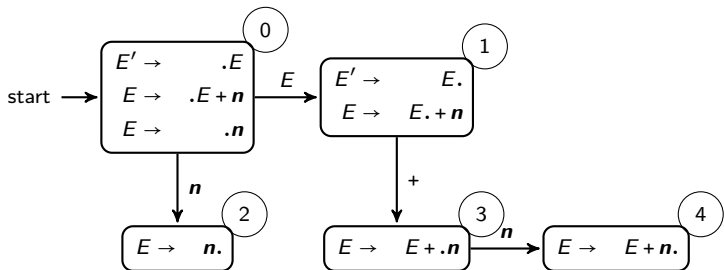


Look at states 0, 2, and 4

DFA addition again: LR(0)?

$$E' \rightarrow E$$
$$E \rightarrow E + \text{number} \mid \text{number}$$


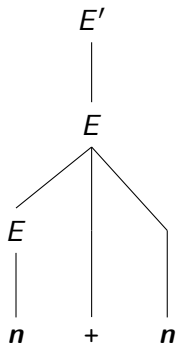
DFA addition again: LR(0)?

$$E' \rightarrow E$$
$$E \rightarrow E + \text{number} \mid \text{number}$$


How to make a decision in state 1?

Decision? If only we knew the ultimate tree already ...

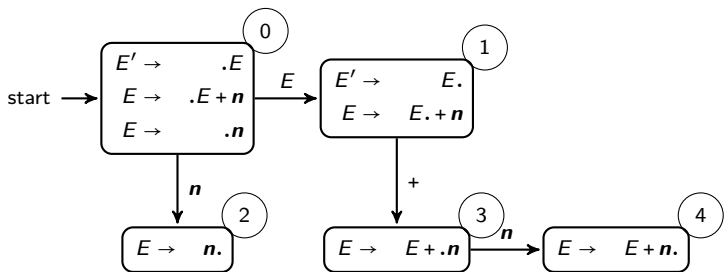
... especially the parts still to come



	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

- current stack: representation of the 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 value of 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 variable)

- LR(0), not very useful, much 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)-DFA's

LR(0) reduce/reduce conflict:

...

$A \rightarrow \alpha.$

...

$B \rightarrow \beta.$

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

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

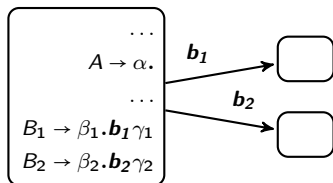
LR(0) reduce/reduce conflict:

...
$A \rightarrow \alpha.$
...
$B \rightarrow \beta.$

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

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

LR(0) shift/reduce conflict:



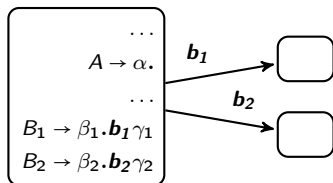
SLR(1) solution: again: use follow sets of non-terms

- If $Follow(A) \cap \{b_1, b_2, \dots\} = \emptyset$

\Rightarrow next symbol (in token) decides!

- if $token \in Follow(A)$ then *reduce* using $A \rightarrow \alpha$, non-terminal A determines new top state
- if $token \in \{b_1, b_2, \dots\}$ then *shift*. Input symbol b_i determines new top state
- \dots

LR(0) shift/reduce conflict:



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

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

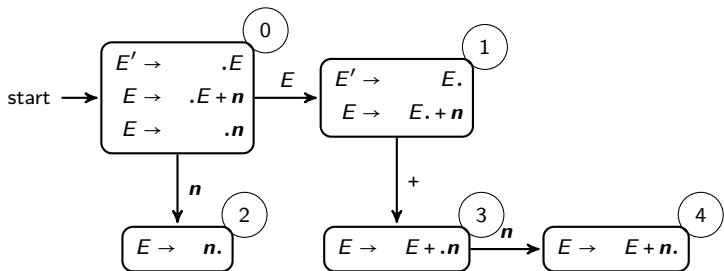
SLR(1) requirement on states (as in the book)

- formulated as conditions on the states (of LR(0)-items)
- given the LR(0)-item DFA as defined

SLR(1) condition, on all states s

1. 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 $X \in \text{Follow}(B)$.
2. For any **two complete** items $A \rightarrow \alpha.$ and $B \rightarrow \beta.$ in s ,
 $\text{Follow}(\alpha) \cap \text{Follow}(\beta) = \emptyset$

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)

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 $s \xrightarrow{X} t^5$
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 empty⁶
 - else:
 - pop*: remove γ (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

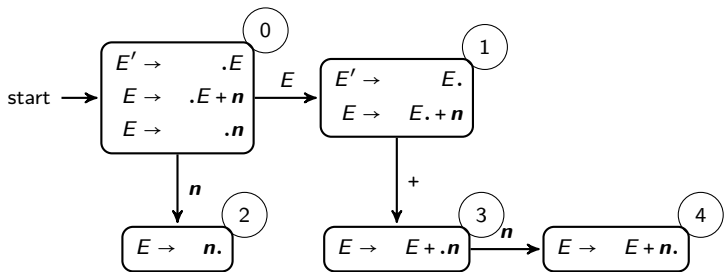
$$u \xrightarrow{A} t$$

3. if next token is such that neither 1. or 2. applies: *error*

⁵Cf. to the LR(0) algo: since we checked the existence of the transition before, the else-part is missing now.

⁶Cf. 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 for SLR(1)



state	input			goto
	n	$+$	$\$$	E
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: $n \notin \text{Follow}(E)$

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 **uniformly**: either shift or else reduce (with given rule)
 - now: non-uniform, **dependent** on the input
- 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)-table⁷

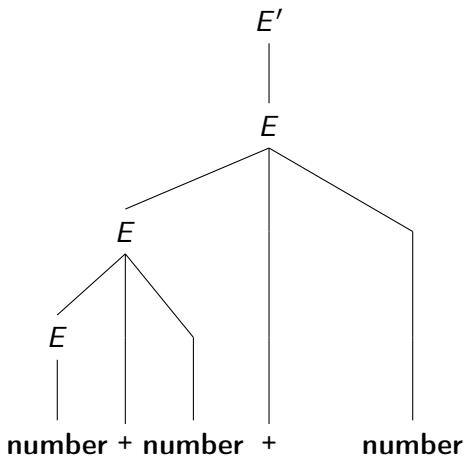
⁷by which it, strictly speaking, would no longer be an SRL(1)-table :)

SLR(1) parser run (= "reduction")

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

stage	parsing stack	input	action
1	$\$_0$	$n + n + n\$$	shift: 2
2	$\$_0 n_2$	$+ n + n\$$	reduce: $E \rightarrow n$
3	$\$_0 E_1$	$+ n + n\$$	shift: 3
4	$\$_0 E_{1+3}$	$n + n\$$	shift: 4
5	$\$_0 E_{1+3} n_4$	$+ n\$$	reduce: $E \rightarrow E + n$
6	$\$_0 E_1$	$n\$$	shift 3
7	$\$_0 E_{1+3}$	$n\$$	shift 4
8	$\$_0 E_{1+3} n_4$	$\$$	reduce: $E \rightarrow E + n$
9	$\$_0 E_1$	$\$$	accept

Corresponding parse tree



Revisit the parentheses again: SLR(1)?

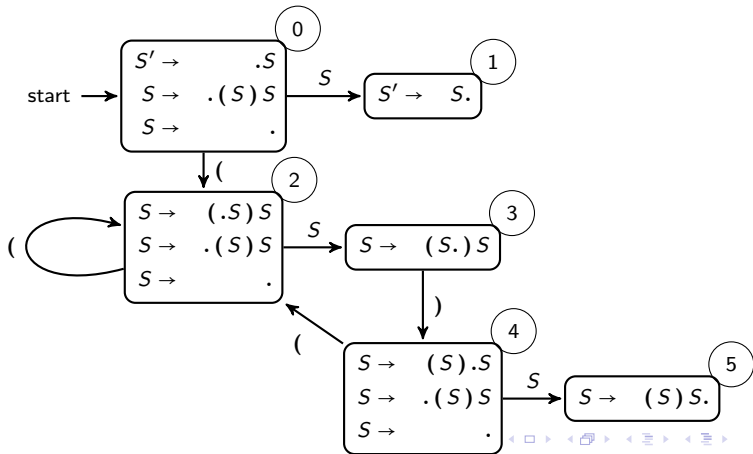
Grammar: parentheses (from before)

$$S' \rightarrow S$$

$$S \rightarrow (S)S \mid \epsilon$$

Follow set

$$\text{Follow}(S) = \{), \$\}$$



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	$\$0(2S3$	$)() \$$	shift: 4
4	$\$0(2S3)_4$	$() \$$	shift: 2
5	$\$0(2S3)_4(2$	$) \$$	reduce: $S \rightarrow \epsilon$
6	$\$0(2S3)_4(2S3$	$) \$$	shift: 4
7	$\$0(2S3)_4(2S3)_4$	$\$$	reduce: $S \rightarrow \epsilon$
8	$\$0(2S3)_4(2S3)_4S5$	$\$$	reduce: $S \rightarrow (S) S$
9	$\$0(2S3)_4S5$	$\$$	reduce: $S \rightarrow (S) S$
10	$\$0S2$	$\$$	accept

- in principle: straightforward: k look-ahead, instead of 1
- rarely used in practice, using $First_k$ and $Follow_k$ 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 $First_k$ and $Follow_k$ 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 $Xw \in First_k(X\beta)$ 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 $w \in Follow_k(A)$ 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 look-ahead)
- However: ambiguous grammar tolerable, if (remaining) conflicts can be solved meaningfully otherwise:

Additional means of disambiguation:

1. by specifying associativity / precedence “outside” the grammar
 2. by “living with the fact” that LR parser (commonly) *prioritizes shifts over reduces*
- for the second point (“let the parser decide according to its preferences”):
 - use sparingly and cautiously
 - typical example: *dangling-else*
 - even if parser makes a decision, programmer may or may not “understand intuitively” the resulting parse tree (and thus AST)
 - a grammar with many S/R-conflicts: go back to the drawing board

Example of an ambiguous grammar

$$\begin{aligned} \textit{stmt} &\rightarrow \textit{if-stmt} \mid \textit{other} \\ \textit{if-stmt} &\rightarrow \mathbf{if} (\textit{exp}) \textit{stmt} \\ &\rightarrow \mathbf{if} (\textit{exp}) \textit{stmt} \mathbf{else} \textit{stmt} \\ \textit{exp} &\rightarrow \mathbf{0} \mid \mathbf{1} \end{aligned}$$

In the following E for exp etc.

Simplified conditionals

Simplified “schematic” if-then-else

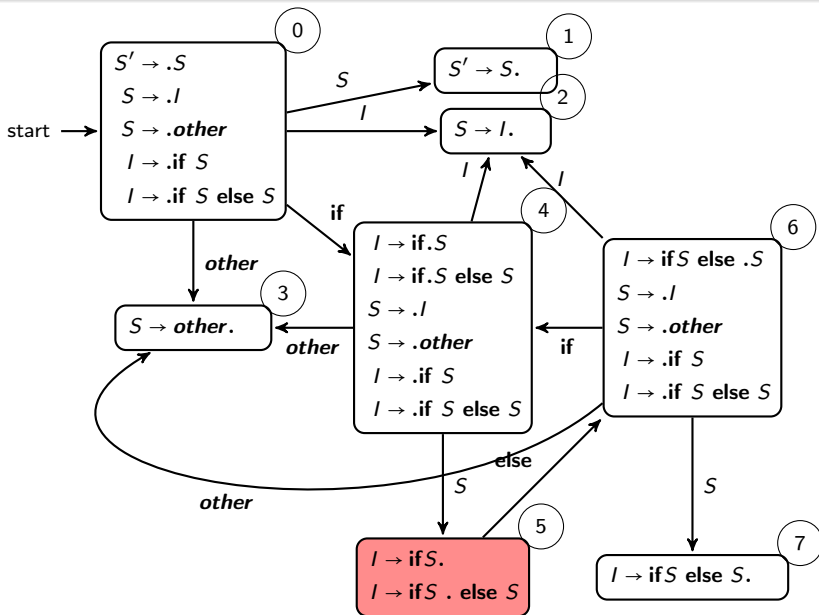
$$S \rightarrow I \mid \mathit{other}$$
$$I \rightarrow \mathit{if } S \mid \mathit{if } S \mathit{ else } S$$

Follow-sets

	<i>Follow</i>
<i>S'</i>	{ $\$$ }
<i>S</i>	{ $\$, \mathit{else}$ }
<i>I</i>	{ $\$, \mathit{else}$ }

- since ambiguous: at least one conflict must be somewhere

DFA of LR(0) items



SLR(1)-parse-table, conflict resolved

Grammar

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

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

- *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	
	if	else	other	\$	/
0	s: 4		s: 3		1 2
1				accept	
2		r: 1		r: 1	
3		r: 2		r: 2	
4	s: 4		s: 3		5 2
5		s: 6		r: 3	
6	s: 4		s: 3		7 2
7		r: 4		r: 4	

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

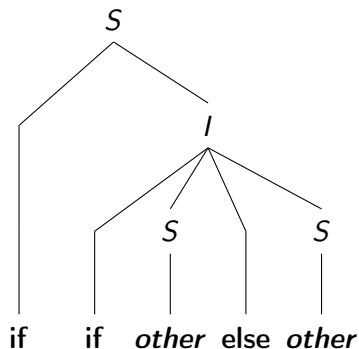
Parser run, different choice

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

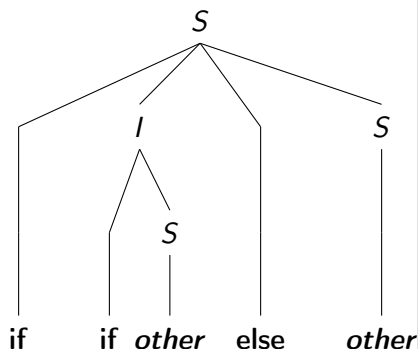
stage	parsing stack	input	action
1	\$ ₀	if if <i>other</i> else <i>other</i> \$	shift: 4
2	\$ ₀ if ₄	if <i>other</i> else <i>other</i> \$	shift: 4
3	\$ ₀ if ₄ if ₄	<i>other</i> else <i>other</i> \$	shift: 3
4	\$ ₀ if ₄ if ₄ <i>other</i> ₃	else <i>other</i> \$	reduce: 2
5	\$ ₀ if ₄ if ₄ S ₅	else <i>other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

Parse trees: simple conditions

shift-precedence: conventional



"wrong" tree



"dangling else"

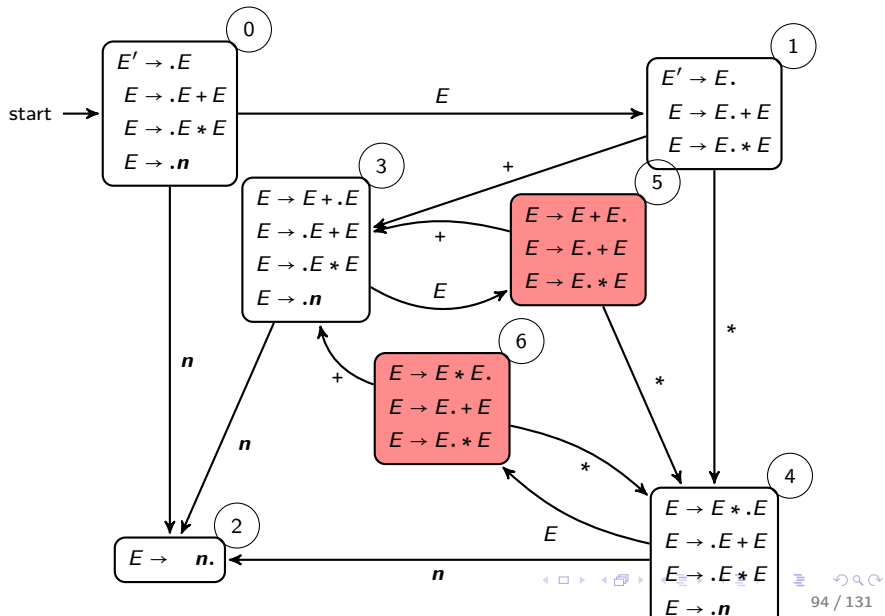
"an **else** belongs to the last previous, still open (= dangling) if-clause"

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 ...

$$\begin{aligned} E' &\rightarrow E \\ E &\rightarrow E + E \mid E * E \mid \text{number} \end{aligned}$$

DFA for + and ×



States with conflicts

- state 5
 - stack contains $\dots E + E$
 - 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 $\dots E * E$
 - 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 \times

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

How about exponentiation (written \uparrow or $**$)?

Defined as *right-associative*. See exercise

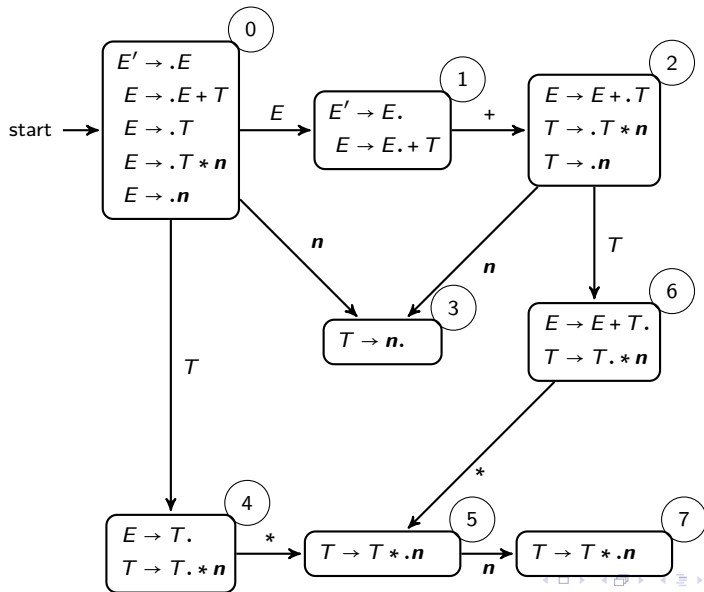
For comparison: unambiguous grammar for + and *

Unambiguous grammar: precedence and left-assoc built in

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

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

DFA for unambiguous + and ×



- the DFA now is SLR(1)
 - check in particular states with complete items
 - state 1: $Follow(E') = \{\$ \}$
 - state 4: $Follow(E) = \{\$, + \}$
 - state 6: $Follow(E) = \{\$, + \}$
 - state 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

- most general form 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)

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)

Assignment grammar fragment^a

^aInspired by Pascal, analogous problems in C ...

$$\begin{aligned} \textit{stmt} &\rightarrow \textit{call-stmt} \mid \textit{assign-stmt} \\ \textit{call-stmt} &\rightarrow \mathbf{id} \\ \textit{assign-stmt} &\rightarrow \textit{var} := \textit{exp} \\ \textit{var} &\rightarrow [\textit{exp}] \mid \mathbf{id} \\ \textit{exp} &\mid \textit{var} \mid \mathbf{number} \end{aligned}$$

Assignment grammar fragment, simplified

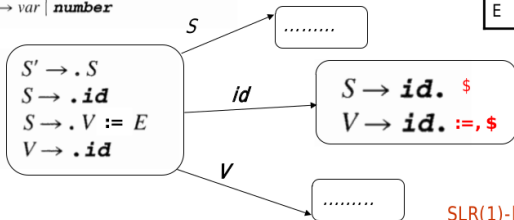
$$\begin{aligned} S &\rightarrow \mathbf{id} \mid V := E \\ V &\rightarrow \mathbf{id} \\ E &\rightarrow V \mid \mathbf{n} \end{aligned}$$

non-SLR(1): Reduce/reduce conflict

$stmt \rightarrow call-stmt \mid assign-stmt$
 $call-stmt \rightarrow identifier$
 $assign-stmt \rightarrow var := exp$
 $var \rightarrow var [exp] \mid identifier$
 $exp \rightarrow var \mid number$

$S \rightarrow id \mid V := E$
 $V \rightarrow id$
 $E \rightarrow V \mid n$

	First	Follow
S	id	\$
V	id	:=, \$
E	id, n	\$



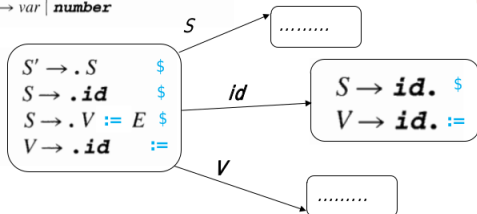
SLR(1)-betragtning: Gir her reduser/reduser-konflikt for input = \$. Se First og Follow over.

Situation can be saved

```
stmt → call-stmt | assign-stmt
call-stmt → identifier
assign-stmt → var := exp
var → var [ exp ] | identifier
exp → var | number
```

$S \rightarrow id \mid V := E$
 $V \rightarrow id$
 $E \rightarrow V \mid n$

Altså, for SLR(1): Gir her reduser/reduser-konflikt for input = \$.
Se First og Follow under.



	First	Follow
S	id	\$
V	id	:=, \$
E	id, n	\$

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 items⁸ due to subset construction
 - the items are LR(0)-items
 - follow-sets as an *after-thought*

Adding 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

⁸That won't change in principle (but the items get more complex)

LR(1) items

- main idea: simply make the look-ahead part of the item
- obviously: proliferation of states⁹

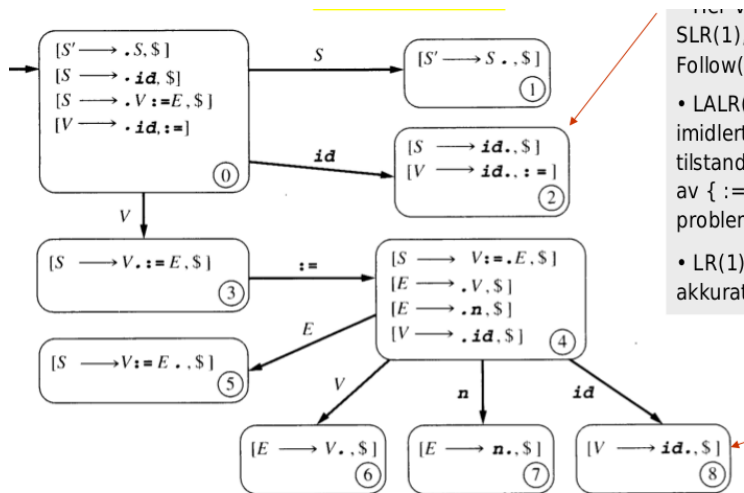
LR(1) items

$$[A \rightarrow \alpha.\beta, \mathbf{a}] \quad (2)$$

- \mathbf{a} : terminal/token, including \$

⁹Not to mention if we wanted look-ahead of $k > 1$, which in practice is not done, though

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



SLR(1),
Follow(
• LALR(
imidert
tilstand
av { :=
proble
• LR(1)
akkurat

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

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

SLR(1)

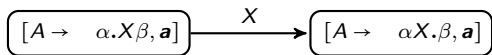
LR(0)-item-based parsing, with *afterwards* adding some extra “pre-compiled” info (about follow-sets) to increase expressivity

LALR(1)

LR(1)-item-based parsing, but *afterwards* throwing away precision by collapsing states, to save space

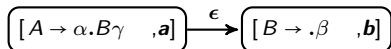
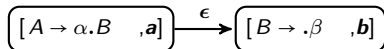
- transitions of the **NFA** (not DFA)

X-transition



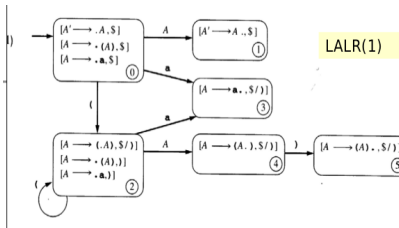
ϵ -transition

for all

 $B \rightarrow \beta_1 \mid \beta_2 \dots$ and all $b \in \text{First}(\gamma a)$ Special case ($\gamma = \epsilon$)for all $B \rightarrow \beta_1 \mid \beta_2 \dots$ 

LALR(1) vs LR(1)

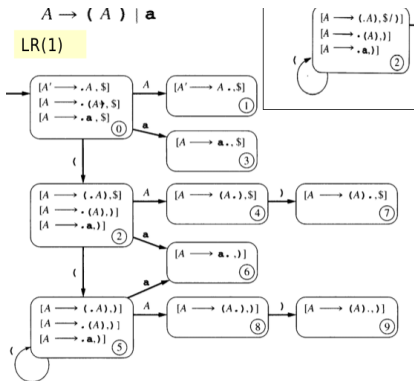
LALR(1)



LR(1)

$A \rightarrow (A) \mid a$

LR(1)



- actually: not done that way in practice
- main idea: *collapse* states with the same *core*

Core of an LR(1) state

= set of *LR(0)*-items (i.e., ignoring the look-ahead)

- 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 as 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 $[A \rightarrow \alpha, \mathbf{a}, \mathbf{b}, \dots]$ 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 language¹⁰
 - 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)

¹⁰If 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.

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

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

Minimal requirement

Upon “stumbling over” an error (= deviation from the grammar): give a *reasonable* & *understandable* error message, indicating also error *location*. Potentially stop parsing

- 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

- 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 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`.

Error recovery in bottom-up parsing

- *panic recovery* in LR-parsing
 - simple form
 - the only one we shortly look at
 - upon error: recovery \Rightarrow
 - 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
- \Rightarrow *heuristic* needed (like panic mode recovery)

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 (₄ d ₅ e ₆	f) gh...\$	no entry for <i>f</i>

state	input				goto			
	...)	<i>f</i>	<i>g</i>	...	<i>A</i>	<i>B</i>	...
...								
3						<i>bcellblueu</i>	<i>bcellv</i>	
4			-			-	-	
5			-			-	-	
6		-	-			-	-	
...								
<i>u</i>		-	-	reduce...				
<i>v</i>		-	-	shift : 7				
...								

Possible error situation

	parse stack	input	action
1	\$ ₀ <i>a</i> ₁ <i>b</i> ₂ <i>c</i> ₃ (₄ <i>d</i> ₅ <i>e</i> ₆	<i>f</i>) <i>gh</i> ...\$	no entry for <i>f</i>
2	\$ ₀ <i>a</i> ₁ <i>b</i> ₂ <i>c</i> ₃ <i>B</i> _v	<i>gh</i> ...\$	back to normal
3	\$ ₀ <i>a</i> ₁ <i>b</i> ₂ <i>c</i> ₃ <i>B</i> _v <i>g</i> ₇	<i>h</i> ...\$...

state	input				goto			
	...)	<i>f</i>	<i>g</i>	...	A	B	...
...								
3						<i>bcellblue</i> _u	<i>bcellv</i>	
4			-			-	-	
5			-			-	-	
6		-	-			-	-	
...								
<i>u</i>		-	-	reduce ...				
<i>v</i>		-	-	shift : 7				
...								

Algo

1. *Pop* states for the stack *until* a state is found with non-empty *goto* entries
2.
 - 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
3. 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	\$ $a_1 b_2 c_3 ($ $d_5 e_6$	$f) gh \dots$ \$	no entry for f

- 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)

Example again

	parse stack	input	action
1	$\$_0 a_1 b_2 c_3 ({}_4 d_5 e_6$	$f) gh \dots \$$	no entry for f
2	$\$_0 a_1 b_2 c_3 B_v$	$gh \dots \$$	back to normal
3	$\$_0 a_1 b_2 c_3 B_v g_7$	$h \dots \$$...

- 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	$\$0$	$(n\ n)\$$	
2	$\$0($	$n\ n)\$$	
3	$\$0(n$	$n)\$$	
4	$\$0(n\ factor_4$	$n)\$$	
6	$\$0(n\ term_3$	$n)\$$	
7	$\$0(n\ exp_{10}$	$n)\$$	panic!
8	$\$0(n\ factor_4$	$n)\$$	been there before: stage 4!

Typical yacc parser table

some variant of the expression grammar again

command → *exp*
exp → *term* * *factor* | *factor*
term → *term* * *factor* | *factor*
factor → **number** | (*exp*)

State	Input							Goto			
	NUMBER	(+	-	*)	\$	<i>command</i>	<i>exp</i>	<i>term</i>	<i>factor</i>
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	s9	r3	r3				
→ 13	r5	r5	r5	r5	r5	r5	r5				
→ 14	r8	r8	r8	r8	r8	r8	r8				

Panicking and looping

	parse stack	input	action
1	$\$0$	$(n\ n)\$$	
2	$\$0(n_6$	$n\ n)\$$	
3	$\$0(n_6n_5$	$n)\$$	
4	$\$0(n_6factor_4$	$n)\$$	
6	$\$0(n_6term_3$	$n)\$$	
7	$\$0(n_6exp_{10}$	$n)\$$	panic!
8	$\$0(n_6factor_4$	$n)\$$	been there before: stage 4!

- error raised in stage 7, no action possible
- panic:
 1. pop-off exp_{10}
 2. state 6: 3 goto's

	<i>exp</i>	<i>term</i>	<i>factor</i>
goto to	10	3	4
with n next: action there	—	reduce r_4	reduce r_6

3. no shift, so we need to decide between the two reduces
4. *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

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).

1. Parsing

Bottom-up parsing

Bibs

[Louden, 1997] Loudon, K. (1997).
Compiler Construction, Principles and Practice.
PWS Publishing.