Lecture 12 **CIS 4521/5521: COMPILERS**

Announcements

- HW3: LLVM Backend
 - Available on the course web pages.
 - Due: Weds., February 26th at 10:00PM
 - Note: test cases should be submitted TONIGHT at 10pm
- Midterm: March 6th
 - In class
 - One-page, letter-sized, double-sided "cheat sheet" of notes permitted
 - Coverage: interpreters, x86, IRs, LLVM IR, calling conventions, lexing, parsing
 - See Ed post for previous exams

Creating an abstract representation of program syntax.

PARSING

Zdancewic CIS 4521/5521: Compilers

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: $LHS \mapsto RHS$
 - LHS is a nonterminal
 - RHS is a *string* of terminals and nonterminals

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

Example: Left- and rightmost derivations

- Leftmost Derivation
- $\mathbf{S} \mapsto \mathbf{E} + \mathbf{S}$ \mapsto (S) + S \mapsto (**E** + S) + S \mapsto (1 + S) + S \mapsto (1 + E + S) + S \mapsto (1 + 2 + S) + S \mapsto (1 + 2 + E) + S \mapsto (1 + 2 + (S)) + S \mapsto (1 + 2 + (E + S)) + S \mapsto (1 + 2 + (3 + S)) + S \mapsto (1 + 2 + (3 + E)) + S \mapsto (1 + 2 + (3 + 4)) + S \mapsto (1 + 2 + (3 + 4)) + E \mapsto (1 + 2 + (3 + 4)) + 5

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

E \low number | (S)

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

 $T \mapsto S\$$ $S \mapsto ES'$ $S' \mapsto \varepsilon$ $S' \mapsto + S$ $E \mapsto number \mid (S)$

	number	+	()	\$ (EOF)
Т	\mapsto S\$		⊷S\$		
S	$\mapsto E S'$		$\mapsto E S'$		
S'		\mapsto + S		$\mapsto \epsilon$	$\mapsto \epsilon$
E	⊢ num.		$\mapsto (S)$		

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

LR GRAMMARS

Zdancewic CIS 4521/5521: Compilers

Bottom-up Parsing (LR Parsers)

- LR(k) parser:
 - <u>L</u>eft-to-right scanning
 - <u>R</u>ightmost 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, ocamlyacc, merlin, etc.)
 - Better error detection/recovery
 - But... harder to implement by hand

Top-down vs. Bottom up

• Consider the leftrecursive grammar:

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

- (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...



Example: Left- and rightmost derivations

• Leftmost Derivation

•	$\underline{S} \mapsto \underline{E} + S$
	→ (<u>§</u>) + S
	→ (<u>E</u> + S) + S
	$\mapsto (1 + \underline{\mathbf{S}}) + \mathbf{S}$
	$\mapsto (1 + \mathbf{\underline{E}} + \mathbf{S}) + \mathbf{S}$
	$\mapsto (1 + 2 + \mathbf{S}) + \mathbf{S}$
	$\mapsto (1 + 2 + \underline{\mathbf{E}}) + \mathbf{S}$
	$\mapsto (1 + 2 + (\underline{\mathbf{S}})) + \mathbf{S}$
	$\mapsto (1 + 2 + (\underline{\mathbf{E}} + \mathbf{S})) + \mathbf{S}$
	\mapsto (1 + 2 + (3 + <u>S</u>)) + S
	\mapsto (1 + 2 + (3 + <u>E</u>)) + S
	\mapsto (1 + 2 + (3 + 4)) + <u>S</u>
	\mapsto (1 + 2 + (3 + 4)) + <u>E</u>
	\mapsto (1 + 2 + (3 + 4)) + 5

• Rightmost derivation:

• $S \mapsto E + S$ → E **+ E** \mapsto E + 5 \mapsto (S) + 5 \mapsto (E + S) + 5 \mapsto (E + E + S) + 5 \mapsto (E + E + E) + 5 \mapsto (E + E + (S)) + 5 $\mapsto (E + E + (E + S)) + 5$ \mapsto (E + E + (E + E)) + 5 $\mapsto (E + E + (E + 4)) + 5$ \mapsto (E + E + (3 + 4)) + 5 \mapsto (**E** + 2 + (3 + 4)) + 5 \mapsto (1 + 2 + (3 + 4)) + 5

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

Progress of Bottom-up Parsing

Reductions	Scanned	Input Remaining
(1 + 2 + (3 + 4)) + 5 ←		(1 + 2 + (3 + 4)) + 5
$(\underline{\mathbf{E}} + 2 + (3 + 4)) + 5 \longleftarrow$	(1 + 2 + (3 + 4)) + 5
$(\underline{S} + 2 + (3 + 4)) + 5 \longleftarrow$	(1	+ 2 + (3 + 4)) + 5
(S + <u>E</u> + (3 + 4)) + 5 ↔	(1 + 2	+ (3 + 4)) + 5
(<u>S</u> + (3 + 4)) + 5 ↔	(1 + 2	+ (3 + 4)) + 5
$(S + (\underline{E} + 4)) + 5 \longleftarrow$	(1 + 2 + (3	+ 4)) + 5
$(S + (\underline{S} + 4)) + 5 \leftarrow 1$	(1 + 2 + (3	+ 4)) + 5
(S + (S + <u>E</u>)) + 5 ↔	(1 + 2 + (3 + 4)))) + 5
(S + (<u>S</u>)) + 5 ↔	(1 + 2 + (3 + 4)))) + 5
(S + <u>E</u>) + 5 ←	(1 + 2 + (3 + 4))) + 5
(<u>S</u>) + 5 ↔	(1 + 2 + (3 + 4))) + 5
<u>E</u> + 5 ↔	(1 + 2 + (3 + 4))	+ 5
<u>S</u> + 5 ↔	(1 + 2 + (3 + 4)) +	5
S + <u>E</u> ←	(1 + 2 + (3 + 4)) + 5	5
S		$S \mapsto S + F \mid F$

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

Shift/Reduce Parsing

- Parser state:
 - 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 γ at top of stack with nonterminal X such that X $\mapsto \gamma$ is a production. (pop γ , push X)

Stack	Input	Action
	(1 + 2 + (3 + 4)) + 5	shift (
(1 + 2 + (3 + 4)) + 5	shift 1
(1	+ 2 + (3 + 4)) + 5	reduce: $E \mapsto number$
(E	+ 2 + (3 + 4)) + 5	reduce: $S \mapsto E$
(S	+ 2 + (3 + 4)) + 5	shift +
(S +	2 + (3 + 4)) + 5	shift <mark>2</mark>
(S + 2	+ (3 + 4)) + 5	reduce: $E \mapsto number$



parser.mly, lexer.mll, range.ml, ast.ml, main.ml

DEMO: BOOLEAN LOGIC

Zdancewic CIS 4521/5521: Compilers

Simple LR parsing with no look ahead.

LR(0) GRAMMARS

Zdancewic CIS 4521/5521: Compilers

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 $\boldsymbol{\sigma}.$
 - Accept states of the DFA correspond to unique reductions that apply.
- Example: LR(0) parsing
 - <u>L</u>eft-to-right scanning, <u>R</u>ight-most derivation, <u>zero</u> 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:

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

• 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 stack + input
- Parsing is a sequence of *shift* and *reduce* operations:
- Shift: move look-ahead token to the stack: e.g.

	()	
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 $S \mapsto id$
(S	, (y, z), w)	reduce $L \mapsto S$

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

Example Run

Stack	Input	Action
	(x, (y, z), w)	shift (
(x, (y, z), w)	shift x
(x	, (y, z), w)	reduce $S \mapsto id$
(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 $S \mapsto id$
(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 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
45 2[/5521: Comp	oilers , w)	shift ,

CIS

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

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, $X \mapsto \varepsilon$ 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 $X \mapsto \gamma$ and $Y \mapsto \beta$, γ and β 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

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

- 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:

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

- Adds items for all productions whose LHS nonterminal occurs in an item in the 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: $CLOSURE({S' \mapsto .S}) = {S' \mapsto .S}, S \mapsto .(L), S \mapsto .id$
- Resulting "closed state" contains the set of all possible productions that might be reduced next.



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



- 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



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

- 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



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

- Finally, for each new state, we take the closure.
- Note that we have to perform two iterations to compute CLOSURE({S → (.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



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

9

- 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.

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 γ and push X.
- Optimization: No need to re-run the DFA from beginning 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:
 - Shift and goto state n
 - Reduce using reduction $X \mapsto \gamma$
 - First pop γ off the stack to reveal the state
 - Look up X in the "goto table" and goto that state



Example Parse Table

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

sx = shift and goto state x
gx = goto state x

Example

• Parse the token stream: (x, (y, z), w)\$

Stack	Stream	Action (according to table)
ϵ_1	(x, (y, z), w)\$	s3
$\varepsilon_1(_3$	x, (y, z), w)\$	s2
$\varepsilon_1(_3X_2$, (y, z), w)\$	Reduce: S⊷id
$\epsilon_1(_3S)$, (y, z), w)\$	g7 (from state 3 follow S)
$\epsilon_1(_3S_7$, (y, z), w)\$	Reduce: L→S
$\epsilon_1(_3L)$, (y, z), w)\$	g5 (from state 3 follow L)
$\epsilon_1(_3L_5$, (y, z), w)\$	s8
$\epsilon_1(_3L_{5'8})$	(y, z), w)\$	s3
$\epsilon_1(_3L_{5,8}(_3$	y, z), w)\$	s2