### Lecture 7: General and Bottom-Up Parsing

#### Administrivia

- Project #1 posted. Due 27 Feb.
- HW #3 posted, due next Monday. HW #4 goes back to Friday schedule.
- Test #1: March 10 (in class).

## Fixing Recursive Descent

• Can formulate top-down parsing analogously to NFAs.

```
parse (A, S):
  """Assuming A is a nonterminal and S = c_1c_2...c_n is a string,
    return integer k such that A can derive the string c_1...c_k."""
  Choose production 'A: \alpha_1\alpha_2\cdots\alpha_m' for A (nondeterministically)
  k = 0
  for x in \alpha_1, \alpha_2, \cdots, \alpha_m:
    if x is a terminal:
        if x = c_{k+1}:
            k += 1
        else:
            GIVE UP
  else:
            k += parse (x, c_{k+1}\cdots c_n)
  return k
```

- $\bullet$  Assume that the grammar contains one production for the start symbol: p:  $\gamma\dashv$  .
- We'll say that a call to parse returns a value if some set of choices for productions (the blue step) would return a value (just like NFA).

• Then if parse(p, S) returns a value, S must be in the language. Last modified: Wed Mar 11 19:41:44 2009 CS164: Lecture #7 2

## Making a Deterministic Algorithm

- If we had an infinite supply of processors, could just spawn new ones at each "Choose" line.
- Some would give up, some loop forever, but on correct programs, at least one process would get through.
- To do this for real (say with one processor), need to keep track of all possibilities systematically.
- This is the idea behind Earley's algorithm:
  - Handles any context-free grammar.
  - Finds all parses of any string.
  - Runs in  $O(N^3)$  time for ambiguous grammars,  $O(N^2)$  time for "non-deterministic grammars", or O(N) time for deterministic grammars (such as accepted by Bison).

## Earley's Algorithm: I

• First, reformulate to ditch the loop. Assume the string  $S = c_1 \cdots c_n$  is fixed.

```
parse (A: \alpha \bullet \beta, s, k):
    """Assumes A: \alpha\beta is a production in the
       grammar, 0 <= s <= k <= n, and \alpha can produce the string
       c_{s+1} \cdots c_k. Returns integer j such that \beta can
       produce c_{k+1} \cdots c_i."""
    if \beta is empty:
       return k
   Assume \beta has the form y\delta
    if y is a terminal:
        if y == c_{k+1}:
             return parse(A: \alpha y \bullet \delta, s, k+1)
       else
             GIVE UP
    else:
       Choose production 'y: \kappa' for y (nondeterministically)
        j = parse(y: \bullet \kappa, k, k)
       return parse (A: \alpha y \bullet \delta, s, j)
```

• Now do all possible choices that result in such a way as to avoid redundant work ("nondeterministic memoization").

Last modified: Wed Mar 11 19:41:44 2009

## **Chart Parsing**

- Idea is to build up a table (known as a *chart* of all calls to parse that have been made.
- Only one entry in chart for each distinct triple of arguments (A:  $\alpha \bullet \beta$ , s, k).
- We'll organize table in columns numbered by the k parameter, so that column k represents all calls that are looking at  $c_{k+1}$  in the input.
- Each column contains entries with the other two parameters: [A:  $\alpha \bullet \beta$ , s], which is called an *item*.
- The columns, therefore, are *item sets*.

### Example

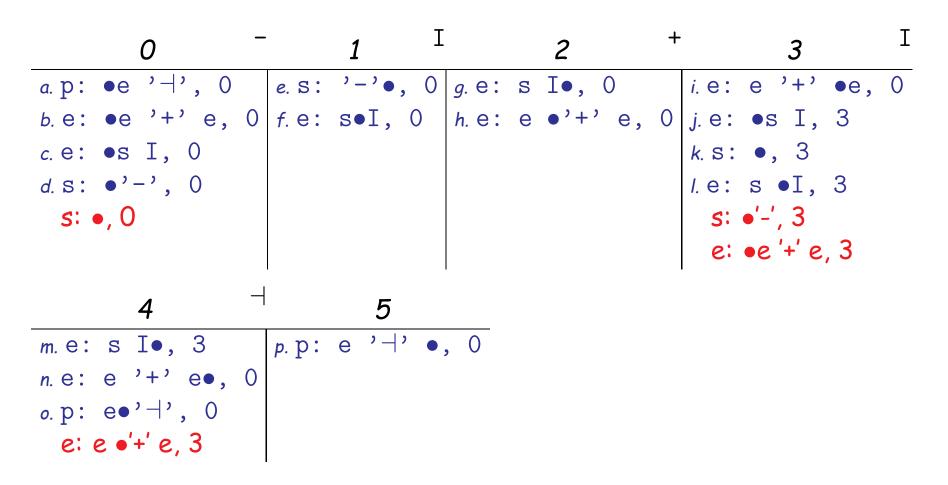
# Grammar Input String p : e '⊣' - I + I ⊣ e : s I | e '+' e - I + I ⊣ s : '-' | - I + I →

**Chart.** Headings are values of k and  $c_{k+1}$  (raised symbols).

0	-	1 <sup>I</sup>		2	+	3	I
a.p: ●e '⊢', 0							
b.e: ●e '+' e,	0 <i>f.</i> e:	s•I, 0	h.e: e	•'+' e,	0 j.e:	•s I, 3	
<i>c</i> .e: ●s I, 0					<i>k.</i> s:	●, 3 s ●I, 3	
d.s: •'-', 0					<i>I.</i> e:	s •I, 3	
4	-	5					
m.e: s I•, 3	<i>p.</i> p:	e '⊣' ●	, 0				
n.e: e '+' e●,	0						
<i>o</i> .p: e●'⊣', 0							

### Example, completed

• Last slide showed only those items that survive and get used. Algorithm actually computes dead ends as well (unlettered, in red).



## Adding Semantic Actions

- Pretty much like recursive descent. The call  $parse(A: \alpha \bullet \beta, s, k)$  can return, in addition to j, the semantic value of the A that matches characters  $c_{s+1} \cdots c_j$ .
- This value is actually computed during calls of the form parse(A:  $\alpha' \bullet$ , s, k) (i.e., where the  $\beta$  part is empty).
- Assume that we have attached these values to the nonterminals in  $\alpha$ , so that they are available when computing the value for A.

## Ambiguity

- Ambiguity only important here when computing semantic actions.
- Rather than being satisfied with a single path through the chart, we look at *all* paths.
- And we attach the set of possible results of  $parse(Y: \bullet \kappa, s, k)$  to the nonterminal Y in the algorithm.