

- We have seen that recursive-descent parsing is a simple and straightforward way to convert a grammar to a program that parses source using that grammar.
- However, because one has to predict which production to take without having seen the source tokens to be produced, it needs workarounds, as we've seen.
- In particular, must eliminate *left-recursion* and perform *left factoring* to make sure that branches are unique.
- So let's see what happens when we put off the decision about what production to use until after we've examined the text to be produced.
- This entails processing the children of a node in the parse tree *before* deciding on the production for that node; we determine the parse tree *from the bottom up*.

A Little Notation

Here and in lectures to follow, we'll often have to refer to general productions or derivations. In these, we'll use various alphabets to mean various things:

- Capital roman letters are nonterminals (A, B, \dots).
- Lower-case roman letters are terminals (or tokens, characters, etc.).
- Lower-case greek letters are sequences of zero or more terminal and nonterminal symbols, such as appear in sentential forms or on the right sides of productions (α, β, \dots).
- Subscripts on lower-case greek letters indicate individual symbols within them, so $\alpha = \alpha_1 \alpha_2 \dots \alpha_n$ and each α_i is a single terminal or nonterminal.

So $A ::= \alpha$ might describe the production $e ::= e '+' t$,

...and $B \Rightarrow \alpha A \gamma \Rightarrow \alpha \beta \gamma$ might describe the derivation steps

$e \Rightarrow e '+' t \Rightarrow e '+' ID$

(α is $e '+'$; A is t ; B is e ; and γ is empty.)

Fixing Recursive Descent

- First, let's define an impractical but simple implementation of a top-down parsing routine.
- For nonterminal A and string $S = c_1 c_2 \dots c_n$, we'll define $\text{parse}(A, S)$ to return the length of a valid prefix of S derivable from A .
- That is, $\text{parse}(A, c_1 c_2 \dots c_n) = k$, where

$$\frac{c_1 c_2 \dots c_k}{A \xRightarrow{*}} c_{k+1} c_{k+2} \dots c_n$$

Abstract body of parse(A, S)

- Can formulate top-down parsing analogously to NFAs.

```

parse (A, S):
    """"Assuming A is a nonterminal and S = c1c2...cn is a string, return
    integer k such that A can derive the prefix string c1...ck of S.""
    Choose production 'A: α1α2...αm' for A (nondeterministically)
    k = 0
    for x in α1, α2, ..., αm:
        if x is a terminal:
            if x == ck+1:
                k += 1
            else:
                GIVE UP
        else:
            k += parse (x, ck+1...cn)
    return k
    
```

- Let the start symbol be **p** with exactly one production: $p ::= \gamma \dashv$.
- We'll say that a call to parse returns a value if *some* set of choices for productions (the blue step) would not give up (just like NFA).
- Then if parse(p, S) returns a value, S must be in the language.

Example

Consider parsing $S = "ID*ID\dashv"$ with a grammar from last time:

```

p ::= e '⊣'
e ::= t
    | e '/' t
    | e '*' t
t ::= ID
    
```

A successful path through the program:

```

parse(p, S):
    parse(p, S):
        Choose p ::= e '⊣':
            parse(e, S):
                parse(e, S):
                    Choose e ::= e '*' t:
                        parse(e, S):
                            parse(e, S):
                                Choose e ::= t:
                                    parse(e, S):
                                        parse(t, S):
                                            choose t ::= ID:
                                                check S[1] == ID; OK, so k3 += 1;
                                                choose t ::= ID:
                                                    return 1 (so k3 added to k2)
                                                check S[1] == ID; OK, return 1
                                            return 1 (and add to k1)
                                        Check S[2] == S[k1+1] == '*', OK, k2 += 1
                                        check S[k2] == '*', OK, k2 += 1
                                        parse(t, S3): # S3 == "ID ⊣" == '*')
                                            choose t ::= ID:
                                                check S3[k3+1] == S3[1] == ID; OK
                                                k3 += 1; return 1 (so k2 += 1)
                                            return 3
                                        Check S[k1+1] == S[4] == '⊣': OK
                                        k1 += 1; return 4
                                
```

k_i means "the variable k in the call to parse that is nested i deep." Outermost k is k_1 . Likewise for S_i .

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 processor 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.
 - Can recognize or reject strings in $O(N^3)$ time for ambiguous grammars, $O(N^2)$ time for "nondeterministic grammars", or $O(N)$ time for deterministic grammars (such as accepted by Bison or CUP).

Earley's Algorithm: I

- First, reformulate to use recursion instead of looping. Assume the string $S = c_1 \dots c_n$ is fixed.
- Redefine **parse**:


```

parse (A: α • β, s, k):
    """"Assumes A: αβ is a production in the grammar,
    0 <= s <= k <= n, and α can produce the string cs+1...ck.
    Returns integer j such that β can produce ck+1...cj.""
            
```
- Or diagrammatically, **parse** returns an integer j such that:

$$c_1 \dots c_s \underbrace{c_{s+1} \dots c_k}_{\alpha \Rightarrow^*} \underbrace{c_{k+1} \dots c_j}_{\beta \Rightarrow^*} c_{j+1} \dots c_n$$

Earley's Algorithm: II

```

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  $c_{s+1} \dots c_k$ .
    Returns integer j such that  $\beta$  can produce  $c_{k+1} \dots c_j$ .""
    if  $\beta$  is empty:
        return k
    Assume  $\beta$  has the form  $x\delta$ 
    if x is a terminal:
        if x ==  $c_{k+1}$ :
            return parse(A ::=  $\alpha x \bullet \delta$ , s, k+1)
        else:
            GIVE UP
    else:
        Choose production ' $x ::= \kappa$ ' for x (nondeterministically)
        j = parse(x ::=  $\bullet \kappa$ , k, k)
        return parse (A ::=  $\alpha x \bullet \delta$ , s, j)
    
```

- Now do all possible choices that result in such a way as to avoid redundant work ("nondeterministic memoization").
- That is, if parse is called with the same three arguments as a previous call, just use the result(s) of the previous call.

Last modified: Mon Feb 11 01:05:40 2019

CS164: Lecture #7 9

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 are called *items*.
- The columns, therefore, are *item sets*.

Last modified: Mon Feb 11 01:05:40 2019

CS164: Lecture #7 10

Example

Grammar

```

p ::= e '+'
e ::= s I | e '+' e
s ::= '-' |
    
```

Input String

- I + I -

Chart. Headings are values of k and c_{k+1} (raised symbols). Item labels ($a-f$) trace the "ancestry" of each item. (Have shortened ' $::=$ ' to ':' for compactness.)

0	-	1	I	2	+	3	I
a.p: $e \bullet '+'$, 0	d.s: $'-\bullet$, 0	c.e: s I \bullet , 0		b.e: e '+' $\bullet e$, 0			
b.e: $e \bullet '+' e$, 0	c.e: s \bullet I, 0	b.e: e $\bullet '+' e$, 0		e.e: $\bullet s$ I, 3			
c.e: $\bullet s$ I, 0				f.s: \bullet , 3			
d.s: $\bullet '-'$, 0				e.e: s \bullet I, 3			
				i.s: $\bullet '-'$, 3			
				j.e: $\bullet e '+' e$, 3			
4	-	5					
e.e: s I \bullet , 3		a.p: e '+' \bullet , 0					
b.e: e '+' e \bullet , 0							
a.p: e '+' \bullet , 0							

Last modified: Mon Feb 11 01:05:40 2019

CS164: Lecture #7 11

Example, completed

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

0	-	1	I	2	+	3	I
a.p: $e \bullet '+'$, 0	d.s: $'-\bullet$, 0	c.e: s I \bullet , 0		b.e: e '+' $\bullet e$, 0			
b.e: $e \bullet '+' e$, 0	c.e: s \bullet I, 0	b.e: e $\bullet '+' e$, 0		e.e: $\bullet s$ I, 3			
c.e: $\bullet s$ I, 0		a.p: e $\bullet '+'$, 0		f.s: \bullet , 3			
d.s: $\bullet '-'$, 0				e.e: s \bullet I, 3			
g.s: \bullet , 0				i.s: $\bullet '-'$, 3			
h.e: s \bullet I, 0				j.e: $\bullet e '+' e$, 3			
4	-	5					
e.e: s I \bullet , 3		a.p: e '+' \bullet , 0					
b.e: e '+' e \bullet , 0							
a.p: e '+' \bullet , 0							
j.e: e $\bullet '+' e$, 3							

Last modified: Mon Feb 11 01:05:40 2019

CS164: Lecture #7 12

Ambiguous Example

Grammar

$p ::= e \text{ '+'}$
 $e ::= I \mid e \text{ '+' } e$

Input String

$I + I + I \text{ -}$

Chart. Only useful items shown.

0	I	1	+	2	I	3	+
a.p: $\bullet e \text{ '+'}$, 0	c.e: $I \bullet$, 0	b.e: $e \text{ '+'} \bullet e$, 0	d.e: $I \bullet$, 2	b.e: $e \text{ '+'} e \bullet$, 0	d.e: $\bullet I$, 2	b.e: $e \text{ '+'} e \bullet$, 0	
b.e: $\bullet e \text{ '+' } e$, 0	b.e: $e \bullet \text{ '+' } e$, 0	d.e: $\bullet I$, 2	e.e: $\bullet e \text{ '+' } e$, 2	e.e: $e \bullet \text{ '+' } e$, 2	e.e: $e \bullet \text{ '+' } e$, 2	b.e: $e \bullet \text{ '+' } e$, 0	
c.e: $\bullet I$, 0							
4	I	5	-	6			
b.e: $e \text{ '+'} \bullet e$, 0	f.e: $I \bullet$, 4	a.p: $e \text{ -} \bullet$, 0					
e.e: $e \text{ '+'} \bullet e$, 2	b.e: $e \text{ '+'} e \bullet$, 0						
f.p: $\bullet I$, 4	e.e: $e \text{ '+'} e \bullet$, 2						
	a.p: $e \bullet \text{-}$, 0						

Adding Semantic Actions

- Using syntax-directed translation to get semantic values is pretty much like recursive descent.
- The call $\text{parse}(A: \alpha \bullet \beta, s, k)$ can return, in addition to j , the semantic value of the A that matches symbols $c_{s+1} \dots c_j$.
- The value is computed during calls of the form $\text{parse}(A: \alpha' \bullet, s, k)$ (i.e., where the β part is empty). For terminal symbols, value is provided by the lexer.

Adding Semantic Actions (II)

- On a chart, when we see an item $A: \alpha \bullet, s$ in column k , it tells us to
 - Perform the semantic action corresponding to the production $A ::= \alpha$, getting a semantic value v for the left-hand side A .
 - For each item $B: \beta \bullet A \gamma, t$ in column s of the chart, when adding the item $B: \beta A \bullet \gamma, t$ to column k , also attach value v to that instance of A in the new item.
 - For all items derived from $B: \beta \bullet A \gamma, t$ as its dot is shifted, also attach v to the same instance of A .

This step is what provides the values of nonterminals needed to compute v values (in Bison notation: $\$1, \2 , etc.; in CUP notation, labels such as $e1$ and $e2$ in the rule $e ::= e: e1 \text{ '+' } e: e2$).

Example with Semantic Values

Grammar

$p: e: a \text{ '+'}$ $\{ : \text{ RESULT} = a; : \}$
 $e: t: b$ $\{ : \text{ RESULT} = b; : \}$
 $e: e: a \text{ '+' } t: b$ $\{ : \text{ RESULT} = a + b; : \}$
 $t: I: a$ $\{ : \text{ RESULT} = a; : \}$
 $t: t: a \text{ '*' } I: b$ $\{ : \text{ RESULT} = a * b; : \}$

Input String

(I's are numerals).
 $1 + 3 * 2 \text{ -}$

Chart. Only useful items shown. Semantic values are subscripts; red items show where they are computed.

0	I ₁	1	+	2	I ₃	3	*
a.p: $\bullet e \text{ '+'}$, 0	d.t ₁ : $I_1 \bullet$, 0	b.e: $e_1 \text{ '+'} \bullet t$, 0	e.t ₃ : $I_3 \bullet$, 2	b.e: $e_1 \text{ '+'} t \bullet$, 0	e.t: $\bullet I$, 2	f.t: $t_3 \bullet \text{ '*' } I$, 2	
b.e: $\bullet e \text{ '+' } t$, 0	c.e ₁ : $t_1 \bullet$, 0	b.e: $e_1 \bullet \text{ '+' } t$, 0	f.t: $\bullet t \text{ '*' } I$, 2				
c.e: $\bullet t$, 0							
d.t: $\bullet I$, 0							
4	I ₂	5	-	6			
f.t: $t_3 \text{ '*' } \bullet I$, 2	f.t ₆ : $t_3 \text{ '*' } I_2 \bullet$, 2	a.p ₇ : $e_7 \text{ -} \bullet$, 0					
	b.e ₇ : $e_1 \text{ '+' } t_6 \bullet$, 0						
	a.p: $e_7 \bullet \text{-}$, 0						

Handling Ambiguity in Semantics (Sketch)

- Ambiguity really only important here when computing semantic actions.
- Rather than being satisfied with a single path through the chart, we look at *all* paths.
- The call `parse(A: $\alpha \bullet \beta$, s, k)` can return a *set* of semantic values.
- Accordingly, we attach sets of semantic values to nonterminals.