# Parsing

## RE macros:

$$digits = [0-9]+$$

$$sum = (digits\ "+")^* digits$$

$$7 + 11 + 42$$

How about:

$$digits = [0-9]+$$

$$sum = (expr\ "+")^* expr$$

$$expr = digits\ |\ "(" sum ")"$$

$$7 + (11 + 42)$$

not regular language

# Context - Free Grammar

Functionality:

  - similar to RE
  - plus recursion

Top-level alternation:

$$expr = a\ b\ (c|d)\ e$$

$$\Downarrow$$

$$aux = c|d$$
$$expr = a\ b\ aux\ e$$

$$\Downarrow$$

$$aux = c$$
$$aux = d$$
$$expr = a\ b\ aux\ e$$

# CFG

Recursion instead of
Kleene closure

$$expr = (a\ b\ c)*$$

$$\Downarrow$$

$$expr = (a\ b\ c)\ expr$$

$$expr = \varepsilon$$

# Context-Free Grammar

- terminal symbols (tokens)

- non-terminal symbols

- start symbol

- rules of the form

$$N \rightarrow X^*$$ where

N nonterminal

X (non-)terminal

## Example:

digit = 0
$$\vdots$$
digit = 9

digits = digit
digits = digit digits

sum = expr "+" expr sum
sum = expr
expr = digits
expr = "(" sum ")"

# Regular Languages

## Right-recursive

$$N \to t$$

$$N \to t N$$

## Left-recursive

$$N \to t$$

$$N \to N t$$

# Example

$$S \rightarrow S ; S$$

$$S \rightarrow id := E$$

$$S \rightarrow print(L)$$

$$E \rightarrow id$$

$$E \rightarrow num$$

$$E \rightarrow E + E$$

$$E \rightarrow (S, E)$$

$$L \rightarrow E$$

$$L \rightarrow L, E$$

## Example (cont.)

$$L \rightarrow E$$
$$L \rightarrow L, E$$

left-recursive

or

$$L \rightarrow E$$
$$L \rightarrow E, L$$

right-recursive

or

$$L \rightarrow E$$
$$L \rightarrow L, L$$

ambiguous

# Styles of rules

left-recursive

doesn't work with
top-down parsing
(e.g., hand-written parser)
$O(n^2)$

right-recursive

may be hard for
bottom-up parsing (e.g., CUP)
$O(n^2)$

ambiguous
are a pain
$O(n^3)$   Earley's algorithm
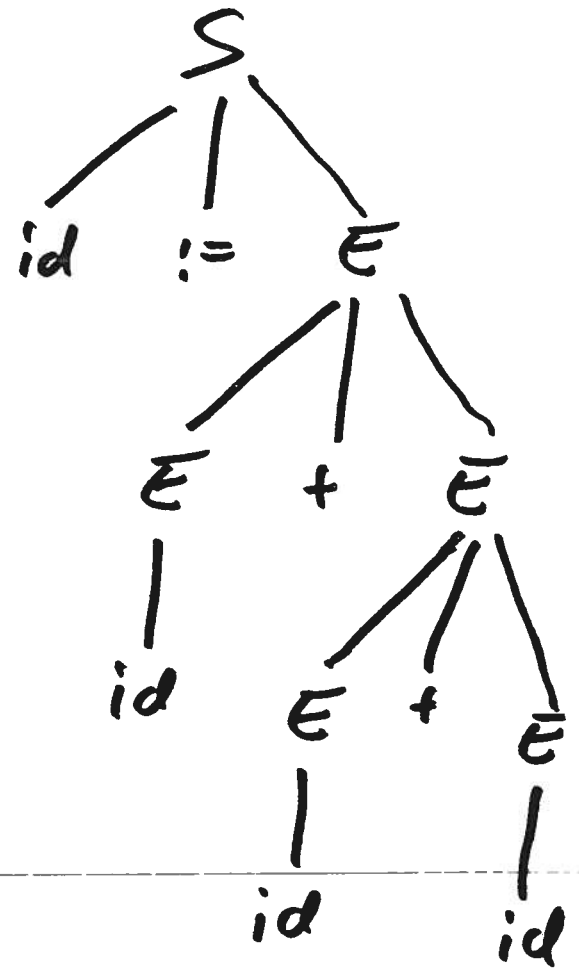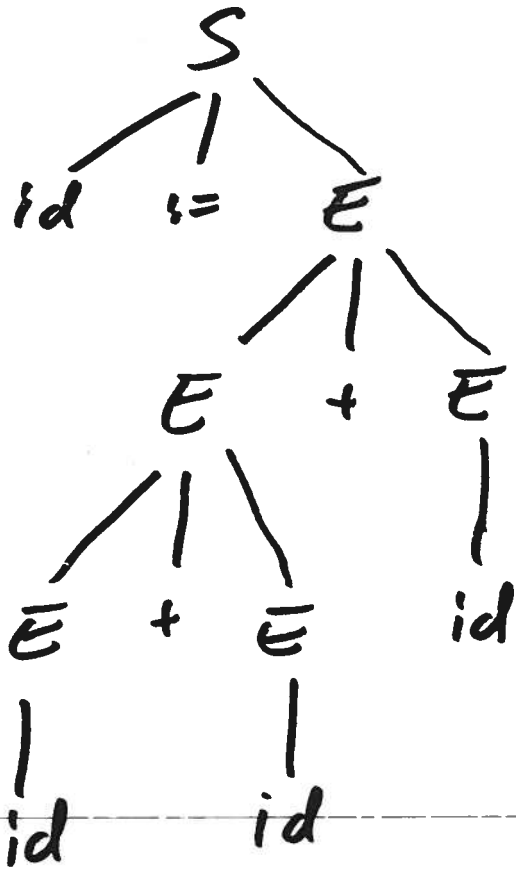
# Parser Types

LL(1)    — recursive descent

LR(1)

LL(k)    — Java CC, ANTLR

LALR(1)   — yacc, bison, CUP

# Parse Trees

id := id + id + id

# Derivations

$\underline{S}$

$S ; \underline{S}$

$\underline{S} ; id := E$

$id := \underline{E} ; \quad id = E$

$id := num ; \quad id := E$

left-most derivation

$\qquad = $ top-down parsing

right-most derivation

$\qquad = $ bottom-up parsing

# Abstract Parse Trees

```
            !=
           /   \
         id      +
                / \
              id    +
                   / \
                 id   id
```

# Ambiguous Grammars

$$E \rightarrow id$$
$$E \rightarrow num$$
$$E \rightarrow E \times E$$
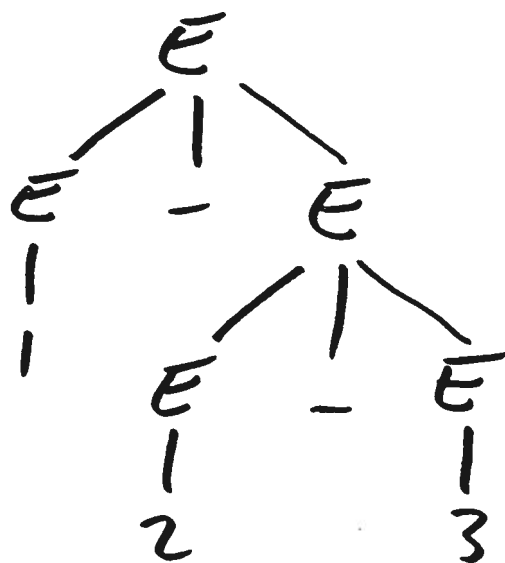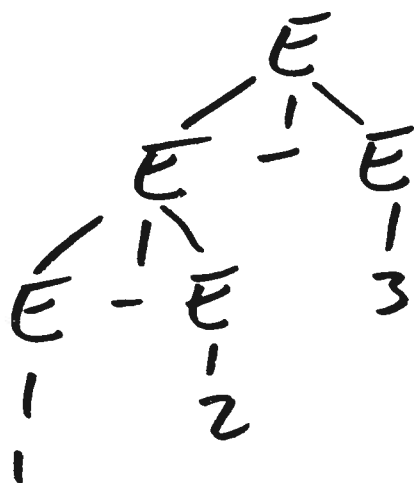$$E \rightarrow E / E$$
$$E \rightarrow E + E$$
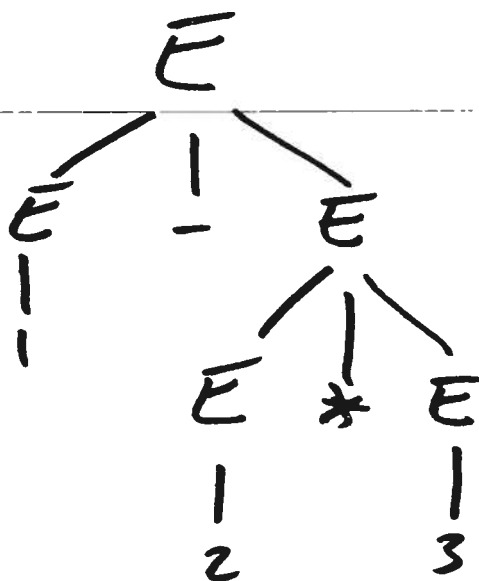$$E \rightarrow E - E$$
$$E \rightarrow (E)$$

# Ambiguous Grammars

## 1 − 2 − 3



## 1 − 2 * 3

# Modify Grammar

$$E \rightarrow E + T$$
$$E \rightarrow E - T$$
$$E \rightarrow T$$

$$T \rightarrow T * F$$
$$T \rightarrow T / F$$
$$T \rightarrow F$$

$$F \rightarrow id$$
$$F \rightarrow num$$
$$F \rightarrow (E)$$

# we don't get

# IF Statement

$S \rightarrow$ if $E$ then $S$ else $S$
$S \rightarrow$ if $E$ then $S$
$S \rightarrow$ other

$$\Downarrow$$

$S \rightarrow M$
$S \rightarrow U$
$M \rightarrow$ if $E$ then $M$ else $M$
$M \rightarrow$ other
$U \rightarrow$ if $E$ then $S$
$U \rightarrow$ if $E$ then $M$ else $U$

# CUP: Precedence Directives

```
precedence  nonassoc  EQ, NEQ;
precedence  left      PLUS, MINUS;
precedence  left      TIMES, DIV;
precedence  right     EXP;
precedence  left      UMINUS;
```

```
exp ::= INT
     |  exp PLUS exp
     |  exp MINUS exp
     |  exp TIMES exp
     |  MINUS exp
        %prec UMINUS
```

# LR - Parsing

## LR (k)

$k$ tokens lookahead

right-most derivations

left-to-right parse

## LALR (1)

lookahead

# Parse Engine

stack

DFA — applied to the stack

— edges labeled with (non) terminals

## Actions

$S_n$     shift and goto state $n$

$g_n$     goto state $n$

$r_k$     reduce by rule $k$

$a$     accept (shift EOF)

—     error

## Example (p. 58)

| Stack | Input | Action |
|---|---|---|
| ₁ | $a := 7\ \$$ | shift 4 |
| ₁ $id_4$ | $:= 7\ \$$ | shift 6 |
| ₁ $id_4 :=_6$ | $7\ \$$ | shift 10 |
| ₁ $id_4 :=_6 num_{10}$ | $\$$ | reduce $E \to num$ |
| ₁ $id_4 :=_6 E_{11}$ | $\$$ | reduce $S \to id := E$ |
| ₁ $S_2$ | $\$$ | accept |

$E \to num$

$E \to id$

$S \to id := E$

# Parsing Table

| | id | num | print | ; | , | + | != | ( | ) | $ | S | E | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | s4 | | s7 | | | | | | | | g2 | | |
| 2 | | | | | s3 | | | | | a | | g5 | |
| 3 | s4 | | s7 | | | | | | | | | g5 | |
| 4 | | | | | | | | s6 | | | | | |
| 5 | | | | | | r1 | r1 | | | r1 | | | |
| 6 | s20 | s10 | | | | | | | s8 | | | | g11 |
| 7 | | | | | | | | | | | | | |
| 8 | | | | | | | | | | | | | |
| 9 | | | | | | | | | | | | | |
| 10 | | | | | | r5 | r5 | r5 | r5 | r5 | | | |
| 11 | | | | | | r2 | r2 | s16 | | r2 | | | |

## Conflicts

shift-reduce      resolved by shift

reduce-reduce      resolved by order of rules

## CUP output

state 17:      shift/reduce conflict

         (shift ELSE, reduce 4)

stm ::= IF ID THEN stm .

stm ::= IF ID THEN stm . ELSE stm

ELSE      shift 19

.      reduce by rule 4

## Actions

shift (n)    - eat one token

             - shift state $n$ onto stack

reduce (k) - pop $n$ states off stack
                   where $n$ is #tokens on
                   RHS of rule $k$

             - in state on top of stack,
                   look up $X$ (LHS of rule $k$)
                   in goto (m)

             - push $m$ onto stack

accept    - stop, report success

error      - stop, report error

# Reduce / Reduce Conflicts

```
precedence    left OR;
precedence    left AND;
precedence    left PLUS;

stm  ::=    ID  ASSIGN  ae
      |     ID  ASSIGN  be;

be   ::=    be  OR  be
      |     be  AND  be
      |     ae  EQUAL  ae
      |     ID;

ae   ::=    ae  PLUS  ae
      |     ID;
```

> R/R conflict

# Reduce/Reduce Conflict

state 5:          R/R conflict
                  between rule 6 and 4
                  on EOF

        be ::=   ID.
        ae ::=   ID.
        PLUS     R6
        AND      R4
        OR       R4
        EQUAL    R6

        EOF      R4   ⟵
        .        error

# Solution: push decision to semantic analysis

```
stm  ::=  ID ASSIGN E;

E    ::=  E OR E
      |   E AND E
      |   E EQUAL E
      |   E PLUS E
      |   ID;
```

Solution: push decision to
scanner

---

g++ :  IDENT

TYPEIDENT

LABELIDENT

C x (a, b, c);

# Shift-Reduce Errors

use precedence declarations
for operators

  % precedence left    PLUS

  % precedence right  ASSIGN

  % precedence   nonassoc  EQUAL


use % prec if no appropriate
operators available

  % precedence left  UMINUS

  E ::=  MINUS E  %prec UMINU

# Error Recovery

On parse error:

- pop stack until there is a shift on error

- shift error token

- discard input until we are in a state with non-error action

- resume parsing

# Error Recovery

$E \rightarrow ID$

$E \rightarrow E + E$

$E \rightarrow (E)$

$E \rightarrow (error)$

$E_s \rightarrow E$

$E_s \rightarrow E_s; E$

$E_s \rightarrow error; E$

# Global Error Repair

Example:

let type a := intArray [10] of 0;
$\uparrow$

Solutions:

- delete type...0 (error prod.)

- replace := with =
  (local repair)

- replace type with var
  (global repair)

# Burke - Fisher Error Repair

- consider possible single token insertions/deletions/substitutions in last $K$ tokens ($K = 15$)

- use the repair that gets us the farthest, preferably at least $R$ tokens ($R = 4$)

# Grammar Rules for Error Reporting

Decl ::= Type ID LBRACK INT RBRACK
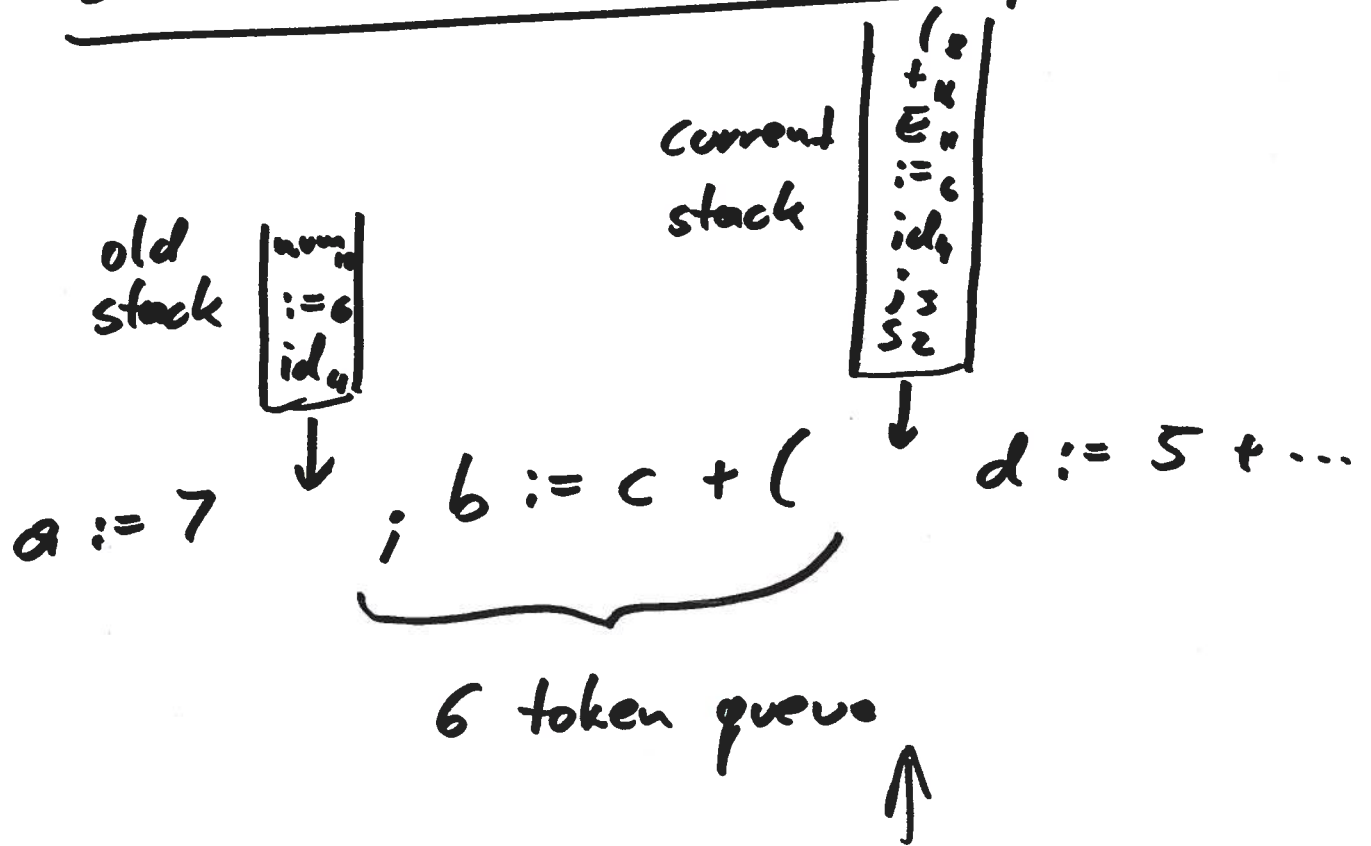ASSIGN LBRACE Explist RBRACE
SEM
{: ... :}

| ....

ASSIGN LBRACE RBRACE SEM
{: error ("...");
RESULT = NULL;
:}

# Burke - Fisher error repair

old stack

| num |
| :=6 |
| id₄ |

current stack

| ( ₈ |
| t ₄ |
| E ₇ |
| := ₆ |
| id₄ |
| ; ₃ |
| S ₂ |

$$a := 7 \quad ; \ b := c + ( \quad d := 5 + \ldots$$

6 token queue

## Cost

for window size $k$ and $N$ tokens

$$k + N \cdot k + N \cdot k$$

deletions

insertions

substitutions

# ML-Yacc

Semantic actions for Insertions

   % value ID ("bogus")

   % value INT (1)

   % value STRING ("")


Programmer-specified substitutions

   % change EQ → ASSIGN

        | ASSIGN → EQ

        | SEM ELSE → ELSE

        | → IN INT END

# Resolving Shift-Reduce Conflicts

Exp ::= Var
     | ID LBRACK Exp RBRACK OF Exp


Var ::= ID
     | Var LBRACK Exp RBRACK

# Example: S-R Conflict

```
Stm ::= VarDec
      | Assign;

VarDec ::= Type ID SEM;

Assign ::= LVal EQ Exp SEM;

Type ::= QualName
       | BuiltinType;

LVal ::= QualName
       | LVal LBRACK Exp RBRACK
       | LVal DOT ID;

QualName ::= ID
           | QualName DOT ID;
```

# LR Parsing

Grammar

↓

$\left.\begin{array}{c} \vdots \\ ? \\ \vdots \end{array}\right\}$ CUP

↓

Parse Tables
+
Shift-Reduce Parser

---

Read pp. 60-68

LR(0)
SCR
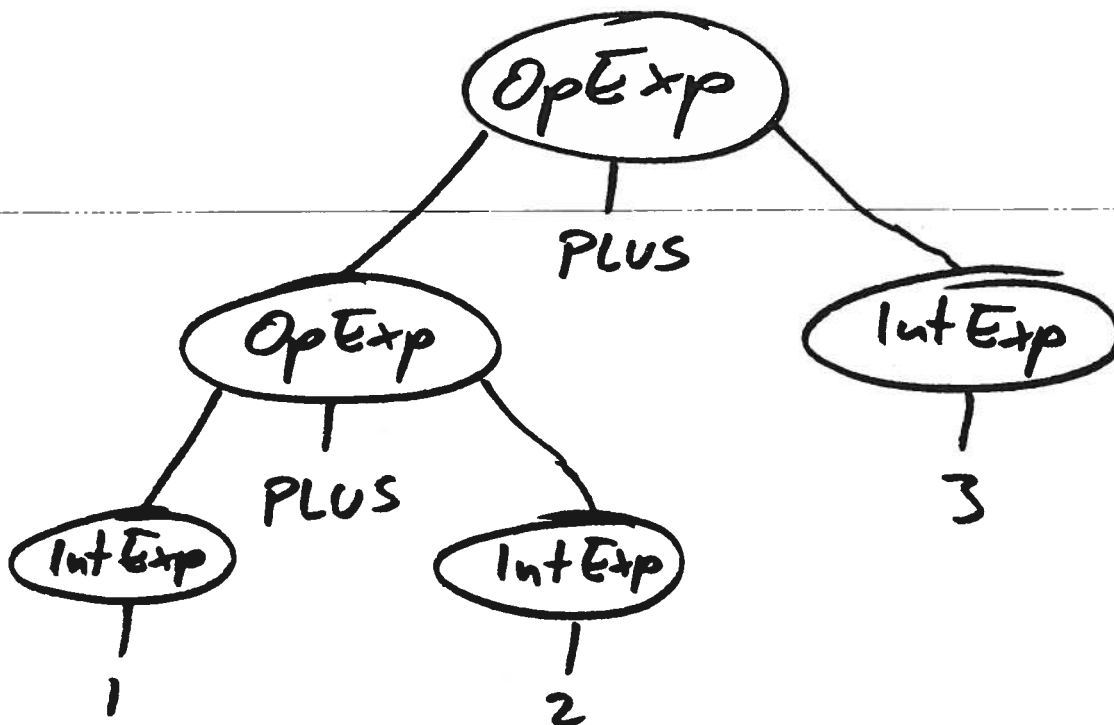→ LALR(1)
LR(1)

# Building Parse Trees

precedence left PLUS

Exp ::= INT:e
    {: RESULT = new IntExp (e); :}

  | Exp:e1 PLUS Exp:e2
    {: RESULT =
       new OpExp (e1, PLUS, e2); :}

---

Input: 1 + 2 + 3

# Building Parse Trees

```
Exp ::=   LPAREN Exp:e RPAREN
              {: RESULT = e; :}
          ;


FooList ::=
              {: RESULT = null; :}
          | Foos:l
              {: RESULT = l; :}
          ;
```

# LL - Parsing

## Recursive Descent

```
void S() {
  switch (tok) {
    case IF:
        eat(IF);
        E();
        eat(THEN);
        S();
        eat(ELSE);
        S();
        break;
    case BEGIN:
        :
        :
        :
        break;
    default: error();
  }
}
```

# Problems in LL Parsing

Left recursion:
$$X \rightarrow X\, y$$
$$\mid z$$

Common left factors:
$$X \rightarrow a\, Y$$
$$\mid a\, Z$$

Rules starting with nonterms:
$$X \rightarrow Y$$
$$\mid Z$$

Empty RHS:
$$X \rightarrow$$
$$\mid Y$$

# LL Parsing

- must chose alternatives (rules) based on lookahead

- must know FIRST, FOLLOW for each rule

# Nullable, FIRST, FOLLOW

## Nullable:

can $X$ derive emty string?

## FIRST:

set of terminals that can begin strings derived from $X$

## FOLLOW:

set of terminals that can follow $X$

## Example

$$Z \to d$$
$$| \ X \ Y \ Z$$

$$Y \to$$
$$| \ c$$

$$X \to Y$$
$$| \ a$$

# Example

| | nullable | FIRST | FOLLOW |
|---|---|---|---|
| X | yes | a c | a c d |
| Y | yes | c | a c d |
| Z | no | a c d | |

# Construction of Predictive Parser

Enter production $X \to p$

in row $X$, column $T$

for each $T \in FIRST(p)$;

if $p$ is nullable, enter
production in row $X$, col. $T$

for each $T \in FOLLOW(X)$

## Example

|   | a | c | d |
|---|---|---|---|
| X | $X \to a$ <br> $X \to Y$ | $X \to Y$ | $X \to Y$ |
| Y | $Y \to$ | $Y \to$ <br> $Y \to c$ | $Y \to$ |
| Z | $Z \to XYZ$ | $Z \to XYZ$ | $Z \to d$ <br> $Z \to XYZ$ |

not LL(1)

# Summary: Parsing

- recognize context-free syntactic structures on token stream

- syntax description using grammar

- bottom-up parsing
  - grammar $\xrightarrow{CUP}$ table-driven parser
  - uses push-down automaton
  - LR(0), SLR, <u>LALR(1)</u>, LR(1), LALR

- top-down parsing
  - recursive descent (hand-written)
  - predictive parser (tool)
  - LL(1), LL(k) (JavaCC)

- one-pass vs. multi-pass