

Lecture 7

CIS 341: COMPILERS

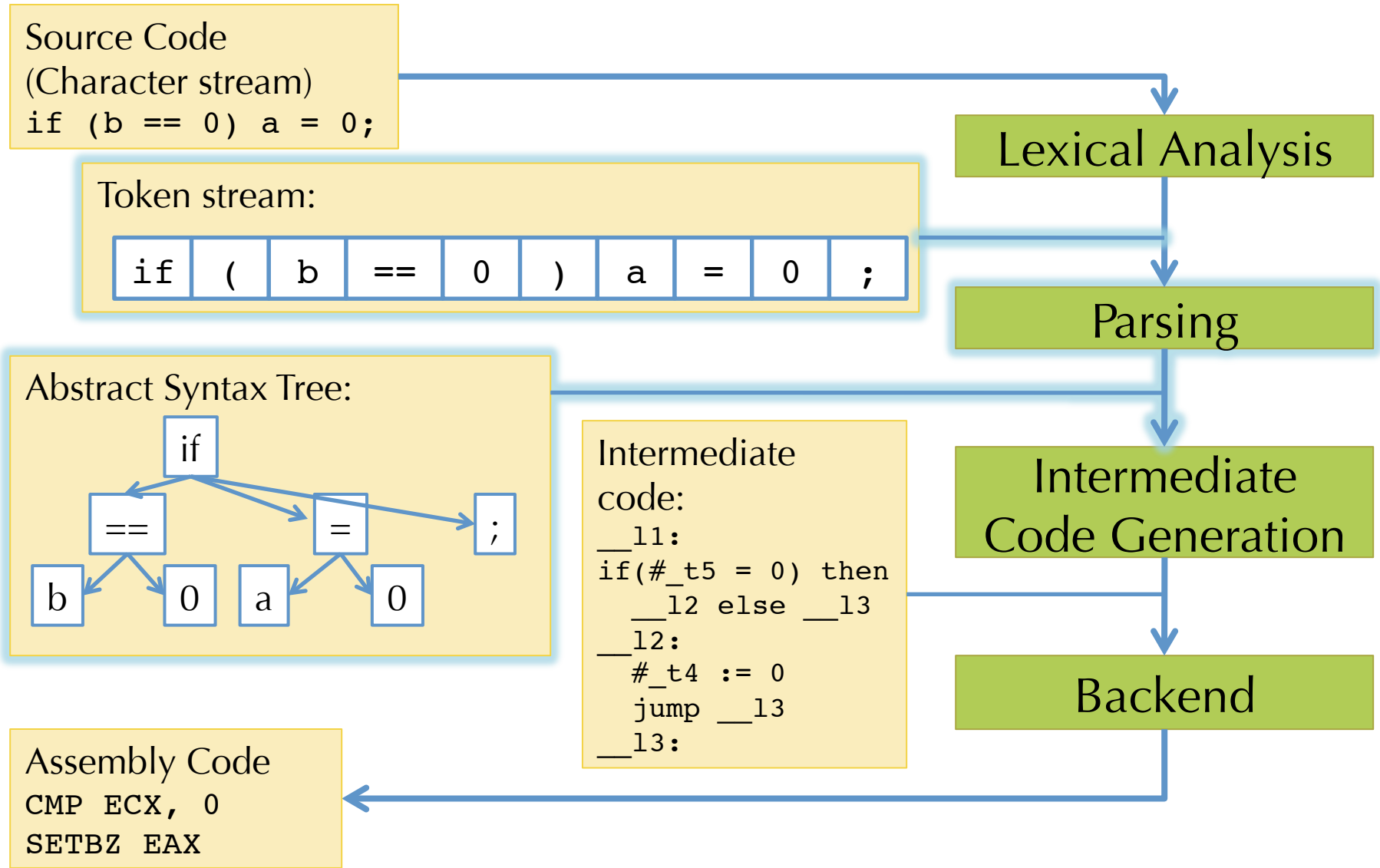
Announcements

- Project 2: Parsing and Compiling Expressions
 - Due: Tuesday, Feb 12th at 11:59:59pm

Searching for derivations.

LL & LR PARSING

Today: Parsing II



CFGs Mathematically

- A Context-free Grammar (CFG) consists of
 - A set of *terminals* (e.g., a token or ϵ)
 - A set of *nonterminals* (e.g., S and other syntactic variables)
 - A designated nonterminal called the *start symbol*
 - A set of productions: $\text{LHS} \mapsto \text{RHS}$
 - LHS is a nonterminal
 - RHS is a *string* of terminals and nonterminals

- Example: The balanced parentheses language:

$$S \mapsto (S)S$$

$$S \mapsto \epsilon$$

- How many terminals? How many nonterminals? Productions?

Consider finding left-most derivations

- Look at only one input symbol at a time.

$$\begin{aligned} S &\mapsto E + S \mid E \\ E &\mapsto \text{number} \mid (S) \end{aligned}$$

Partly-derived String	Look-ahead	Parsed /Unparsed Input
<u>S</u>	((1 + 2 + (3 + 4)) + 5
\mapsto <u>E</u> + S	((1 + 2 + (3 + 4)) + 5
\mapsto (<u>S</u>) + S	1	(1 + 2 + (3 + 4)) + 5
\mapsto (<u>E</u> + S) + S	1	(1 + 2 + (3 + 4)) + 5
\mapsto (1 + <u>S</u>) + S	2	(1 + 2 + (3 + 4)) + 5
\mapsto (1 + <u>E</u> + S) + S	2	(1 + 2 + (3 + 4)) + 5
\mapsto (1 + 2 + <u>S</u>) + S	((1 + 2 + (3 + 4)) + 5
\mapsto (1 + 2 + <u>E</u>) + S	((1 + 2 + (3 + 4)) + 5
\mapsto (1 + 2 + (<u>S</u>)) + S	3	(1 + 2 + (3 + 4)) + 5
\mapsto (1 + 2 + (<u>E</u> + S)) + S	3	(1 + 2 + (3 + 4)) + 5
\mapsto ...		

There is a problem

- We want to decide which production to apply based on the look-ahead symbol.
- But, there is a choice:

$$\begin{array}{l} S \mapsto E + S \mid E \\ E \mapsto \text{number} \mid (S) \end{array}$$

(1) $S \mapsto E \mapsto (S) \mapsto (E) \mapsto (1)$

vs.

(1) + 2 $S \mapsto E + S \mapsto (S) + S \mapsto (E) + S \mapsto (1) + S \mapsto (1) + E$
 $\mapsto (1) + 2$

- Given the look-ahead symbol: '(' it isn't clear whether to pick $S \mapsto E$ or $S \mapsto E + S$ first.

LL(1) GRAMMARS

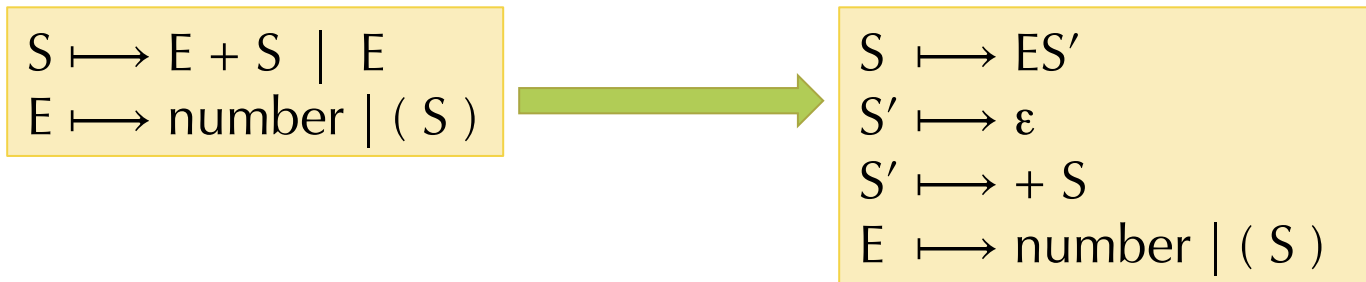
Grammar is the problem

- Not all grammars can be parsed “top-down” with only a single lookahead symbol.
- *Top-down*: starting from the start symbol (root of the parse tree) and going down
- LL(1) means
 - Left-to-right scanning
 - Left-most derivation,
 - 1 lookahead symbol
- This language isn’t “LL(1)”
- Is it LL(k) for some k?
- What can we do?

$$\begin{array}{l} S \mapsto E + S \mid E \\ E \mapsto \text{number} \mid (S) \end{array}$$

Making a grammar LL(1)

- *Problem:* We can't decide which S production to apply until we see the symbol after the first expression.
- *Solution:* "Left-factor" the grammar. There is a common S prefix for each choice, so add a new non-terminal S' at the decision point:



- Also need to eliminate left-recursion somehow. Why?
- Consider:

$$\begin{aligned} S &\mapsto S + E \mid E \\ E &\mapsto \text{number} \mid (S) \end{aligned}$$

LL(1) Parse of the input string

- Look at only one input symbol at a time.

$$\begin{aligned} S &\mapsto ES' \\ S' &\mapsto \varepsilon \\ S' &\mapsto + S \\ E &\mapsto \text{number} \mid (S) \end{aligned}$$

Partly-derived String	Look-ahead	Parsed/Unparsed Input
<u>S</u>	((1 + 2 + (3 + 4)) + 5
\mapsto <u>E</u> S'	((1 + 2 + (3 + 4)) + 5
\mapsto (<u>S</u>) S'	1	(1 + 2 + (3 + 4)) + 5
\mapsto (<u>E</u> S') S'	1	(1 + 2 + (3 + 4)) + 5
\mapsto (1 <u>S'</u>) S'	+	(1 + 2 + (3 + 4)) + 5
\mapsto (1 + <u>S</u>) S'	2	(1 + 2 + (3 + 4)) + 5
\mapsto (1 + <u>E</u> S') S'	2	(1 + 2 + (3 + 4)) + 5
\mapsto (1 + 2 <u>S'</u>) S'	+	(1 + 2 + (3 + 4)) + 5
\mapsto (1 + 2 + <u>S</u>) S'	((1 + 2 + (3 + 4)) + 5
\mapsto (1 + 2 + <u>E</u> S') S'	((1 + 2 + (3 + 4)) + 5
\mapsto (1 + 2 + (<u>S</u>)S') S'	3	(1 + 2 + (3 + 4)) + 5

Predictive Parsing

- Given an LL(1) grammar:
 - For a given nonterminal, the lookahead symbol uniquely determines the production to apply.
 - Top-down parsing = predictive parsing
 - Driven by a predictive parsing table:
nonterminal * input token \rightarrow production

$$\begin{aligned} T &\mapsto S\$ \\ S &\mapsto ES' \\ S' &\mapsto \epsilon \\ S' &\mapsto + S \\ E &\mapsto \text{number} \mid (S) \end{aligned}$$

	number	+	()	\$ (EOF)
T	$\mapsto S\$$		$\mapsto S\$$		
S	$\mapsto E S'$		$\mapsto E S'$		
S'		$\mapsto + S$		$\mapsto \epsilon$	$\mapsto \epsilon$
E	$\mapsto \text{num.}$		$\mapsto (S)$		

- Note: it is convenient to add a special *end-of-file* token \$ and a start symbol T (top-level) that requires \$.


How do we construct the parse table?

- Consider a given production: $A \rightarrow \gamma$
- Construct the set of all input tokens that may appear *first* in strings that can be derived from γ
 - Add the production $\rightarrow \gamma$ to the entry (A,token) for each such token.
- If γ can derive ϵ (the empty string), then we construct the set of all input tokens that may *follow* the nonterminal A in the grammar.
 - Add the production $\rightarrow \gamma$ to the entry (A, token) for each such token.
- Note: if there are two different productions for a given entry, the grammar is not LL(1)

Example

- $\text{First}(T) = \text{First}(S)$
- $\text{First}(S) = \text{First}(E)$
- $\text{First}(S') = \{ + \}$
- $\text{First}(E) = \{ \text{number}, '(' \}$
- $\text{Follow}(S') = \text{Follow}(S)$
- $\text{Follow}(S) = \{ \$, ')' \} \cup \text{Follow}(S')$

$T \mapsto S\$$
 $S \mapsto ES'$
 $S' \mapsto \epsilon$
 $S' \mapsto + S$
 $E \mapsto \text{number} \mid (S)$



	number	+	()	\$ (EOF)
T	$\mapsto S\$$		$\mapsto S\$$		
S	$\mapsto E S'$		$\mapsto E S'$		
S'		$\mapsto + S$		$\mapsto \epsilon$	$\mapsto \epsilon$
E	$\mapsto \text{num.}$		$\mapsto (S)$		

Converting the table to code

- Define n mutually recursive functions
 - one for each nonterminal A : `parse_A`
 - The type of `parse_A` is `unit -> ast` if A is *not* an auxiliary nonterminal
 - Parse functions for auxiliary nonterminals (e.g. S') take extra `ast`'s as inputs, one for each nonterminal in the “factored” prefix.
- Each function “peeks” at the lookahead token and then follows the production rule in the corresponding entry.
 - Consume terminal tokens from the input stream
 - Call `parse_X` to create sub-tree for nonterminal X
 - If the rule ends in an auxiliary nonterminal, call it with appropriate `ast`'s. (The auxiliary rule is responsible for creating the `ast` after looking at more input.)
 - Otherwise, this function builds the `ast` tree itself and returns it.

	number	+	()	\$ (EOF)
T	$\mapsto S\$$		$\mapsto S\$$		
S	$\mapsto E S'$		$\mapsto E S'$		
S'		$\mapsto + S$		$\mapsto \epsilon$	$\mapsto \epsilon$
E	$\mapsto \text{num.}$		$\mapsto (S)$		

Hand-generated LL(1) code for the table above.

DEMO: PARSER.ML

LL(1) Summary

- Top-down parsing that finds the leftmost derivation.
- Language Grammar \Rightarrow LL(1) grammar \Rightarrow prediction table \Rightarrow recursive-descent parser
- Problems:
 - Grammar must be LL(1)
 - Can extend to LL(k) (it just makes the table bigger)
 - Grammar cannot be left recursive (parser functions will loop!)
- Is there a better way?

LR GRAMMARS

Bottom-up Parsing (LR Parsers)

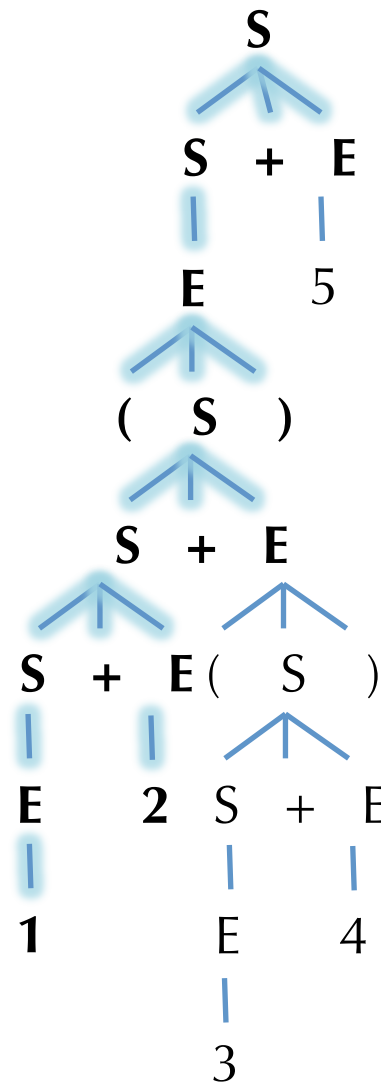
- LR(k) parser:
 - Left-to-right scanning
 - Rightmost derivation
 - k lookahead symbols
- LR grammars are more expressive than LL
 - Can handle left-recursive (and right recursive) grammars; virtually all programming languages
 - Easier to express programming language syntax (no left factoring)
- Technique: “Shift-Reduce” parsers
 - Work bottom up instead of top down
 - Construct right-most derivation of a program in the grammar
 - Used by many parser generators (e.g. yacc, CUP, ocamllyacc, etc.)
 - Better error detection/recovery

Top-down vs. Bottom up

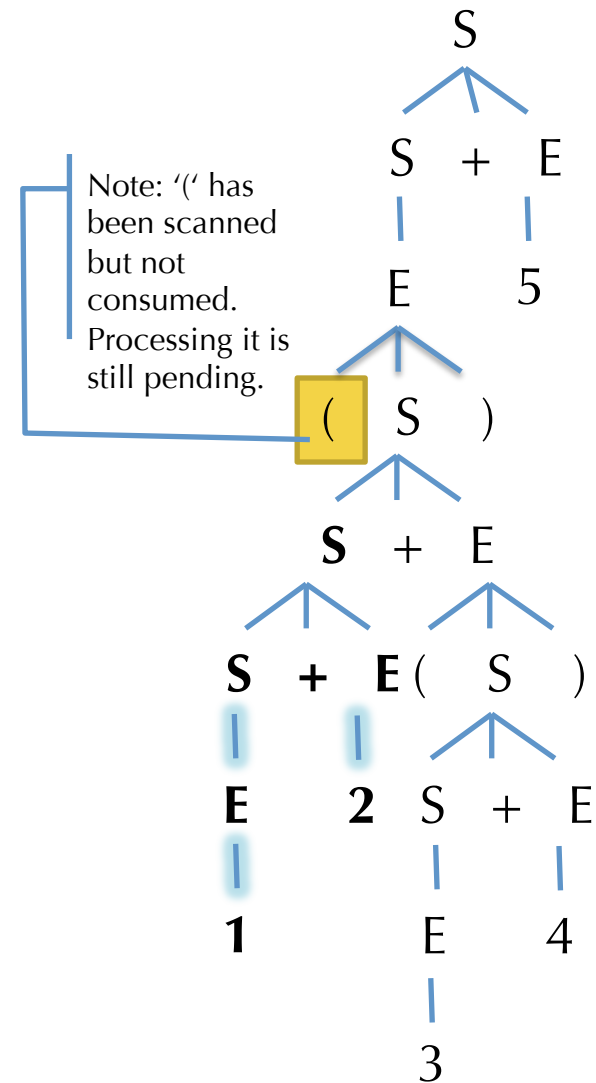
- Consider the left-recursive grammar:

$S \mapsto S + E \mid E$
 $E \mapsto \text{number} \mid (S)$

- $(1 + 2 + (3 + 4)) + 5$
- What part of the tree must we know after scanning just $(1 + 2$
- In top-down, must be able to guess which productions to use...




Top-down



Bottom-up

Progress of Bottom-up Parsing

	Reductions	Scanned	Input Remaining
 Rightmost derivation	$(1 + 2 + (3 + 4)) + 5 \leftarrow$		$(1 + 2 + (3 + 4)) + 5$
	$(\underline{\mathbf{E}} + 2 + (3 + 4)) + 5 \leftarrow$	$($	$+ 2 + (3 + 4)) + 5$
	$(\underline{\mathbf{S}} + 2 + (3 + 4)) + 5 \leftarrow$	$(1$	$+ 2 + (3 + 4)) + 5$
	$(\mathbf{S} + \underline{\mathbf{E}} + (3 + 4)) + 5 \leftarrow$	$(1 + 2$	$+ (3 + 4)) + 5$
	$(\underline{\mathbf{S}} + (3 + 4)) + 5 \leftarrow$	$(1 + 2$	$+ (3 + 4)) + 5$
	$(\mathbf{S} + (\underline{\mathbf{E}} + 4)) + 5 \leftarrow$	$(1 + 2 + (3$	$+ 4)) + 5$
	$(\mathbf{S} + (\underline{\mathbf{S}} + 4)) + 5 \leftarrow$	$(1 + 2 + (3$	$+ 4)) + 5$
	$(\mathbf{S} + (\mathbf{S} + \underline{\mathbf{E}})) + 5 \leftarrow$	$(1 + 2 + (3 + 4$	$)) + 5$
	$(\mathbf{S} + (\underline{\mathbf{S}})) + 5 \leftarrow$	$(1 + 2 + (3 + 4$	$)) + 5$
	$(\mathbf{S} + \underline{\mathbf{E}}) + 5 \leftarrow$	$(1 + 2 + (3 + 4)$	$) + 5$
	$(\underline{\mathbf{S}}) + 5 \leftarrow$	$(1 + 2 + (3 + 4)$	$) + 5$
	$\underline{\mathbf{E}} + 5 \leftarrow$	$(1 + 2 + (3 + 4))$	$+ 5$
	$\underline{\mathbf{S}} + 5 \leftarrow$	$(1 + 2 + (3 + 4))$	$+ 5$
	$\mathbf{S} + \underline{\mathbf{E}} \leftarrow$	$(1 + 2 + (3 + 4)) + 5$	
	\mathbf{S}		

$S \mapsto S + E \mid E$
 $E \mapsto \text{number} \mid (S)$

Shift/Reduce Parsing

- Parser state:
 - Stack of terminals and nonterminals.
 - Unconsumed input is a string of terminals
 - Current derivation step is $\text{stack} + \text{input}$
- Parsing is a sequence of *shift* and *reduce* operations:
- Shift**: move look-ahead token to the stack
- Reduce**: Replace symbols γ at top of stack with nonterminal X such that $X \mapsto \gamma$ is a production. (pop γ , push X)

$$S \mapsto S + E \mid E$$

$$E \mapsto \text{number} \mid (S)$$

Stack	Input	Action
	(1 + 2 + (3 + 4)) + 5	shift (
(1 + 2 + (3 + 4)) + 5	shift 1
(1	+ 2 + (3 + 4)) + 5	reduce: $E \mapsto \text{number}$
(E	+ 2 + (3 + 4)) + 5	reduce: $S \mapsto E$
(S	+ 2 + (3 + 4)) + 5	shift +
(S +	2 + (3 + 4)) + 5	shift 2
(S + 2	+ (3 + 4)) + 5	reduce: $E \mapsto \text{number}$

Simple LR parsing with no look ahead.

LR(0) GRAMMARS

LR Parser States

- Goal: know what set of reductions are legal at any given point.
- Idea: Summarize all possible stack prefixes α as a finite parser state.
 - Parser state is computed by a DFA that reads the stack σ .
 - Accept states of the DFA correspond to unique reductions that apply.
- Example: LR(0) parsing
 - Left-to-right scanning, Right-most derivation, zero look-ahead tokens
 - Too weak to handle many language grammars (e.g. the “sum” grammar)
 - But, helpful for understanding how the shift-reduce parser works.

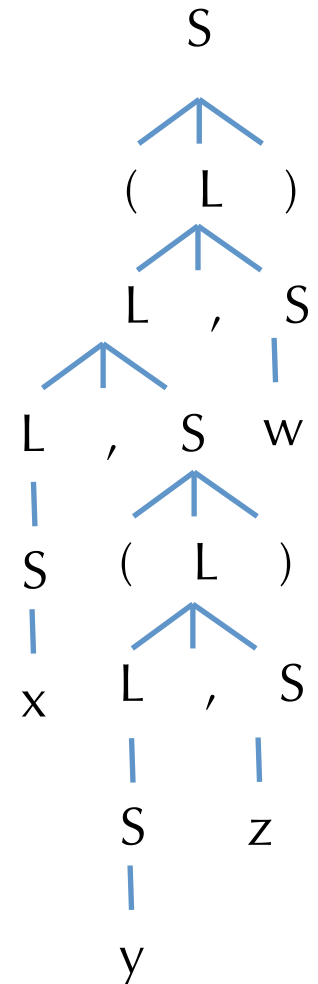
Example LR(0) Grammar: Tuples

- Example grammar for non-empty tuples and identifiers:

$$\begin{array}{lcl} S & \longmapsto & (L) \mid \text{id} \\ L & \longmapsto & S \mid L, S \end{array}$$

- Example strings:
 - x
 - (x, y)
 - $((((x))))$
 - $(x, (y, z), w)$
 - $(x, (y, (z, w)))$

Parse tree for:
(x, (y, z), w)



Shift/Reduce Parsing

- Parser state:
 - Stack of terminals and nonterminals.
 - Unconsumed input is a string of terminals
 - Current derivation step is $\text{stack} + \text{input}$
- Parsing is a sequence of *shift* and *reduce* operations:
- Shift**: move look-ahead token to the stack: e.g.

$$\begin{array}{l} S \mapsto (L) \mid \text{id} \\ L \mapsto S \mid L , S \end{array}$$

Stack	Input	Action
	(x, (y, z), w)	shift (
(x, (y, z), w)	shift x

- Reduce**: Replace symbols γ at top of stack with nonterminal X such that $X \mapsto \gamma$ is a production. (pop γ , push X): e.g.

Stack	Input	Action
(x	, (y, z), w)	reduce $S \mapsto \text{id}$
(S	, (y, z), w)	reduce $L \mapsto S$

Example Run

Stack	Input	Action
	(x, (y, z), w)	shift (
(x, (y, z), w)	shift x
(x	, (y, z), w)	reduce $S \mapsto \text{id}$
(S	, (y, z), w)	reduce $L \mapsto S$
(L	, (y, z), w)	shift ,
(L,	(y, z), w)	shift (
(L, (y, z), w)	shift y
(L, (y	, z), w)	reduce $S \mapsto \text{id}$
(L, (S	, z), w)	reduce $L \mapsto S$
(L, (L	, z), w)	shift ,
(L, (L,	z), w)	shift z
(L, (L, z), w)	reduce $S \mapsto \text{id}$
(L, (L, S), w)	reduce $L \mapsto L, S$
(L, (L), w)	shift)
(L, (L)	, w)	reduce $S \mapsto (L)$
(L, S	, w)	reduce $L \mapsto L, S$
(L, S	, w)	shift ,
(L, S,	w)	shift w
(L, S, w)	reduce $S \mapsto \text{id}$

$S \mapsto (L) \mid \text{id}$
 $L \mapsto S \mid L, S$

Action Selection Problem

- Given a stack σ and a look-ahead symbol b , should the parser:
 - Shift b onto the stack (new stack is σb)
 - Reduce a production $X \mapsto \gamma$, assuming that $\sigma = \alpha\gamma$ (new stack is αX)?
- Sometimes the parser can reduce but shouldn't
 - For example, $X \mapsto \epsilon$ can *always* be reduced
- Sometimes the stack can be reduced in different ways
- Main idea: decide what to do based on a *prefix* α of the stack plus the look-ahead symbol.
 - The prefix α is different for different possible reductions since in productions $X \mapsto g$ and $Y \mapsto b$, g and b might have different lengths.
- Main goal: know what set of reductions are legal at any point.
 - How do we keep track?

LR(0) States

- An LR(0) *state* is a set of *items* keeping track of progress on possible upcoming reductions.
- An LR(0) *item* is a production from the language with an extra separator “.” somewhere in the right-hand-side

$$\begin{array}{l} S \mapsto (L) \mid id \\ L \mapsto S \mid L , S \end{array}$$


- Example items: $S \mapsto . (L)$ or $S \mapsto (. L)$ or $L \mapsto S .$
- Intuition:
 - Stuff before the ‘.’ is already on the stack (beginnings of possible γ 's to be reduced)
 - Stuff after the ‘.’ is what might be seen next
 - The prefixes α are represented by the state itself

Constructing the DFA: Start state & Closure

- First step: Add a new production $S' \mapsto S\$$ to the grammar
- Start state of the DFA = empty stack, so it contains the item:
 $S' \mapsto .S\$$
- Closure of a state:
 - Adds items for all productions whose LHS nonterminal occurs in an item in the state just after the $'.'$
 - The added items have the $'.'$ located at the beginning (no symbols for those items have been added to the stack yet)
 - Note that newly added items may cause yet more items to be added to the state... keep iterating until a *fixed point* is reached.
- Example: $\text{CLOSURE}(\{S' \mapsto .S\$\}) = \{S' \mapsto .S\$, S \mapsto .(L), S \mapsto .id\}$
- Resulting “closed state” contains the set of all possible productions that might be reduced next.

$S' \mapsto S\$$
$S \mapsto (L) \mid id$
$L \mapsto S \mid L , S$

Example: Constructing the DFA


 $S' \mapsto .S\$$

$$\begin{array}{l} S' \mapsto S\$ \\ S \mapsto (L) \mid \text{id} \\ L \mapsto S \mid L , S \end{array}$$

- First, we construct a state with the initial item $S' \mapsto .S\$$

Example: Constructing the DFA

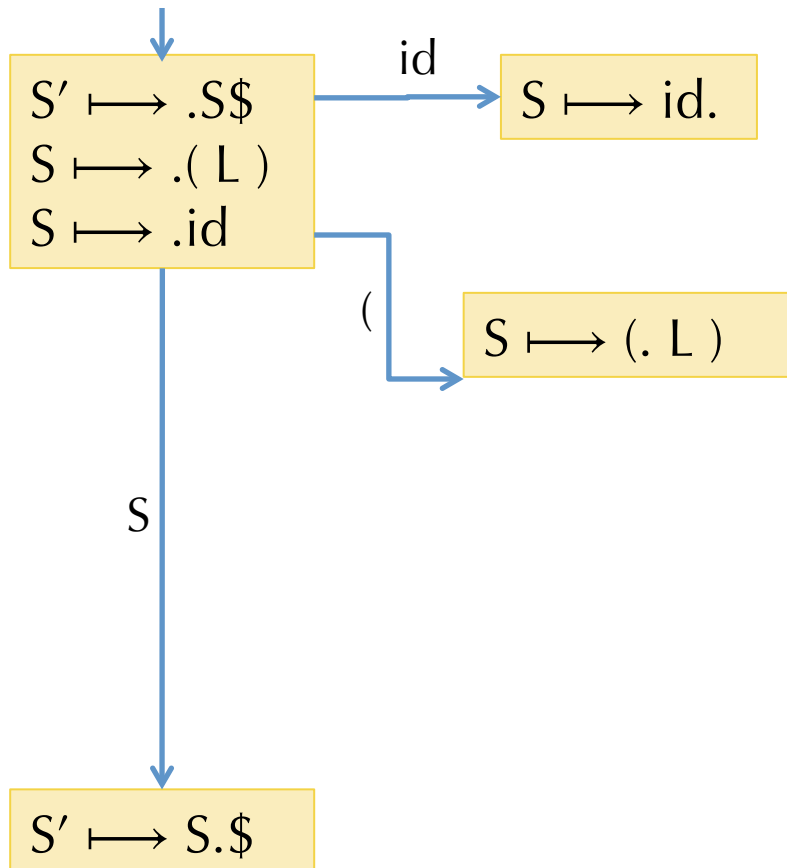
↓
 $S' \mapsto .S\$$
 $S \mapsto .(L)$
 $S \mapsto .id$

$S' \mapsto S\$$
 $S \mapsto (L) \mid id$
 $L \mapsto S \mid L , S$

- Next, we take the closure of that state:
 $CLOSURE(\{S' \mapsto .S\$\}) = \{S' \mapsto .S\$, S \mapsto .(L), S \mapsto .id\}$
- In the set of items, the nonterminal S appears after the $'.'$
- So we add items for each S production in the grammar

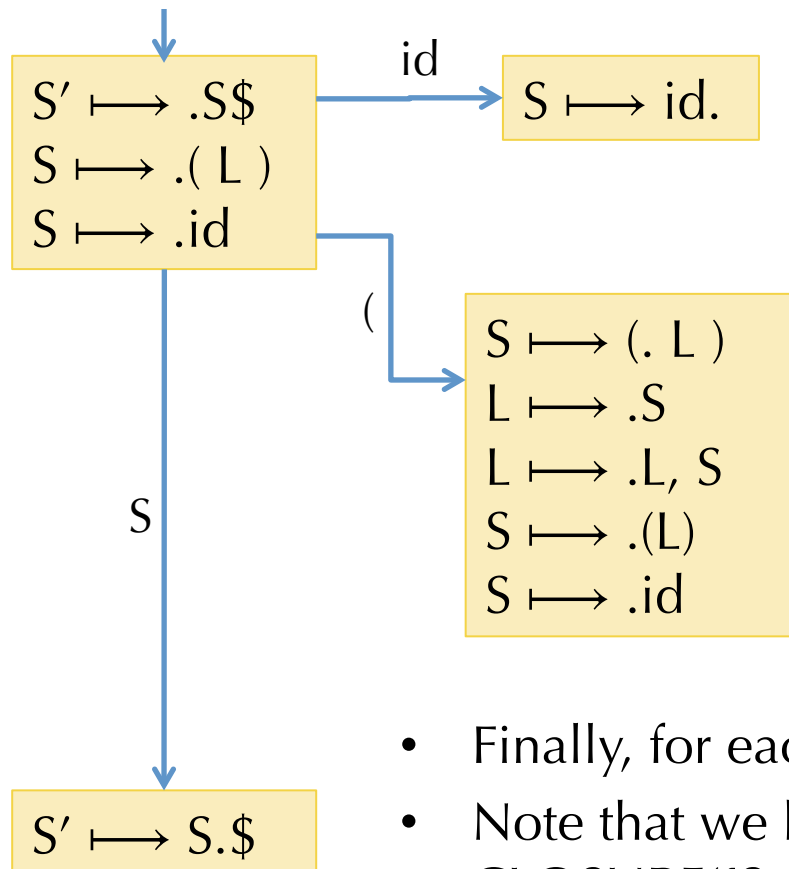
Example: Constructing the DFA

$S' \mapsto S\$$
 $S \mapsto (L) \mid id$
 $L \mapsto S \mid L , S$



- Next we add the transitions:
- First, we see what terminals and nonterminals can appear after the '.' in the source state.
 - Outgoing edges have those label.
- The target state (initially) includes all items from the source state that have the edge-label symbol after the '.', but we advance the '.' (to simulate shifting the item onto the stack)

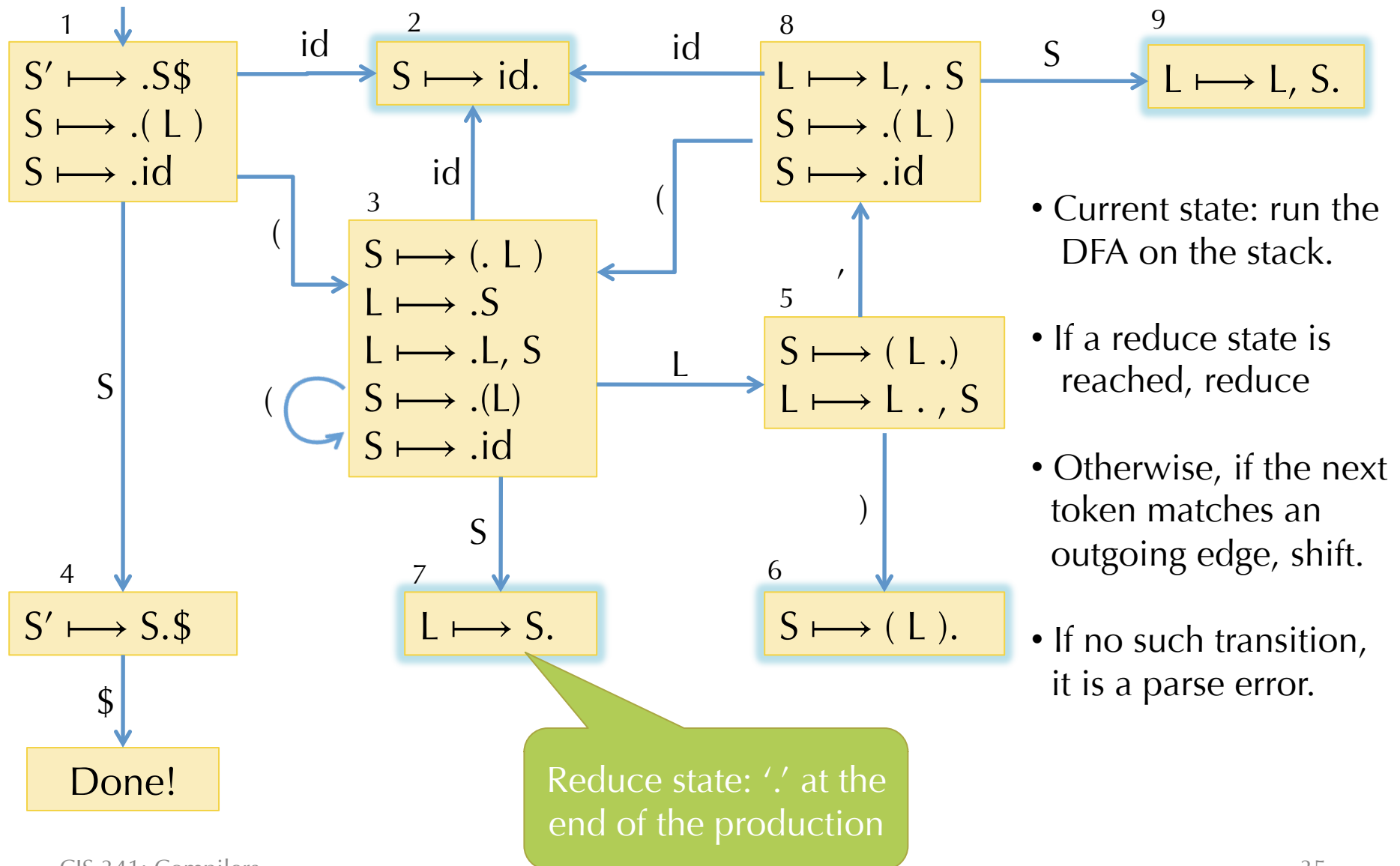
Example: Constructing the DFA



$S' \mapsto S\$$
 $S \mapsto (L) \mid id$
 $L \mapsto S \mid L, S$

- Finally, for each new state, we take the closure.
- Note that we have to perform two iterations to compute $CLOSURE(\{S \mapsto (.L)\})$
 - First iteration adds $L \mapsto .S$ and $L \mapsto .L, S$
 - Second iteration adds $S \mapsto .(L)$ and $S \mapsto .id$

Full DFA for the Example



- Current state: run the DFA on the stack.
- If a reduce state is reached, reduce
- Otherwise, if the next token matches an outgoing edge, shift.
- If no such transition, it is a parse error.