

Recursive-Descent Parsing

First, a digression on lexing

Let's assume the **get-token** function will give me the next token

```
(define lex
  (lexer
    ; skip spaces:
    [#\space      (lex input-port)]
    ; skip newline:
    [#\newline   (lex input-port)]

    [#\+         'plus]
    [#\-         'minus]
    [#\*         'times]
    [#\/         'div]

    [(:: (?: #\-) (:+ (char-range #\0 #\9)))
     (string->number lexeme)]
    ; an actual character:
    [any-char    (string-ref lexeme 0)]))
```

Assume current token is `curtok`

`(accept c)` matches character `c`

```
(define curtok (next-tok))
```

```
(define (accept c)  
  (if (not (equal? curtok c))  
      (raise 'unexpected-token)  
      (begin  
        (printf "Accepting ~a\n" c)  
        (set! curtok (next-tok)))))
```

L left to right

L left derivation

1 token of lookahead

Let's say I want to parse the following grammar

$$S \rightarrow aSa \mid bb$$

First, a few questions

$$S \rightarrow aSa \mid bb$$

Is this grammar ambiguous?

If I were matching the string **bb**, what would my derivation look like?

If I were matching the string **abba**, what would my derivation look like?

First, a few questions

$$S \rightarrow aSa \mid bb$$

Key idea: if I look at the next input, at most one of these productions can “fire”

If I see an a I know that I must use the first production

If I see a b, I know I must be in second production

This is called a **predictive** parser. It uses lookahead to determine which production to choose

(My friend Tom points out that **predictive** is a dumb name because it is really “determining”, no guess)

In this class, we'll restrict ourselves to grammars that require only **one** character of lookahead

Generalizing to k characters is straightforward

I need two characters of lookahead

$S \rightarrow aaS \mid abS \mid c$

I need three characters of lookahead

$S \rightarrow aaaS \mid aabS \mid c$

I need four characters of lookahead

$S \rightarrow aaaaS \mid aaabS \mid c$

...

Slight transformation..

$S \rightarrow A \mid B$

$A \rightarrow aSa$

$B \rightarrow bb$

Slight transformation..

$S \rightarrow A \mid B$

$A \rightarrow aSa$

$B \rightarrow bb$

Now, I write out **one function** to parse **each** nonterminal

$S \rightarrow A \mid B$

$A \rightarrow aSa$

$B \rightarrow bb$

Intuition: when I see **a**, I call parse-A

when I see **b**, I call parse-B

```
(define (parse-A)
  (match curtok
    [#\a
     (begin
      (accept #\a)
      (parse-A)
      (accept #\a))]
    [#\b (parse-B)]))
```



```
(define (parse-B)
  (begin
    (accept #\b)
    (accept #\b)))
```

Livecoding this parser in class

Three parsing-related pieces of trivia

FIRST(A)

FIRST(A) is the **set** of terminals that could occur **first** when I recognize A

NULLABLE

Is the set productions which could generate ϵ

FOLLOW(A)

FOLLOW(A) is the set of terminals that appear immediately to the right of A in some form

Why learn these?

A: They help your intuition for building parsers
(as we'll see)

$S \rightarrow A \mid B$

$A \rightarrow aAa$

$B \rightarrow bb$

What is FIRST for each nonterminal

What is NULLABLE for the grammar

What is FOLLOW for each nonterminal

More practice...

$$E \rightarrow TE'$$

$$E' \rightarrow +TE'$$

$$E' \rightarrow \varepsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT'$$

$$T' \rightarrow \varepsilon$$

$$F \rightarrow (E)$$

$$F \rightarrow id$$

What is FIRST for each nonterminal

What is NULLABLE for the grammar

What is FOLLOW for each nonterminal

We use the **FIRST** set to help us design our recursive-descent parser!

LL(1)

A grammar is LL(1) if we only have to look at the **next** token to decide which production will match!

I.e., if $S \rightarrow A \mid B$, $\text{FIRST}(A) \cap \text{FIRST}(B)$ must be empty

Recursive-descent is called **top-down** parsing because you build a parse tree from the root down to the leaves

There are also **bottom-up** parsers,
which produce the rightmost derivation

Won't talk about them, in general they're
impossibly-hard to write / understand, easier to use

UNIX Programming Tools

2nd Edition



lex & yacc

O'REILLY®

*John R. Levine,
Tony Mason & Doug Brown*

Basically everyone uses lex and yacc to write real parsers

Recursive-descent is easy to implement, but requires lots of messing around with grammar

More practice with parsers

This one is more tricky!!

Plus \rightarrow num MoreNums

MoreNums \rightarrow + num MoreNums | ϵ

How would you do it?

(Hint: Think about NULLABLE)

Code up
collectively.....

```
(define (parse-Plus)
  (begin
    (parse-num)
    (parse-MorePlus)))
```

```
(define (parse-MorePlus)
  (match curtok
    ['plus
     (begin
       (accept 'plus)
       (parse-num)
       (parse-MorePlus))]
    ['eof (void)]))
```

Key rule: At each step of the way, if I see some token next, what rule production **must** I choose

Now yet another....

This will use the intuition from
FOLLOW

Add \rightarrow Term MoreTerms

MoreTerms \rightarrow + Term MoreTerms

MoreTerms \rightarrow ϵ

Term \rightarrow num MoreNums

MoreNums \rightarrow * num MoreNums | ϵ

Consider how we would implement MoreTerms

Add \rightarrow Term MoreTerms

MoreTerms \rightarrow + Term MoreTerms

MoreTerms \rightarrow ϵ

Term \rightarrow num MoreNums

MoreNums \rightarrow * num MoreNums | ϵ

If you're at the beginning of MoreTerms you **have** to see a +

Add \rightarrow Term MoreTerms

MoreTerms \rightarrow + Term MoreTerms

MoreTerms \rightarrow ϵ

Term \rightarrow num MoreNums

MoreNums \rightarrow * num MoreNums | ϵ

If you've just seen a + you have to see FIRST(Term)

Add \rightarrow Term MoreTerms

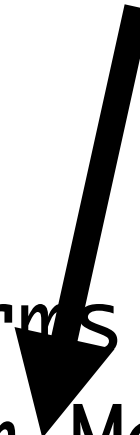
MoreTerms \rightarrow + Term MoreTerms

MoreTerms \rightarrow ϵ

Term \rightarrow num MoreNums

MoreNums \rightarrow * num MoreNums | ϵ

After Term you recognize something in FOLLOW(Term)



Add \rightarrow Term MoreTerms

MoreTerms \rightarrow + Term MoreTerms

MoreTerms \rightarrow ϵ

Term \rightarrow num MoreNums

MoreNums \rightarrow * num MoreNums | ϵ

Because MoreTerms is NULLABLE, have to account for null

Add \rightarrow Term MoreTerms

MoreTerms \rightarrow + Term MoreTerms

MoreTerms \rightarrow ϵ

Term \rightarrow num MoreNums

MoreNums \rightarrow * num MoreNums | ϵ

Code up
collectively.....

Let's say I want to generate an AST

Model my AST...

```
(struct add (left right) #:transparent)  
(struct times (left right) #:transparent)
```

More Recursive-descent practice...

Write recursive-descent parsers for the following....

A grammar for S-Expressions

Parsing mini-Racket / Scheme

```
datum ::= number
       | string
       | identifier
       | 'SExpr
SExpr ::= (SExprs)
       | datum
SExprs ::= SExpr SExprs
        | ε
```

$S \rightarrow a C H I b H C$

$H \rightarrow b H I d$

$C \rightarrow e C I f C$

$E \rightarrow A$

$E \rightarrow L$

$A \rightarrow n$

$A \rightarrow i$

$L \rightarrow (S)$

$S \rightarrow E S'$

$S' \rightarrow , S$

$S' \rightarrow \epsilon$

So far, I've given you grammars that are amenable to LL(1) parsers...

(Many grammars are **not**)

(But you can manipulate them to be!)

What about this grammar?

$$E \rightarrow E - T \mid T$$
$$T \rightarrow \text{number}$$

This grammar is **left recursive**

$$E \rightarrow E - T \mid T$$
$$T \rightarrow \text{number}$$

What happens if we try to write recursive-descent parser?

This grammar is **left recursive**

$$E \rightarrow E - T \mid T$$
$$T \rightarrow \text{number}$$

We really **want** this grammar, because it corresponds to the **correct** notion of associativity

$E \rightarrow E - T \mid T$

$T \rightarrow \text{number}$

5 - 3 - 1

Infinite loop!

$E \rightarrow E - T \mid T$

$T \rightarrow \text{number}$

5 - 3 - 1

A recursive descent parser will first call parse-E

And then crash

$E \rightarrow E - T \mid T$

$T \rightarrow \text{number}$

5 - 3 - 1

Draw the **rightmost derivation** for this string

If we could only have the **rightmost** derivation, our problem would be solved

The problem is, a recursive-descent parser needs to look at the **next input immediately**

Recursive descent parsers work by looking at the next token and making a decision / prediction

Rightmost derivations require us to delay making choices about the input until later

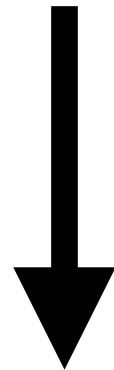
As humans, **we** naturally guess which derivation to use (for small examples)

Thus, LL(k) parsers cannot generate rightmost derivations :(

We can remove left recursion

$E \rightarrow E - T \mid T$

$T \rightarrow \text{number}$



Factor!

$E \rightarrow T E'$

$E' \rightarrow - T E'$

$E' \rightarrow \epsilon$

In general, if we have

$$A \rightarrow Aa \mid bB$$

Rewrite to...

$$A \rightarrow bB A'$$

$$A' \rightarrow a A' \mid \varepsilon$$

Generalizes even further

https://en.wikipedia.org/wiki/LL_parser#Left_Factoring

But this still doesn't give us what we want!!!

$$E \rightarrow T E'$$

$$E' \rightarrow - T E'$$

$$E' \rightarrow \varepsilon$$

$$E \rightarrow T E'$$

$$\rightarrow T - T E'$$

$$\rightarrow T - T - T E'$$

$$\rightarrow T - T - T$$

So how do we get left associativity?

Answer: Basically, hack in implementation

Sub \rightarrow num Sub'

Sub' \rightarrow + num Sub' | epsilon

Is basically...

Sub \rightarrow num Sub' (+ num)*

Intuition: treat this as while loop, then when building parse tree, put in left-associative order

$$\text{Sub} \rightarrow \text{num Sub}' (+ \text{num})^*$$

Sub \rightarrow num Sub'

Sub' \rightarrow + num Sub' | epsilon

If you want to get **rightmost** derivation, you need to use an LR parser

```

input:    /* empty */
         | input line
;

line:    '\n'
         | exp '\n' { printf ("\t%.10g\n", $1); }
;

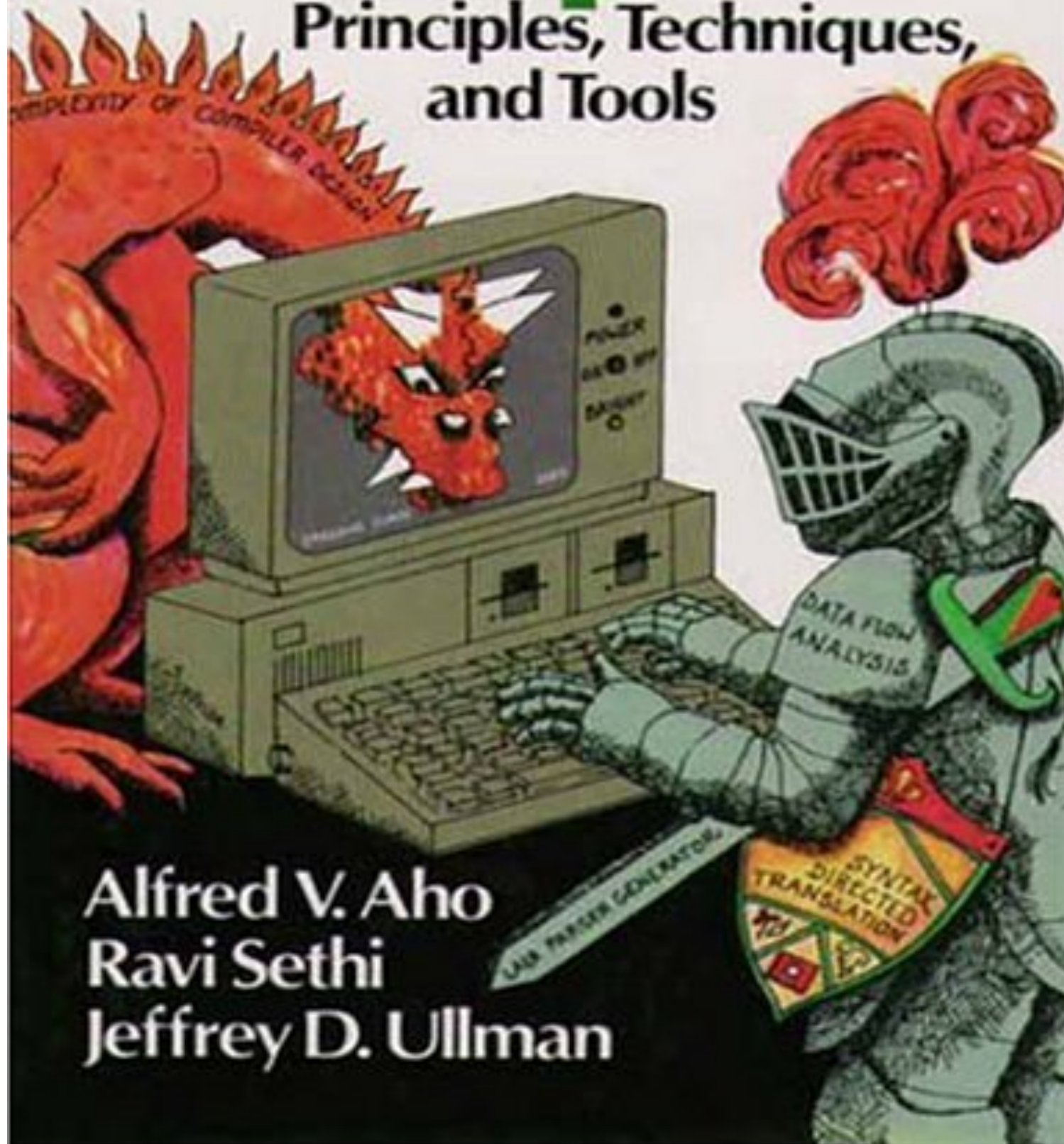
exp:     NUM          { $$ = $1;          }
         | exp exp '+' { $$ = $1 + $2;    }
         | exp exp '-' { $$ = $1 - $2;    }
         | exp exp '*' { $$ = $1 * $2;    }
         | exp exp '/' { $$ = $1 / $2;    }
         /* Exponentiation */
         | exp exp '^' { $$ = pow ($1, $2); }
         /* Unary minus */
         | exp 'n'    { $$ = -$1;          }
;

```

Parsing is lame, it's 2017

Compilers

Principles, Techniques,
and Tools



Alfred V. Aho
Ravi Sethi
Jeffrey D. Ullman

If you can, just use something like JSON /
protobufs / etc...

Inventing your own format is probably wrong

For small / prototypical things, recursive-descent

For real things, use yacc / bison / ANTLR