

Lecture 16+17

Bottom-Up Parsing

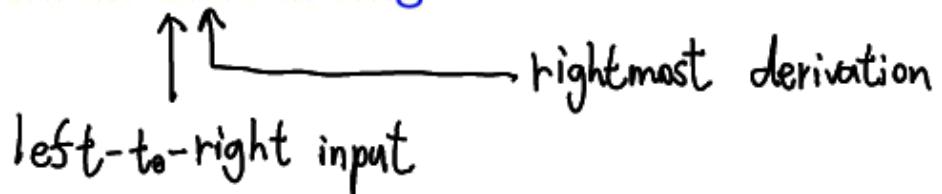
CS 241: Foundations of Sequential Programs
Fall 2014

Troy Vasiga et al
University of Waterloo

Example CFG

1. $S' \rightarrow \vdash S \dashv$
2. $S \rightarrow AyzB$
3. $A \rightarrow ab$
4. $A \rightarrow cd$
5. $B \rightarrow z$
6. $B \rightarrow wz$

Stacks in LR Parsing



- ▶ Recall that a stack in LL/top-down parsing is used in the following way:

$$\text{input processed} + \text{stack} = \text{current derivation}$$

(Note that the stack here is read from the top to bottom)

- ▶ For LR/bottom-up parsing, we have

$$\text{stack} + \text{input to be read} = \text{current derivation}$$

(stack is read from bottom to top here)

A trace

top of the stack

Derivation	Stack	Input read	Unread Input	Action
$\vdash \text{abywz} \dashv$	\vdash	\vdash	$\text{abywz} \dashv$	initialize (Shift \vdash)
$\vdash \text{abywz} \dashv$	$\vdash \text{a}$	$\vdash \text{a}$	$\text{bywz} \dashv$	Shift a
$\vdash \text{abywz} \dashv$	$\vdash \text{a b}$	$\vdash \text{ab}$	$\text{ywz} \dashv$	Shift b
$\vdash \text{Aywz} \dashv$	$\vdash \text{A}$	$\vdash \text{ab}$	$\text{ywz} \dashv$	Reduce $\text{A} \rightarrow \text{ab}$
$\vdash \text{Aywz} \dashv$	$\vdash \text{A y}$	$\vdash \text{aby}$	$\text{wz} \dashv$	Shift y
$\vdash \text{Aywz} \dashv$	$\vdash \text{A y w}$	$\vdash \text{abyw}$	$\text{z} \dashv$	Shift w
$\vdash \text{Aywz} \dashv$	$\vdash \text{A y w z}$	$\vdash \text{abywz}$	\dashv	Shift z
$\vdash \text{AyB} \dashv$	$\vdash \text{A y B}$	$\vdash \text{abywz}$	\dashv	Reduce $\text{B} \rightarrow \text{w z}$
$\vdash \text{S} \dashv$	$\vdash \text{S}$	$\vdash \text{abywz}$	\dashv	Reduce $\text{S} \rightarrow \text{AyB}$
$\vdash \text{S} \dashv$	$\vdash \text{S} \dashv$	$\vdash \text{abywz} \dashv$	ϵ	Shift \dashv

Shift: shifting a token from one place to another (push)

Reduce: size of the stack may be reduced (pop RHS, push LHS)

Shift/Reduce

- ▶ Somehow, we shifted at just the right time, and reduced just at the right time
- ▶ How did we know this?
 - ▶ Recall that for LL(1) parsing, we had a predictor table
 - ▶ For LR(1) parsing, we have an oracle, in the form of a DFA


transducer

Constructing DFA oracle for LR(1) grammars

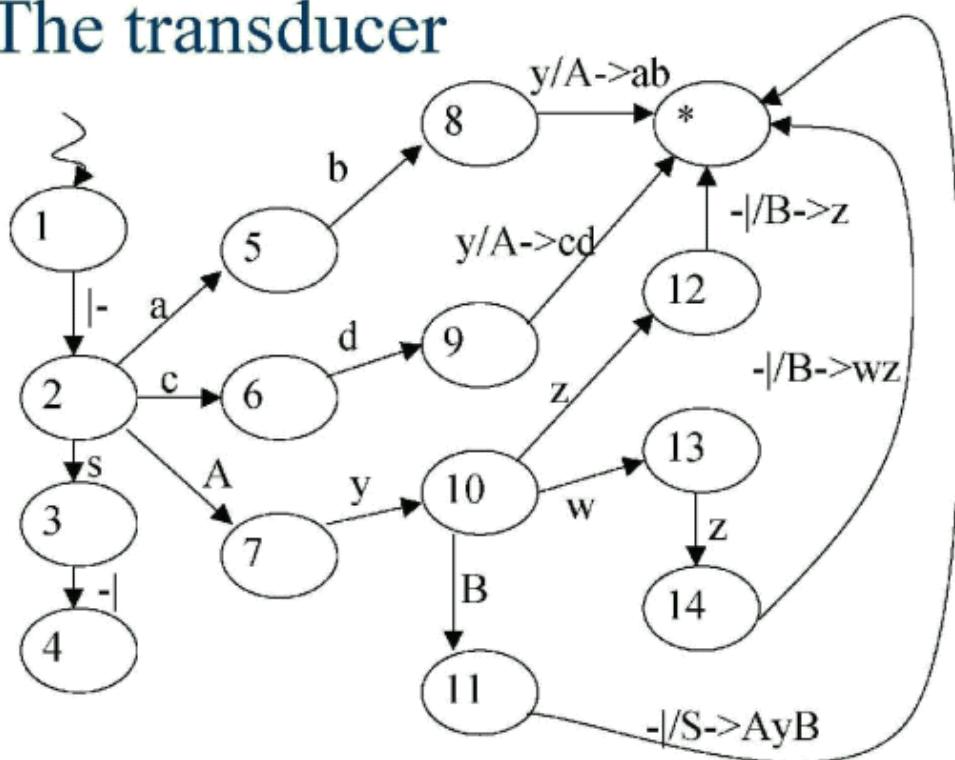


Take CS 462

- ▶ This is difficult to do
 - ▶ Donald Knuth proved a theorem that we can construct a DFA (really, a transducer) for LR(1) grammars (1965)
- ▶ You should know how to use the transducer

The transducer

The transducer



Understanding the transducer

Indicates whether to:

- ▶ shift (default)
- ▶ reduce (indicates rule)
- ▶ reject
 - ▶ All states except * and ERROR are final states
- ▶ Let's use this on input $\vdash \text{abywz} \dashv$



Using the transducer

Stack	States visited	Input read	Unread Input	Action
ϵ	1	\vdash	abywz \dashv	shift \vdash
\vdash	1 2	$\vdash a$	bywz \dashv	shift a
$\vdash a$	1 2 5	$\vdash ab$	ywz \dashv	shift b
$\vdash a b$	1 2 5 8	$\vdash ab$	ywz \dashv	Reduce A \rightarrow ab
$\vdash A$	1 2 7	$\vdash aby$	wz \dashv	shift y
$\vdash A y$	1 2 7 10	$\vdash abyw$	$z \dashv$	shift w
$\vdash A y w$	1 2 7 10 13	$\vdash abywz$	\dashv	shift z
$\vdash A y w z$	1 2 7 10 13 14	$\vdash abywz$	\dashv	Reduce B \rightarrow wz
$\vdash A y B$	1 2 7 10 11	$\vdash abywz$	\dashv	Reduce S \rightarrow AyB
$\vdash S$	1 2 3	$\vdash abywz$	\dashv	shift
$\vdash S \dashv$	1 2 3 4	$\vdash abywz \dashv$	ϵ	ACCEPT

↑
read stack
first

4 is not error
not ϵ

pop B
pop y
pop A
push S

A basic algorithm

- ▶ for each input token
 - ▶ Start in the start state
 - ▶ Read the stack (from the bottom up) and read the current input, and do the action indicated for the current input
 - ▶ If at any point, we hit the error state, reject input
- ▶ At the end, if we read $\vdash S \dashv$ on the stack, accept

Making this more efficient

Current running time of this algorithm:

Suppose N tokens: $O(n^2)$

Instead of scanning the stack each time...

push (input, state)

Start the transducer in....

the state at the top of the stack

Running time:

top: $O(1)$

Overall $\rightarrow O(N)$

Reduce $A \rightarrow ab$

(b, 8)	↓
(a, 5)	(A, 7)
(†, 2)	(†, 2)
(ε, 1)	(ε, 1)

Outputting a derivation

↳ 'proof'

- ▶ Easy: each time we do a reduction, output the rule
- ▶ But, this isn't quite right. Derivations should start with the start symbol. Bottom-up parsing doesn't.

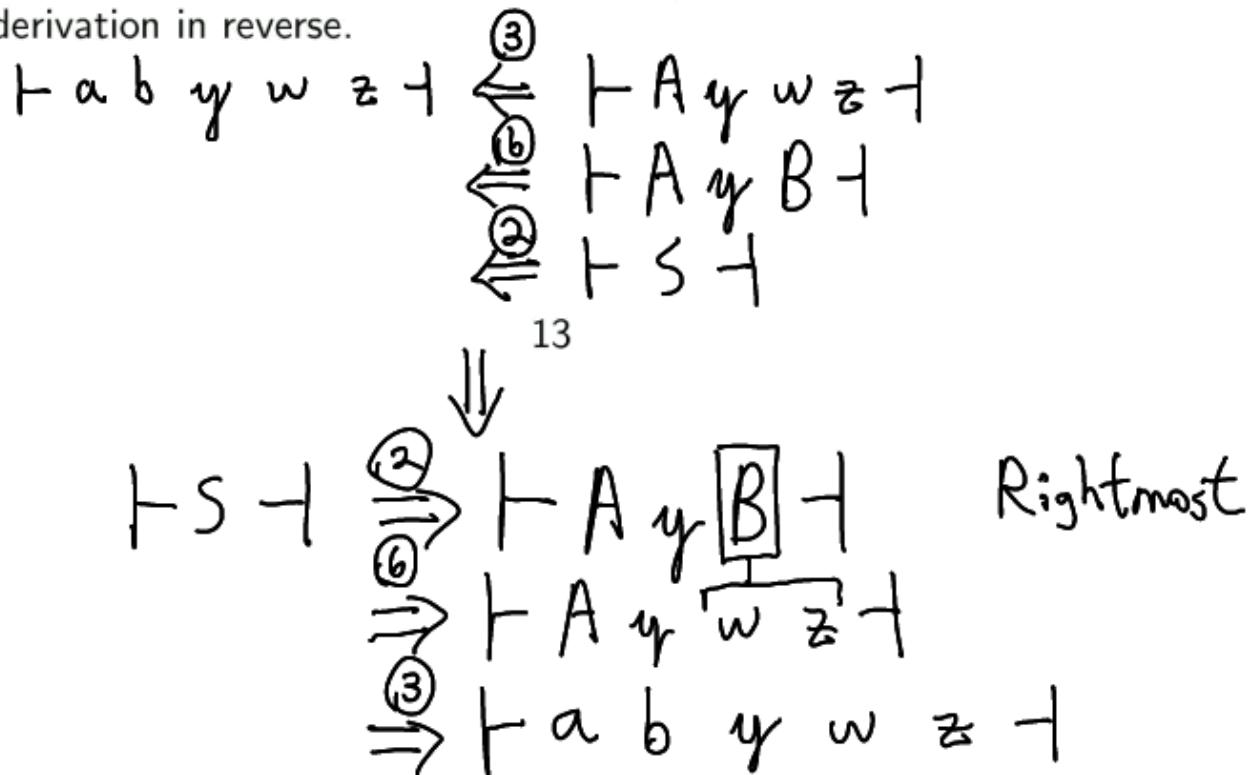
$$\begin{cases} A \rightarrow ab \\ B \rightarrow wZ \\ S \rightarrow AyB \end{cases}$$

A simple observation

- ▶ Didn't we say that this was LR(1) parsing? **yup**
- ▶ Doesn't the "R" mean rightmost derivation? **Yep**
- ▶ Aren't we always reducing the leftmost nonterminal?

Uh, yeah: $A \rightarrow ab$ $S \rightarrow \underline{A}yB$
 $B \rightarrow wz$

- ▶ But notice the direction we are creating the derivation. Write the derivation in reverse.

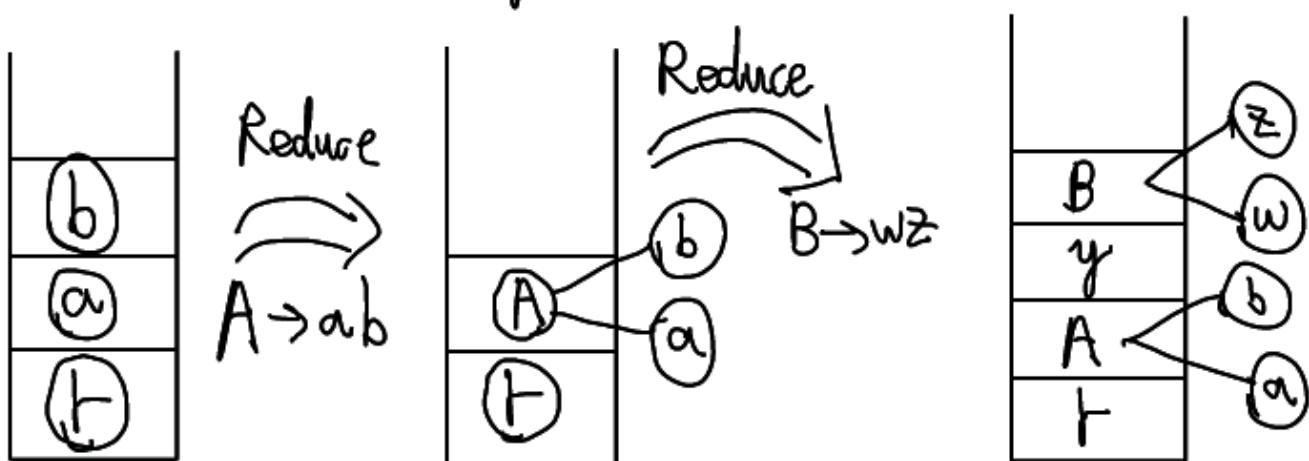


Outputting the parse tree

Algorithm

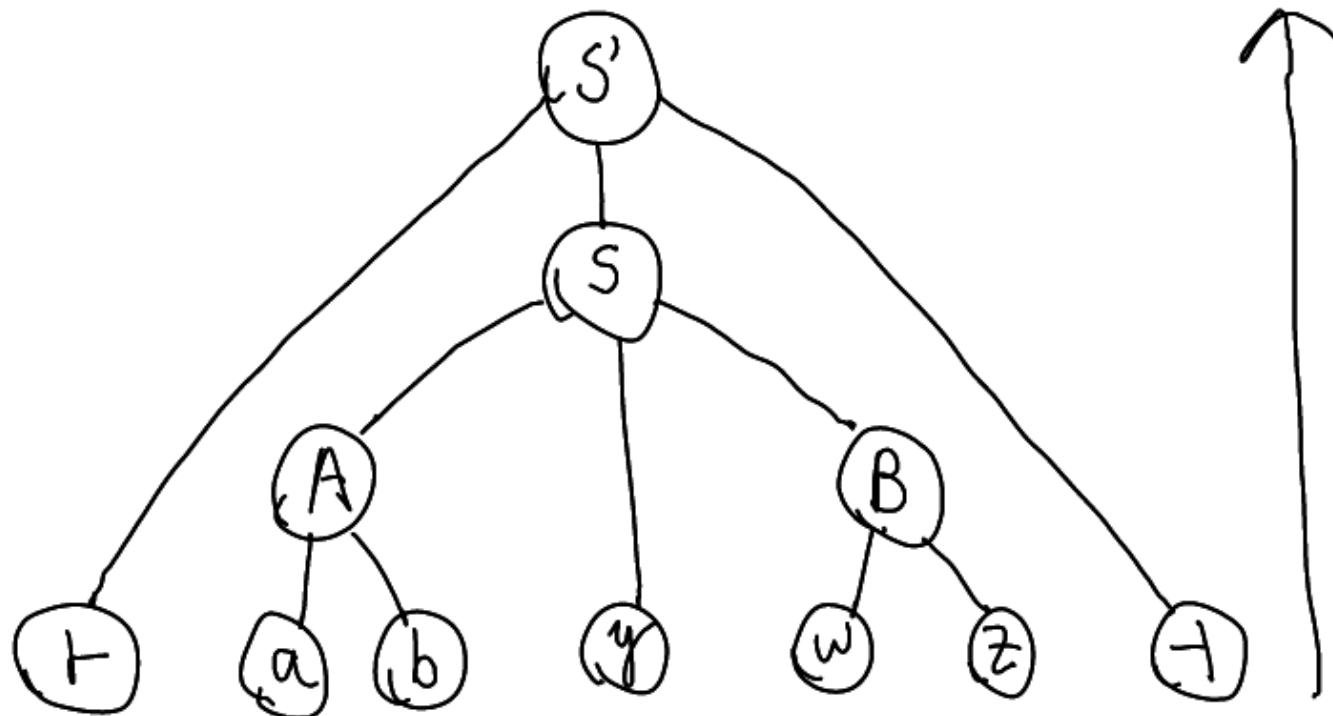
- ▶ Create a "tree stack"
- ▶ Each time we reduce, pop the right hand side nodes from tree stack
- ▶ Push the left hand side node and make its children the nodes we just popped
- ▶ Example: $\vdash a b y w z \dashv$

In sync with the (input, state) stack.

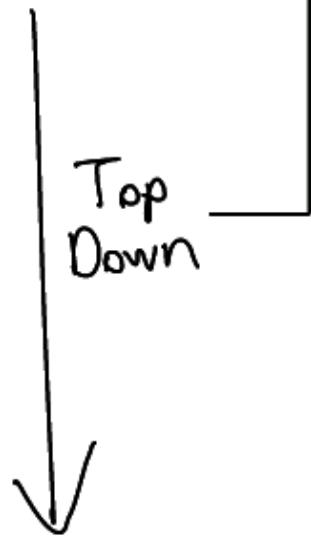
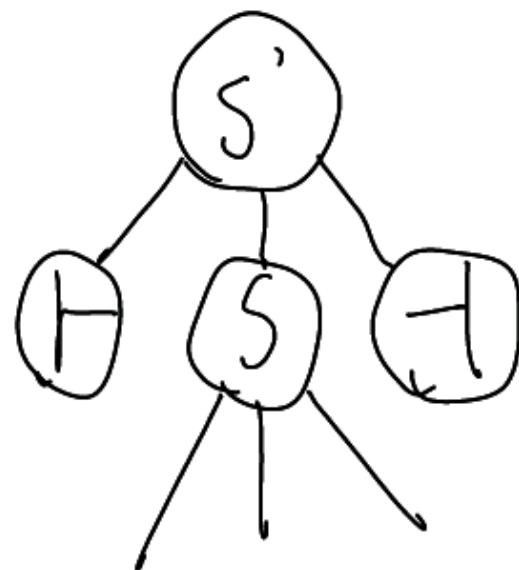
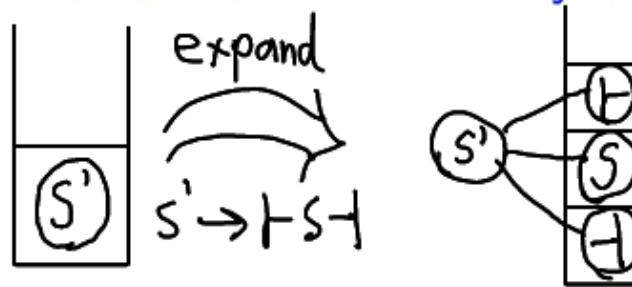


How the tree is actually built LR parsing

Bottom Up



How the tree is actually built LL parsing



Assignment hints

A7

P1, P2: write a cfg-r derivation by hand

RTFS
spec

P3: Write a parser

- ▶ Read a CFG, the DFA and input
- ▶ Output cfg-r (derivation)

→ the oracle

P4: Write a parser for ~~WLPP~~ WL~~P~~4

- ▶ Your scanner reads characters as input and outputs tokens
- ▶ Your parser will read tokens, builds a parse tree and outputs a left-most derivation
- ▶ Don't read the ~~WLPP~~ grammar from a separate file. Find a way to embed it in your program. Don't do this by typing it into your program directly!

]} A6

Going back

Looking at: $L = \{a^n b^m : n \geq m \geq 0\}$ (non-LL(1) language)

1. $S' \rightarrow \vdash S \dashv$] and non-LL(k)
2. $S \rightarrow a S$
3. $S \rightarrow T$
4. $T \rightarrow a T b$
5. $T \rightarrow \epsilon$
- } this is LR(1)

What to do when you see the symbol:

- ▶ $\vdash \Rightarrow$ Shift
- ▶ $a \Rightarrow$ Shift if no t on stack.
- ▶ $b \Rightarrow$ reduce by rule 5, shift, reduce by rule 4
- ▶ $\dashv \Rightarrow$ reduce by rule 3, repeatedly reduce by rule 2 until reach t, shift

$|L(G)|=1$
 $S \Rightarrow T$
 $T \Rightarrow Q$
 $Q \Rightarrow h_i$
 only one rule for each non terminal

Final fun facts

zero tokens of lookahead!!

$$S' \Rightarrow T \quad S \dashv$$

Notice: LL(0) grammars are useless

- ▶ Theorem: For any augmented LR(1) grammar, there is an equivalent LR(0) grammar. See CS 462
- ▶ Theorem: The class of languages that can be parsed deterministically with a stack can be represented with an LR(1) grammar.
- ▶ Comparing LL(1) vs. LR(1)

