



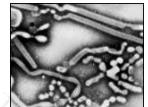
COMP 181

Lecture 9
LR parsing


October 3, 2006

Prelude



- What is this?
Micrograph of influenza A
- Where does influenza A come from?
 - Birds – aka “Avian Flu”
 - BUT**, it's extremely rare for humans to be infected
- So, how do people get influenza A
 - Pigs can get both avian and human flu strains
 - Virus recombines in pigs



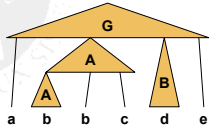

2

Bottom-up parsing

- Start with input stream
 - “Leaves” of parse tree
- Build up towards goal symbol
 - Construct the reverse derivation

#	Production rule
1	$G \rightarrow a A B e$
2	$A \rightarrow A b c$
3	$\mid b$
4	$B \rightarrow d$

Rule	Sentential form
-	$abcde$
3	$aAbcde$
2	$aAde$
4	$aABe$
1	G

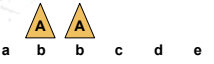
3

Bottom-up parsing


- Problem:**
 - Not good enough to simply find production right-hand sides and reduce
 - Example:

#	Production rule
1	$G \rightarrow a A B e$
2	$A \rightarrow A b c$
3	$\mid b$
4	$B \rightarrow d$

Rule	Sentential form
-	$abcde$
3	$aAbcde$
2	$aAAcde$
?	$\dots now what?$



- “ $aAAcde$ ” is not part of any sentential form




4

LR parsing

- State of the parser:
 - $\alpha \mid \gamma$
 - α is a stack of terminals and non-terminals
 - γ is string of unexamined terminals
- Two operations:
 - Shift** – read next terminal, push on stack
 $E + (\mid int) \rightarrow E + (int \mid)$
 - Reduce** – pop RHS symbols off stack, push LHS
 $E + (E + (E \mid)) \rightarrow E + (E \mid)$

#	Production rule
1	$E \rightarrow E + (E)$
2	$\mid int$




5

LR parsing

```

repeat
  if top symbols on stack match  $\beta$  for some  $A \rightarrow \beta$ 
    Reduce: “found an A”
    Pop those symbols off
    Push A on stack
  else Get next token from scanner
    if token is useful
      Shift: “still working on something”
      Push on stack
    else error
until stack contains goal and no more input
  
```



6

Example

1. | int + (int) + (int) Nothing on stack, get next token

Stack

Tufts University Computer Science 7

Example

1. | int + (int) + (int) Nothing on stack, get next token
2. int | + (int) + (int) *Shift*: push int

Top of stack matches $E \rightarrow int$

Stack int

Tufts University Computer Science 8

Example

1. | int + (int) + (int) Nothing on stack, get next token
2. int | + (int) + (int) *Shift*: push int
3. int | + (int) + (int) *Reduce*: pop int

Top of stack matches $E \rightarrow int$

Stack int

Tufts University Computer Science 9

Example

1. | int + (int) + (int) Nothing on stack, get next token
2. int | + (int) + (int) *Shift*: push int
3. int | + (int) + (int) *Reduce*: pop int, push E

Top of stack matches $E \rightarrow int$

Stack E int

Tufts University Computer Science 10

Example

1. | int + (int) + (int) Nothing on stack, get next token
2. int | + (int) + (int) *Shift*: push int
3. int | + (int) + (int) *Reduce*: pop int, push E
4. int + | (int) + (int) *Shift*: push +
5. int + (| int) + (int) *Shift*: push (
6. int + (int |) + (int) *Shift*: push int

Top of stack matches $E \rightarrow int$

Stack E + (int

Tufts University Computer Science 11

Example

1. | int + (int) + (int) Nothing on stack, get next token
2. int | + (int) + (int) *Shift*: push int
3. int | + (int) + (int) *Reduce*: pop int, push E
4. int + | (int) + (int) *Shift*: push +
5. int + (| int) + (int) *Shift*: push (
6. int + (int |) + (int) *Shift*: push int
7. int + (int |) + (int) *Reduce*: pop int, push E

Stack E + (E

Tufts University Computer Science 12

Example

1. | int + (int) + (int) Nothing on stack, get next token
2. int | + (int) + (int) *Shift:* push int
3. int | + (int) + (int) *Reduce:* pop int, push E
4. int + | (int) + (int) *Shift:* push +
5. int + (| int) + (int) *Shift:* push (
6. int + (int |) + (int) *Shift:* push int
7. int + (int |) + (int) *Reduce:* pop int, push E
8. int + (int) | + (int) *Shift:* push)

Top of stack matches E → E + (int)

Stack E + (E)



Example

1. | int + (int) + (int) Nothing on stack, get next token
2. int | + (int) + (int) *Shift:* push int
3. int | + (int) + (int) *Reduce:* pop int, push E
4. int + | (int) + (int) *Shift:* push +
5. int + (| int) + (int) *Shift:* push (
6. int + (int |) + (int) *Shift:* push int
7. int + (int |) + (int) *Reduce:* pop int, push E
8. int + (int) | + (int) *Shift:* push)
9. int + (int) | + (int) *Reduce:* pop x 5, push E

Stack E



Example

-
9. int + (int) | + (int) *Reduce:* pop x 5, push E
 10. int + (int) + | (int) *Shift:* push +
 11. int + (int) + (| int) *Shift:* push (
 12. int + (int) + (int |) *Shift:* push int

Stack E + (int



Example

-
9. int + (int) | + (int) *Reduce:* pop x 5, push E
 10. int + (int) + | (int) *Shift:* push +
 11. int + (int) + (| int) *Shift:* push (
 12. int + (int) + (int |) *Shift:* push int
 13. int + (int) + (int |) *Reduce:* pop int, push E
 14. int + (int) + (int) | *Shift:* push)

Stack E + (E)



Example

-
9. int + (int) | + (int) *Reduce:* pop x 5, push E
 10. int + (int) + | (int) *Shift:* push +
 11. int + (int) + (| int) *Shift:* push (
 12. int + (int) + (int |) *Shift:* push int
 13. int + (int) + (int |) *Reduce:* pop int, push E
 14. int + (int) + (int) | *Shift:* push)
 15. int + (int) + (int) | *Reduce:* pop x 5, push E

DONE!

Stack E



Key problems

- (1) Will this work?
How do we know that shifting and reducing using a stack is sufficient to compute the reverse derivation?
- (2) How do we know when to shift and reduce?
 - Can we efficiently match top symbols on the stack against productions?
 - *Right-hand sides of productions may have parts in common*
 - Will shifting a token move us closer to a reduction?
 - *Are we making progress?*
 - *How do we know when an error occurs?*



Key problems

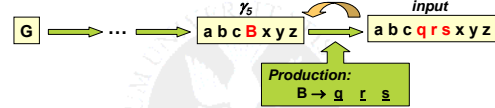
- (1) Will this always work?
 - Yes, for unambiguous grammars
- Why?
 - Unambiguous:
 - Unique right-most derivation for every string
 - At each parsing step, one possible reduction

$$G \rightarrow \gamma_1 \rightarrow \gamma_2 \rightarrow \gamma_3 \rightarrow \gamma_4 \rightarrow \gamma_5 \rightarrow \text{input}$$



Shift-reduce parsing

- Consider last step:



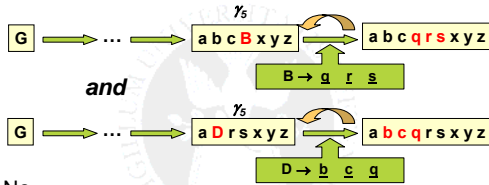
- To reverse this step:
 - Shift until \underline{g} , \underline{r} , \underline{s} on top of stack
 - Reduce: pop \underline{g} , \underline{r} , \underline{s} , push \underline{B}
- Parsing state:

Input:	a b c q r s x y z
Stack:	a b c B



Right-most derivation

- Could there be an alternative reduction?

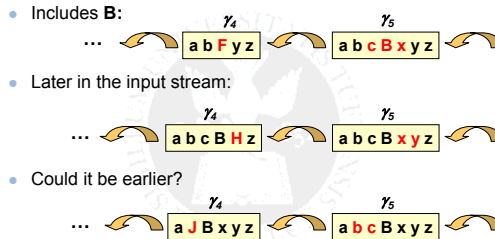


- No
 - Two right-most derivations for the same string
 - I.e., the grammar would be ambiguous



Reductions

- Where is the next reduction?



- No – this is not the right-most derivation!



Implications

- Cases:



- Parsing state:

Input:	a b c q r s x y z
Stack:	a b c B

- **Key:** next reduction must consume top of stack
Possibly after shifting more terminal symbols

- How does this help?
 - Can consume terminal symbols in order
 - Never need to search inside the stack

We can perform LR parsing using only stack operations



Key problems

- (2) How do we know when to shift or reduce?

- Reductions
 - Good news:
 - At any given step, reduction is unique
 - Matching production occurs at top of stack
 - Problem:
 - How to efficiently find the right production
- Shifts
 - Default behavior: shift when there's no reduction
 - Still need to handle errors



When to shift or reduce

- What is on the stack? $\left\{ \begin{array}{l} \text{Input: } a b c q r s \mid x y z \\ \text{Stack: } \underline{a} \underline{b} \underline{c} B \end{array} \right.$
- Either:**
 - A right sentential form
 - A prefix of a right sentential form (missing some terminals)
Called a *viable prefix*
- Idea:** a DFA that recognizes viable prefixes
 - Input: stack contents (a mix of terminals, non-terminals)
 - Each state represents either
 - A right sentential form – labeled with the reduction to apply
 - A viable prefix – labeled with tokens to expect next



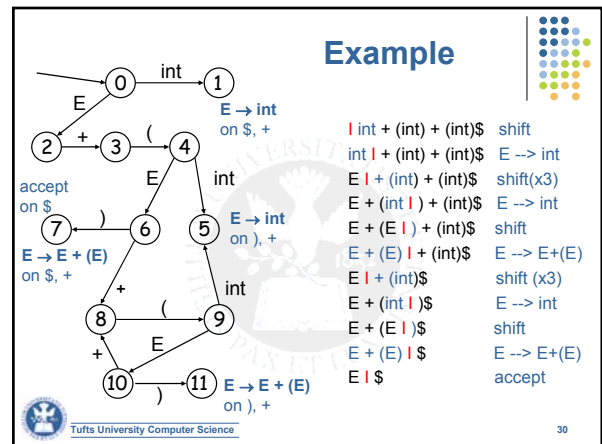
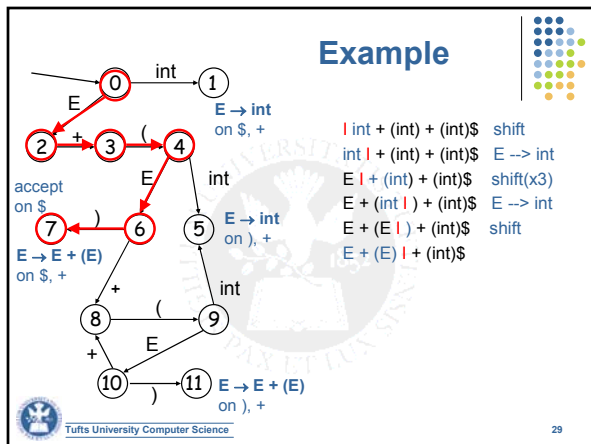
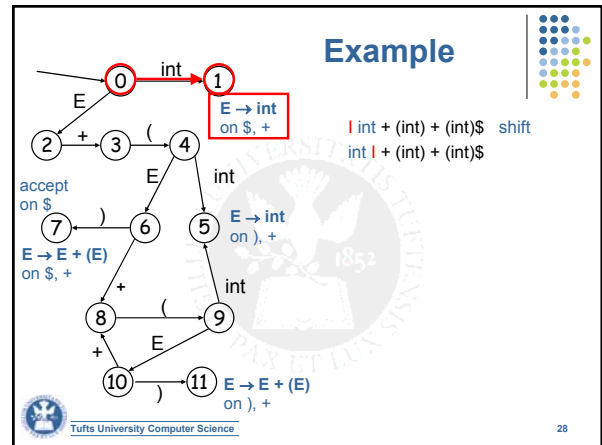
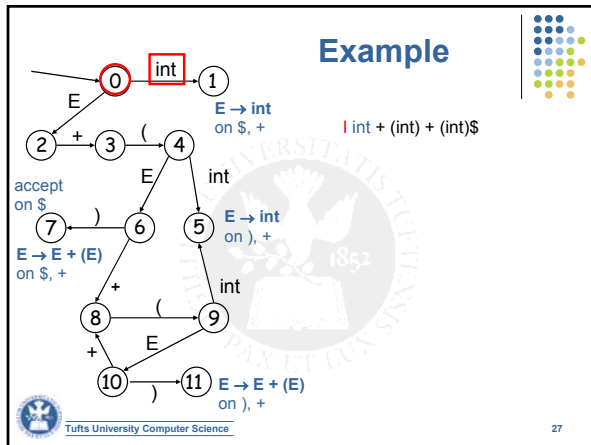
Shift/reduce DFA

- Using the DFA
 - At each parsing step run DFA on stack contents
 - Examine the resulting state X and the token t immediately following | in the input stream
 - If X has an outgoing edge labeled t, then **shift**
 - if X is labeled "A → β on t", then **reduce**

Example:

#	Production rule
1	$E \rightarrow E + (E)$
2	$\mid \text{int}$

- First, we'll look at how to use such a DFA...

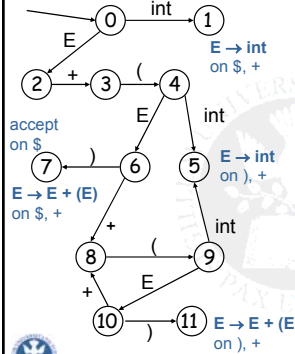


Improvements

- Each DFA state represents stack contents
 - At each step, we rerun the DFA to compute the new state
 - Can we avoid this?
 - Two actions:
 - Shift: Push a new token
 - Reduce: Pop some symbols off, push a new symbol
- **Idea:**
 - For each symbol on the stack, remember the DFA state that represents the contents up to that point
 - Push a new token = go forward in DFA
 - Pop a sequence of symbols = "unwind" DFA to previous state



Example



$I \text{ int} + (\text{int}) + (\text{int})\$$ shift
 $\text{int } I + (\text{int}) + (\text{int})\$$ $E \rightarrow \text{int}$
 $E I + (\text{int}) + (\text{int})\$$ shift(x3)
 At state 2
 go forward in DFA
 $2 \rightarrow 3 \rightarrow 4 \rightarrow 5$
 $E + (\text{int } I) + (\text{int})\$$ $E \rightarrow \text{int}$
 Back up to state 4
 Go forward with E
 $4 \rightarrow 6$



Algorithm components

- Stack
 - String of the form : $\langle \text{sym}_1, \text{state}_1 \rangle, \dots, \langle \text{sym}_n, \text{state}_n \rangle$
 - **sym_i**: grammar symbol (left part of string)
 - **state_i**: DFA state
 - Intuitively: represents what we've seen so far
 - state_k is the final state of the DFA on sym₁ ... sym_k
 - And, captures what we're looking for next
- Represent as two tables:
 - **action** – whether to shift, reduce, accept, error
 - **goto** – next state



Tables

- **Action**
 Given state and the next token, **action**[s,a] =
 - Shift s', where s' is a state
 - Reduce by a grammar production $A \rightarrow \beta$
 - Accept
 - Error
- **Goto**
 Given a state and a grammar symbol, **goto**[s,X] =
 - Transition to next state (after recognizing an X)



Algorithm

```

push s0 on stack
token = scanner.next_token()
repeat
  s = state at top of stack
  if action[s, token] = reduce A -> beta then
    pop |beta| pairs (Xi, sm) off the stack
    s' = top of stack
    push A on stack
    push goto[s', A] on stack
  else if action[s, token] = shift s' then
    push token on stack
    push s on stack
    token = scanner.next_token()
  else if action(s, token) = accept then
    return true
  else error()
    
```

Top of stack is handle
 $A \rightarrow \beta$

- Work
 - Shift each token
 - Pop each token
- Errors
 - Input exhausted
 - Error entry in table

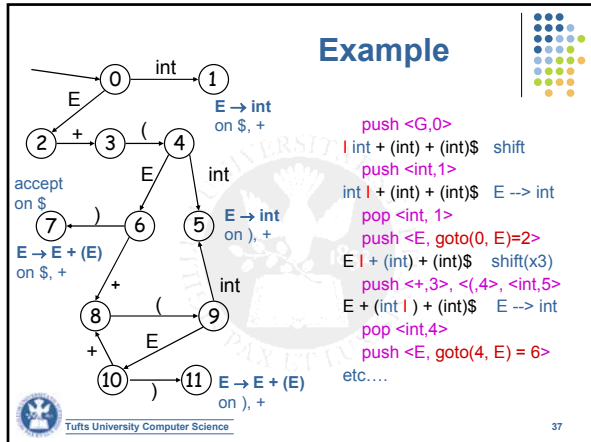


Representing the DFA

- Combined table:

	action(state, token)					goto
	int	+	()	\$	E	
...						
3			s4			
4	s5					g6
5	r _{E→int}		r _{E→int}			
6	s8		s7			
7	r _{E→E+(E)}				r _{E→E+(E)}	
...						





How is the DFA Constructed?

- The stack describes the context of the parse
 - What non-terminal we are looking for
 - What production rhs we are looking for
 - What we have seen so far from the rhs
- Each DFA state represents several such contexts
 - E.g., when we are looking for non-terminal **E**, we might be looking either for an **int** or a **E + (E)** on the RHS
 - We want to capture all the productions we *could be* working on

Tufts University Computer Science 38

LR Items

- An **LR(1) item** is a pair:

$$[A \rightarrow \alpha \cdot \beta, a]$$
 - $A \rightarrow \alpha\beta$ is a production
 - a is a terminal (the lookahead terminal)
 - LR(1) means 1 lookahead terminal
- $[A \rightarrow \alpha \cdot \beta, a]$ describes a context of the parser
 - We are trying to find an **A** followed by an **a**, and
 - We have α on top of the stack
 - We need to see a prefix derived from βa

Tufts University Computer Science 39

LR Items

- In context containing

$$[E \rightarrow E + \cdot (E), +]$$
 - If (next then we can a **shift** to context containing

$$[E \rightarrow E + (\cdot E), +]$$
- In context containing

$$[E \rightarrow E + (E) \cdot, +]$$
 - We can **reduce** with $E \rightarrow E + (E)$
 - But only if a **+** follows

Tufts University Computer Science 40

LR Items

- Consider the item

$$E \rightarrow E + (\cdot E), +$$
 - What could we see next?
 - We expect a string derived from **E** +
 - There are two productions for **E**

$$E \rightarrow \text{int} \quad \text{and} \quad E \rightarrow E + (E)$$
- We extend the context with two more items:

$$E \rightarrow \cdot \text{int},)$$

$$E \rightarrow \cdot E + (E),)$$

Tufts University Computer Science 41

Closure operation

- The operation of extending the context with items is called the **closure** operation
- Closure function:
 - Given a set of items
 - Compute all other items that could represent the current parsing state
- Observation:
 - At $A \rightarrow \alpha \cdot B\beta$ we expect to see **B** next
 - If $B \rightarrow \gamma$ is a production, then we could also see a γ

Tufts University Computer Science 42

Closure operation

- Algorithm:

$\text{closure}(\text{Items}) =$
 repeat
 for each $[A \rightarrow \alpha \cdot B\beta, \underline{a}]$ in Items
 for each production $B \rightarrow \gamma$
 add $[B \rightarrow \cdot \gamma, \underline{?}]$ to Items
 until Items is unchanged

What is the lookahead?



Closure operation

- Algorithm:

$\text{closure}(\text{Items}) =$
 repeat
 for each $[A \rightarrow \alpha \cdot B\beta, \underline{a}]$ in Items
 for each production $B \rightarrow \gamma$
 for each $\underline{b} \in \text{FIRST}(\beta a)$
 add $[B \rightarrow \cdot \gamma, \underline{b}]$ to Items
 until Items is unchanged



Building the DFA – part 1

- Starting context = $\text{closure}(\{S \rightarrow \cdot E, \$\})$

$S \rightarrow \cdot E, \$$
 $E \rightarrow \cdot E+(E), \$$
 $E \rightarrow \cdot \text{int}, \$$
 $E \rightarrow \cdot E+(E), +$
 $E \rightarrow \cdot \text{int}, +$

- Abbreviated:

$S \rightarrow \cdot E, \$$
 $E \rightarrow \cdot E+(E), \$/+$
 $E \rightarrow \cdot \text{int}, \$/+$



Building the DFA – part 2

- DFA states

- Each DFA state is a closed set of LR(1) items
- Start state: $\text{closure}(\{S \rightarrow \cdot E, \$\})$

- Reductions

- Label each item $[A \rightarrow \alpha\beta \cdot, \underline{x}]$ with “Reduce with $A \rightarrow \alpha\beta$ on \underline{x} ”

- What about transitions?



DFA transitions

- Idea:

- If the parser was in state $[A \rightarrow \alpha \cdot X\beta]$ and then recognized an instance of X, then the new state is $[A \rightarrow \alpha X \cdot \beta]$
- Note: X could be a terminal or non-terminal

- Algorithm:

$\text{transition}(I, X) =$
 $J = \{\}$
 for each $[A \rightarrow \alpha \cdot X\beta, \underline{b}] \in I$
 add $[A \rightarrow \alpha X \cdot \beta, \underline{b}]$ to J
 return $\text{closure}(J)$



DFA construction

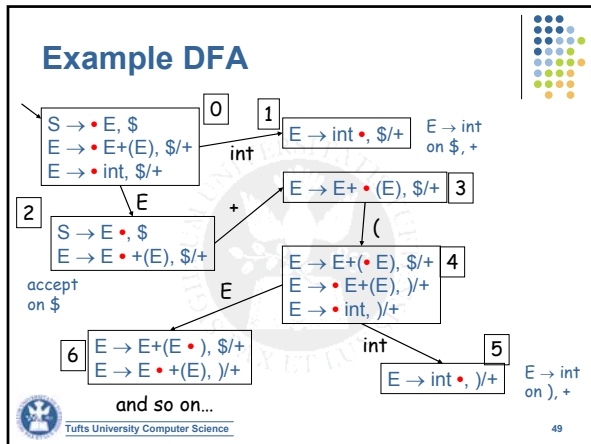
- Data structure:

- T – set of states (each state is a set of items)
- E – edges of the form $I \xrightarrow{X} J$
 where $I, J \in T$ and X is a terminal or non-terminal

- Algorithm:

$T = \{\text{closure}(\{S \rightarrow \cdot E, \$\}), E = \{\}$
 repeat
 for each state I in T
 for each item $[A \rightarrow \alpha \cdot X\beta, \underline{b}] \in I$
 let $J = \text{transition}(I, X)$
 $T = T \cup J$
 $E = E \cup \{I \xrightarrow{X} J\}$
 until E and T no longer change





- ### To form into tables
- Two tables
 - action(I, token)
 - goto(I, symbol)
 - Layout:
 - One row for each states – each I in T
 - One columns for each symbol
 - Entries:
 - For each edge $I \xrightarrow{X} J$
 - If X is a terminal, add shift J at position (I, X) in action
 - if X is a non-terminal, add goto J at position (I, X) goto
 - For each state $[A \rightarrow \alpha\beta \cdot, \underline{x}]$ in I
 - Add reduce n at position (I, \underline{x}) in action (where n is |rhs|)
- Tufts University Computer Science 50

- ### Next time...
- Some issues with bottom-up parsing
 - Parser-generator tools
 - New programming assignment
- Tufts University Computer Science 51