## **EECS 483: Compiler Construction** Lecture 24: LR Grammars, Bottom-Up Parsing

April 14 Winter Semester 2025







# Announcements

- Reminder: Assignment 5 Due Sunday
- Midterm regrades done. Grades updated on Canvas
- Course Evaluations are now open. Please fill them out!

# updated on Canvas en. Please fill them out!

# **Final Exam Info**

- Final Exam on Wednesday 04/30 4-6pm
  - Location:

DOW 1010 (uniquame starts with A-L) DOW 1014 (uniquame starts with M-Z)

- Topics: Assignments 4 and 5, lecture material after spring break.
- Exam Review on Monday, 04/21
- 1 page of notes, double sided ok, printed or written ok.
- Practice material:

https://maxsnew.com/teaching/eecs-483-wn24/syllabus.html questions about lexing/parsing/analysis/optimization

- Appel book, Dragon book linked on webpage have exercises as well.

#### LR GRAMMARS



4

# **The Parsing Problem**

- The Parsing Problem: ullet
  - Input: a context-free grammar G
  - Output: a **parser** that takes in a string and outputs a parse tree of that string in G or raises an exception if there is no parse tree.
  - Notice that an ambiguous grammar may be parsed in multiple ways
- In practice: fuse the generation of the parse tree with semantic actions that construct the abstract syntax tree
  - The parse tree is usually never "materialized" in memory
- Another "mini-compiler" for a DSL lacksquare
- Bad news: best algorithms are O(n^3)
  - CYK, Earley, GLR algorithms
- Compromise: find restrictions on CFGs that allow for O(n) parsing  $\bullet$ Intuition: parsing is a *search problem*, find restrictions that limit the amount
- of backtracking needed.
  - Cost: more burden on the programmer (i.e., **you**) to adapt their grammar to fit the restriction

- Top-down parsing that finds the leftmost derivation.  $\bullet$
- Language Grammar  $\bullet$ 
  - $\Rightarrow$  LL(1) grammar (manual rewrite)
  - ⇒ prediction table (intermediate representation)
  - $\Rightarrow$  recursive-descent parser (code generation)
- **Problems:** ullet
  - 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?  $\bullet$ 

#### LL(1) Summary

### **Bottom-up Parsing (LR Parsers)**

- LR(k) parser:
  - Left-to-right scanning
  - <u>Rightmost derivation</u>
  - k lookahead symbols
- LR grammars are more expressive than LL
   Can handle left-recursive (and right recursive) grammars; virtually all
  - Can handle left-recursive (and ri 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, ocamlyacc, lalrpop, etc.)
  - Better error detection/recovery

E

Consider the left- $\bullet$ recursive grammar:

> $S \mapsto S + E \mid E$  $E \mapsto number \mid (S)$

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

#### **Top-down vs. Bottom up**



### **Progress of Bottom-up Parsing**

	Reductions	Scanned	Input Remaining
1	$(1 + 2 + (3 + 4)) + 5 \leftarrow 4$		(1 + 2 + (3 + 4)) + 5
	$(\underline{\mathbf{E}} + 2 + (3 + 4)) + 5 \leftarrow \mathbf{I}$	(	1 + 2 + (3 + 4)) + 5
	$(\underline{\mathbf{S}} + 2 + (3 + 4)) + 5 \leftarrow 1$	(1	+2+(3+4))+5
	$(S + \underline{E} + (3 + 4)) + 5 \longleftarrow$	(1 + 2	+(3+4))+5
	$(\underline{\mathbf{S}} + (3 + 4)) + 5 \leftarrow 1$	(1 + 2)	+(3+4))+5
Ι<α	$(S + (\underline{E} + 4)) + 5 \longleftarrow$	(1 + 2 + (3 + 3))	+ 4)) + 5
סט	$(S + (\underline{S} + 4)) + 5 \longleftarrow$	(1 + 2 + (3 + 3))	(+ 4)) + 5
U C C C C	$(S + (S + \underline{E})) + 5 \longleftarrow$	(1 + 2 + (3 + 4))	)) + 5
	$(S + (\underline{S})) + 5 \leftarrow \mathbf{I}$	(1 + 2 + (3 + 4))	)) + 5
	$(S + \underline{E}) + 5 \leftarrow \cdot$	(1 + 2 + (3 + 4))	) + 5
	$(\underline{\mathbf{S}}) + 5 \longleftarrow$	(1 + 2 + (3 + 4))	) + 5
	<u><b>E</b></u> + 5 ← · ·	(1 + 2 + (3 + 4))	+ 5
	$\underline{\mathbf{S}}$ + 5 $\leftarrow$	(1 + 2 + (3 + 4))	+ 5
	S + <u>E</u> ← →	(1 + 2 + (3 + 4)) + 5	
	S		

**Rightmost derivation** 

# $S \mapsto S + E \mid E$ $E \mapsto number \mid (S)$

### Shift/Reduce Parsing

• Parser state:

(1

(E

(S

(S

(S +

(S + 2)

(S + E)

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

(1 + 2 + (3 + 4)) +1 + 2 + (3 + 4)) +

- +2+(3+4))+
  - +2+(3+4))+
  - +2+(3+4))+
    - 2 + (3 + 4)) +
      - +(3+4))+
      - + (3 + 4)) +

+(3+4))+

 $S \mapsto S + E \mid E$  $E \mapsto number \mid (S)$ 

	Action
5	shift (
5	shift 1
5	reduce: $E \mapsto number$
5	reduce: $S \mapsto E$
5	shift +
5	shift 2
5	reduce: $E \mapsto number$
5	reduce: $S \mapsto S + E$
5	shift +

#### **Shift/Reduce Parsing**

- Parser state:  ${\bullet}$ 
  - Stack of terminals and nonterminals.
  - Unconsumed input is a string of terminals
  - Current derivation step is stack + input
- $\bullet$ in reverse

Stack	Input
	(1 + 2 + (3 +
(	1 + 2 + (3 + 3)
(1	+ 2 + (3 +
(E	+ 2 + (3 +
(S	+ 2 + (3 +
(S +	2 + (3 +
(S + 2	+ (3 +
(S + E	+ (3 +

(S

 $S \mapsto S + E \mid E$  $E \mapsto number \mid (S)$ 

Invariant: Stack plus input is a step in building the Rightmost derivation

Derivation steps (1 + 2 + (3 + 4)) + 5+ 4)) + 5 (+ 4)) + 5(+ 4)) + 5 $(\underline{E} + 2 + (3 + 4)) + 5$ ( $\underline{E} + 2 + (3 + 4)) + 5$ ( $\underline{S} + 2 + (3 + 4)) + 5$ derivation Rightmost (+ 4)) + 5(4)) + 5 $(S + \underline{E} + (3 + 4)) + 5$ (+ 4)) + 5 $(\underline{S} + (3 + 4)) + 5$ +(3+4))+5

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  $\alpha$  as a finite parser state. – Parser state is computed by a DFA that reads the stack  $\sigma$ . – Accept states of the DFA correspond to unique reductions that apply.
- lacksquare $\bullet$

- Example: LR(0) parsing  $\bullet$ 
  - Too weak to handle many language grammars (e.g. the "sum" grammar)

  - <u>Left-to-right scanning</u>, <u>**R**ight-most derivation</u>, <u>**zero** look-ahead tokens</u> \_\_\_\_\_ – But, helpful for understanding how the shift-reduce parser works.

### **Example LR(0) Grammar: Tuples**

- Example grammar for non-empty tuples and identifiers: •

- Example strings: •
  - X
  - (x,y)
  - ((((x))))
  - (x, (y, z), w)
  - (x, (y, (z, w)))

# $S \mapsto (L) | id$ $L \mapsto S | L, S$

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



### **Shift/Reduce Parsing**

Parser state:

(S

- Stack of terminals and nonterminals.
- Unconsumed input is a string of terminals
- Current derivation step is stack + input
- Parsing is a sequence of *shift* and *reduce* operations:

Shift: move look-ahead toker		
Stack	Inp	out
	(x,	(y, z), w)
(	Х,	(y, z), w)

Reduce: Replace symbols  $\gamma$  at top of stack with nonterminal X such that  $X \mapsto \gamma$  is a production. (pop  $\gamma$ , push X): e.g. Stack Action Input , (y, z), w) reduce  $S \mapsto id$  $(\mathbf{X})$ 

, (y, z), w)

to the stack: e.g.

_	
	Action
	shift (
	shift x

reduce  $L \mapsto S$ 

 $\begin{array}{c|c} S \longmapsto (L) & | id \\ L \longmapsto S & | L, S \end{array}$ 

Stack	Input	
	(x, (y, z), w)	
(	x, (y, z), w)	
(x	, (y, z), w)	
(S	, (y, z), w)	
(L	, (y, z), w)	
(L,	(y, z), w)	
(L, (	y, z), w)	
(L, (y	, z), w)	
(L, (S	, z), w)	
(L, (L	, z), w)	
(L, (L,	z), w)	
(L, (L, z	), w)	
(L, (L, S	), w)	
(L, (L	), W)	
(L, (L)	, W)	
(L, S	, w)	
(L	, w)	
(L,	W)	
(L, w	)	
(L, S	)	
(L	)	
(L)		
S		

# **Example Run**

- Action
- shift (
- shift x
- reduce  $S \mapsto id$
- reduce  $L \mapsto S$
- shift ,
- shift (
- shift y
- reduce  $S \mapsto id$
- reduce  $L \mapsto S$
- shift ,
- shift z
- reduce  $S \mapsto id$
- reduce  $L \mapsto L$ , S
- shift )
- reduce  $S \mapsto (L)$
- reduce  $L \mapsto L$ , S
- shift ,
- shift w
- reduce  $S \mapsto id$
- reduce  $L \mapsto L$ , S
- shift)
- reduce  $S \mapsto (L)$

$$S \mapsto (L) | id$$
$$L \mapsto S | L, S$$

### **Action Selection Problem**

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



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

$$\begin{array}{c} S \longmapsto (L) \\ L \longmapsto S \end{array} \right|$$

- Example items:  $S \mapsto .(L)$  or  $S \mapsto (.L)$  or  $L \mapsto S$ .  $\bullet$
- Intuition:  ${\color{black}\bullet}$ 
  - Stuff before the '.' is already on the stack (beginnings of possible  $\gamma$ 's to be reduced)
  - Stuff after the '.' is what might be seen next
  - The prefixes  $\alpha$  are represented by the state itself

#### LR(0) States

| id L , S

#### **Constructing the DFA: Start state & Closure**

- First step: Add a new production  $\bullet$  $S' \mapsto S$  to the grammar
- Start state of the DFA = empty stack, ulletso it contains the item:  $S' \mapsto .S\$$
- Closure of a state:
  - the state just after the '.'
  - items have been added to the stack yet)
  - state... keep iterating until a *fixed point* is reached.
- might be reduced next.

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

– Adds items for all productions whose LHS nonterminal occurs in an item in

– The added items have the '.' located at the beginning (no symbols for those

– Note that newly added items may cause yet more items to be added to the

Example:  $CLOSURE({S' \mapsto .S}) = {S' \mapsto .S}, S \mapsto .(L), S \mapsto .id}$ 

Resulting "closed state" contains the set of all possible productions that

$$S' \mapsto .S$$

First, we construct a state with the initial item  $S' \rightarrow .S$ lacksquare

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

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

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

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



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

- 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)



- $CLOSURE(\{S \mapsto ( \ . \ L \ )\})$

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

Finally, for each new state, we take the closure.

Note that we have to perform two iterations to compute

– 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

id

8

5

 $L \mapsto L, . S$ 

 $S \mapsto .(L)$ 

 $S \mapsto (L)$ 

 $L \mapsto L . , S$ 

 $S \mapsto .id$ 



) 6  $S \mapsto (L).$ Reduce state: '.' at the end of the production • Current state: run the DFA on the stack.

 $L \mapsto L, S.$ 

Q

S

- 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.

# Using the DFA

- Run the parser stack through the DFA. • The resulting state tells us which productions might be
- reduced next.
  - If not in a reduce state, then shift the next symbol and transition according to DFA.
  - If in a reduce state,  $X \mapsto \gamma$  with stack  $\alpha \gamma$ , pop  $\gamma$  and push X.
- Optimization: No need to re-run the DFA from beginning  $\bullet$ every step
  - Store the state with each symbol on the stack: e.g.  $_1(_3(_3L_5)_6)$
  - On a reduction  $X \mapsto \gamma$ , pop stack to reveal the state too: e.g. From stack  $_1(_3(_3L_5)_6)$  reduce  $S \mapsto (L)$  to reach stack  $_1(_3)$
  - Next, push the reduction symbol: e.g. to reach stack  $_1(_3S)$
  - Then take just one step in the DFA to find next state:  $_1(_3S_7)$

# **Implementing the Parsing Table**

Represent the DFA as a table of shape: state \* (terminals + nonterminals)

- Entries for the "action table" specify two kinds of actions:  $\bullet$ 
  - Shift and goto state n
  - Reduce using reduction  $X \mapsto \gamma$ 
    - First pop  $\gamma$  off the stack to reveal the state
    - Look up X in the "goto table" and goto that state ●





Nonterminal Symbols

Action table

#### **Example Parse Table**

	(	)	id	,	\$	S	L
1	s3		s2			g4	
2	S⊷id	S⊷id	S⊷id	S⊷id	S⊷id		
3	<b>s</b> 3		s2			g7	g5
4					DONE		
5		s6		s8			
6	$S \mapsto (L)$						
7	$L \mapsto S$						
8	<b>s</b> 3		s2			g9	
9	L → L,S	$L \mapsto L,S$	$L \mapsto L,S$	L → L,S	$L \mapsto L,S$		

sx = shift and goto state x gx = goto state x (used when we reduce)



• Parse the toke	en stream: (x, (y,
Stack	Stream
$\epsilon_1$	(x, (y, z), w)\$
$\epsilon_1(3)$	x, (y, z), w)\$
$\epsilon_1(_3x_2$	, (y, z), w)\$
$\epsilon_1(_3S)$	, (y, z), w)\$
$\epsilon_1(_3S_7$	, (y, z), w)\$
$\epsilon_1(_3L)$	, (y, z), w)\$
$\epsilon_1(_3L_5$	, (y, z), w)\$
$\epsilon_1(_3L_5,_8)$	(y, z), w)\$
$\epsilon_1(_3L_5,_8(_3$	y, z), w)\$

#### Example

- z), w)\$
- Action (according to table)
- **s**3
- s2
- Reduce:  $S \mapsto id$
- g7 (from state 3 follow S)
- Reduce:  $L \mapsto S$
- g5 (from state 3 follow L)
- s8
- **s**3
- s2

# LR(0) Limitations

- An LR(0) machine only works if states with reduce actions have a *single* reduce action.
  - In such states, the machine *always* reduces (ignoring lookahead)
- With more complex grammars, the DFA construction will yield states with shift/reduce and reduce/reduce conflicts: shift/reduce reduce/reduce OK

$$S \mapsto (L).$$

- Such conflicts can often be resolved by using a look-ahead symbol: SLR(1) or LR(1)
- $S \mapsto (L).$  $L \mapsto .L, S$

$$S \mapsto L, S.$$
  
 $S \mapsto , S.$ 



Consider the left associative and right associative "sum" grammars: ullet



- One is LR(0) the other isn't... which is which and why? lacksquare
- What kind of conflict do you get? Shift/reduce or Reduce/reduce?  $\bullet$
- Ambiguities in associativity/precedence usually lead to shift/reduce ulletconflicts.

#### Examples

right  $S \mapsto E + S \mid E$  $E \mapsto number \mid (S)$ 



Consider the left associative and right associative "sum" grammars: ullet

left

$$S \mapsto S + E \mid E$$
$$E \mapsto number \mid (S)$$

- One is LR(0) the other isn't... which is which and why? ullet
- What kind of conflict do you get? Shift/reduce or Reduce/reduce? ullet

If the stack is a single E, then the state is



LR(0) parser can't decide

#### Examples

right

$$S \mapsto E + S \mid E$$
  
 $E \mapsto number \mid (S)$ 



shift-reduce conflict: we can either shift the + or reduce the E to an S.

# SLR(1) ("simple" LR) Parsers

- What conflicts are there in LR(0) parsing? ullet
  - reduce/reduce conflict: an LR(0) state has two reduce actions
  - shift/reduce conflict: an LR(0) state mixes reduce and shift actions
- Can we use lookahead to disambiguate?
- SLR(1) uses the same DFA construction as LR(0)ullet– modifies the actions based on lookahead. More powerful
- Suppose reducing an A nonterminal is possible in some state: •
  - compute Follow(A) for the given grammar

  - disjoint

– if the lookahead symbol is in Follow(A), then reduce, otherwise shift

- can disambiguate between reduce/reduce conflicts if the follow sets are

- Yet more powerful than SLR(1)
- Algorithm is similar to LR(0) DFA construction:
  - LR(1) state = set of LR(1) items
  - An LR(1) item is an LR(0) item + a set of look-ahead symbols:  $A \mapsto \alpha.\beta$  ,  $\mathcal{L}$
- LR(1) closure is a little more complex: ullet
- Form the set of items just as for LR(0) algorithm.
- Whenever a new item  $C \mapsto .\gamma$  is added because  $A \mapsto \beta.C\delta$ ,  $\mathcal{L}$  is • already in the set, we need to compute its look-ahead set  $\mathcal{M}$ :
  - 1. The look-ahead set  $\mathcal{M}$  includes FIRST( $\delta$ ) (the set of terminals that may start strings derived from  $\delta$ )
  - 2. If  $\delta$  is itself  $\epsilon$  or can derive  $\epsilon$  (i.e. it is nullable), then the look-ahead  $\mathcal{M}$  also contains  $\mathcal{L}$



### **Example Closure**

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

- Start item:  $S' \mapsto .S$ , {}
- Since S is to the right of a '.', add:  $S \mapsto .E + S , \{\$\}$  $S \mapsto .E$  , {\$}
- Need to keep closing, since E appears to the right of a '.' in '.E + S':

E →	.number	,	$\{+\}$
E →	.(S)	,	{+}

- Because E also appears to the right of '.' in '.E' we get: •  $E \mapsto .number$ , {\$}  $\mathsf{E} \longmapsto .(\mathsf{S}) \qquad , \quad \{\$\}$
- All items are distinct, so we're done •

 $E \mapsto number \mid (S)$ 

```
Note: {$} is FIRST($)
```

```
Note: + added for reason 1
    FIRST(+ S) = \{+\}
Note: $ added for reason 2
    δisε
```



- The behavior is determined if:  $\bullet$ 
  - There is no overlap among the look-ahead sets for each reduce item, and
  - None of the look-ahead symbols appear to the right of a '.'



	+	\$	E
1			g2
2	s3	$S \mapsto E$	

Fragment of the Action & Goto tables



- LR(1) gives maximal power out of a 1 look-ahead symbol parsing table lacksquare– DFA + stack is a push-down automaton
- In practice, LR(1) tables are big. lacksquare
  - Modern implementations (e.g., menhir) directly generate code
- LALR(1) = "Look-ahead LR"ullet
  - Merge any two LR(1) states whose items are identical except for the look-ahead sets:

$S' \mapsto .S\$$	{}
$S \mapsto .E + S$	{\$}
$S \mapsto .E\{\$\}$	
E → .num	{+}
$E \mapsto .(S)$	{+}
E ↦ .num	{\$}
$E \mapsto .(S)$	{\$}

- This is the usual technology for automatic parser generators: yacc, ocamlyacc
- Such merging can lead to nondeterminism (e.g., reduce/reduce conflicts), but – Results in a much smaller parse table and works well in practice \_\_\_\_\_
- GLR = "Generalized LR" parsing

  - Efficiently compute the set of *all* parses for a given input Later passes should disambiguate based on other context

#### LR variants



#### **Classification of Grammars**



Debugging parser conflicts. Disambiguating grammars.

# LALRPOP DEMO

