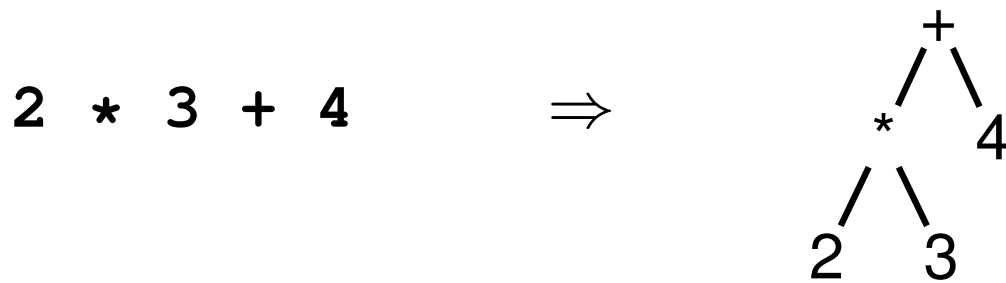# Parsing

# Parsing

Objective: build an abstract syntax tree (AST) for the token sequence from the scanner.

```
2 * 3 + 4        ⇒            +
                            / \
                           *   4
                          / \
                         2   3
```

Goal: discard irrelevant information to make it easier for the next stage.

Parentheses and most other forms of punctuation removed.

# Grammars

Most programming languages described using a
*context-free grammar*.

Compared to regular languages, context-free languages
add one important thing: recursion.

Recursion allows you to count, e.g., to match pairs of
nested parentheses.

Which languages do humans speak? I'd say it's regular: I
do not not not not not not not not not not understand this
sentence.

# Languages

Regular languages ($t$ is a terminal):

$$A \to t_1 \ldots t_n B$$

$$A \to t_1 \ldots t_n$$

Context-free languages ($P$ is terminal or a variable):

$$A \to P_1 \ldots P_n$$

Context-sensitive languages:

$$\alpha_1 A \alpha_2 \to \alpha_1 B \alpha_2$$

"$B \to A$ only in the 'context' of $\alpha_1 \cdots \alpha_2$"

# Issues

Ambiguous grammars

Precedence of operators

Left- versus right-recursive

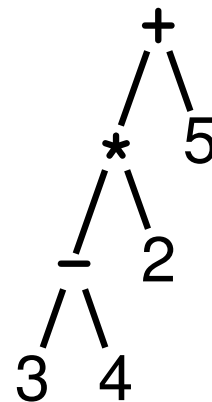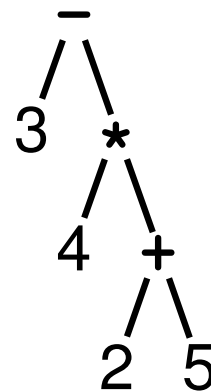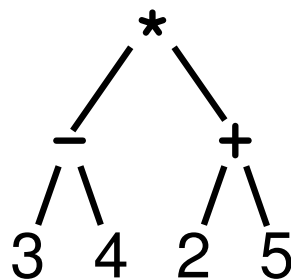Top-down vs. bottom-up parsers

Parse Tree vs. Abstract Syntax Tree

# Ambiguous Grammars
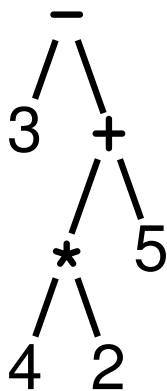
A grammar can easily be ambiguous. Consider parsing

$$3 - 4 * 2 + 5$$

with the grammar

$$e \rightarrow e + e \mid e - e \mid e * e \mid e / e \mid N$$

# Operator Precedence and Associativity

Usually resolve ambiguity in arithmetic expressions

Like you were taught in elementary school:

"My Dear Aunt Sally"

Mnemonic for multiplication and division before addition and subtraction.

# Operator Precedence

Defines how "sticky" an operator is.

`1 * 2 + 3 * 4`

`*` at higher precedence than `+`:

`(1 * 2) + (3 * 4)`

`+` at higher precedence than `*`:

`1 * (2 + 3) * 4`

# Associativity

Whether to evaluate left-to-right or right-to-left

Most operators are left-associative

`1 - 2 - 3 - 4`

```
            _                              _
           / \                            / \
          _   4                          1   _
         / \                                / \
        _   3                              2   _
       / \                                    / \
      1   2                                  3   4
    ((1 - 2) - 3) - 4                  1 - (2 - (3 - 4))

      left associative                   right associative
```

# Fixing Ambiguous Grammars

Original ANTLR grammar specification

```
expr
    : expr '+' expr
    | expr '-' expr
    | expr '*' expr
    | expr '/' expr
    | NUMBER
    ;
```

Ambiguous: no precedence or associativity.

# Assigning Precedence Levels

Split into multiple rules, one per level

```
expr : expr '+' expr
     | expr '-' expr
     | term ;


term : term '*' term
     | term '/' term
     | atom ;


atom : NUMBER ;
```

Still ambiguous: associativity not defined

# Assigning Associativity

Make one side or the other the next level of precedence

```
expr : expr '+' term
     | expr '-' term
     | term ;


term : term '*' atom
     | term '/' atom
     | atom ;


atom : NUMBER ;
```

# Parsing Context-Free Grammars

There are $O(n^3)$ algorithms for parsing arbitrary CFGs, but most compilers demand $O(n)$ algorithms.

Fortunately, the LL and LR subclasses of CFGs have $O(n)$ parsing algorithms. People use these in practice.

# The CYK algorithm (Cocke-Younger-Kasami)

Inputs: a string $w = a_1 a_2 \ldots a_n$ and a grammar in Chomsky Normal Form (only productions $X \to a$ and $X \to YZ$)

Construct an $n \times n$ table $T$, where $T[i, j]$ is the set of nonterminals that generate the substring $a_i a_{i+1} \ldots a_j$. Using dynamic programming, each table entry can be filled in $O(n)$ time. The algorithm runs in $O(n^3)$ time.

Question: how to determine that $a_1 a_2 \ldots a_n$ is in the language?

# Parsing LL(k) Grammars

LL: Left-to-right, Left-most derivation

k: number of tokens to look ahead

Parsed by top-down, predictive, recursive parsers

Basic idea: look at the next token to predict which production to use

ANTLR builds recursive LL(k) parsers

Almost a direct translation from the grammar.

# Implementing a Top-Down Parser

```
stmt : 'if' expr 'then' expr
     | 'while' expr 'do' expr
     | expr ':=' expr ;
expr : NUMBER | '(' expr ')' ;
```

```
stmt() {
  switch (next-token) {
  case IF:
      match(IF); expr(); match(THEN); expr();          break;
  case WHILE:
      match(WHILE); expr(); match(DO); expr();          break;
  case NUMBER or LPAREN:
      expr(); match(COLEQ); expr();                     break;
}}
```

# Writing LL(k) Grammars

Cannot have left-recursion

```
expr : expr '+' term | term ;
```

becomes

```
AST expr() {
  switch (next-token) {
  case NUMBER : expr(); /* Infinite Recursion */
```

# Writing LL(1) Grammars

Cannot have common prefixes

```
expr : ID '(' expr ')'
     | ID '=' expr
```

becomes

```
expr() {
  switch (next-token) {
  case ID:
    match(ID); match(LPAR); expr(); match(RPAR);  break;
  case ID:
    match(ID); match(EQUALS); expr();                 break;
```

# Eliminating Common Prefixes

Consolidate common prefixes:

```
expr
    : expr '+' term
    | expr '-' term
    | term
    ;
```

becomes

```
expr
    : expr ('+' term | '-' term )
    | term
    ;
```

# Eliminating Left Recursion

Understand the recursion and add tail rules

```
expr
   : expr ('+' term | '-' term )
   | term
   ;
```

becomes

```
expr : term exprt ;
exprt : '+' term exprt
      | '-' term exprt
      | /* nothing */
      ;
```

# Using ANTLR's EBNF

ANTLR makes this easier since it supports `*` and `+`:

```
expr : expr '+' term
     | expr '-' term
     | term ;
```

becomes

```
expr : term ('+' term | '-' term)* ;
```

# The Dangling Else Problem

Who owns the *else*?

```
if (a) if (b) c(); else d();
```



Grammars are usually ambiguous; manuals give disambiguating rules such as C's:

> As usual the "else" is resolved by connecting an else with the last encountered elseless if.

# The Dangling Else Problem

```
stmt : "if" expr "then" stmt iftail
     | other-statements ;


iftail
   : "else" stmt
   | /* nothing */
   ;
```

Problem comes when matching "iftail."

Normally, an empty choice is taken if the next token is in the "follow set" of the rule. But since "else" can follow an iftail, the decision is ambiguous.

# The Dangling Else Problem

ANTLR can resolve this problem by making certain rules "greedy." If a conditional is marked as greedy, it will take that option even if the "nothing" option would also match:

```
stmt
    : "if" expr "then" stmt
        ( options {greedy = true;}
        : "else" stmt
        )?
    | other-statements
    ;
```

# The Dangling Else Problem

Some languages resolve this problem by insisting on nesting everything.

E.g., Algol 68:

```
if a < b then a else b fi;
```

"fi" is "if" spelled backwards. The language also uses do–od and case–esac.

# Statement separators/terminators

C uses ; as a statement terminator.

```
if (a<b) printf("a less");
else {
  printf("b"); printf(" less");
}
```

Pascal uses ; as a statement separator.

```
if a < b then writeln('a less')
else begin
  write('a'); writeln(' less')
end
```

Pascal later made a final ; optional.

# Table-driven
# Top-Down Parsing

# Nomenclature

$a, b, c, \ldots$ represent terminal symbols

$A, B, C, \ldots$ represent nonterminal symbols

$S$ represents the initial nonterminal symbol

$X, Y, Z, \ldots$ represent terminal or non-terminal symbols

$u, v, w, \ldots$ represent words of terminal symbols

$\alpha, \beta, \gamma, \ldots$ represent words of terminal and non-terminal symbols

$\Rightarrow_G$ represents a one-step derivation with grammar G

$\overset{*}{\Rightarrow}_G$ represents zero or more steps in a derivation

# First and follow

$$\text{First}(\alpha) \equiv$$

$$\{a : \alpha \overset{*}{\Rightarrow} a\beta\} \cup$$

$$(\textbf{if } \alpha \overset{*}{\Rightarrow} \varepsilon \textbf{ then } \{\varepsilon\} \textbf{ else } \emptyset)$$

$$\text{Follow}(A) \equiv$$

$$\{a : S \overset{+}{\Rightarrow} \alpha A a\beta\} \cup$$

$$(\textbf{if } S \overset{*}{\Rightarrow} \alpha A \textbf{ then } \{\varepsilon\} \textbf{ else } \emptyset)$$

# Calculating first and follow

FIRST:

    For all terminals $a$: FIRST$(a) = \{a\}$;

    For all nonterminals $X$: FIRST$(X) = \emptyset$;

    For all productions $X \to \varepsilon$, add $\varepsilon$ to FIRST$(X)$;

    **repeat**

        For all productions $X \to Y_1 Y_2 \ldots Y_k$      ⟨outer loop⟩

           For $i$ in $1 \ldots k$     ⟨inner loop⟩

               add $(\text{FIRST}(Y_i) \setminus \{\varepsilon\})$ to FIRST$(X)$;

               **if**  $\varepsilon \notin$ FIRST$(Y_i)$  **continue outer loop**;

           add $\varepsilon$ to FIRST$(X)$;

    **until** no further progress

# Calculating first and follow

FOLLOW:

   FOLLOW$(S) = \{\varepsilon\}$, where $S$ is the start symbol

   For all other symbols $X$, FOLLOW$(X) = \emptyset$

   **repeat**

      For all productions $A \to \alpha B \beta$

         add (FIRST$(\beta) \setminus \{\varepsilon\}$) to FOLLOW$(B)$

      For all productions $A \to \alpha B$

               or $A \to \alpha B \beta$, where $\varepsilon \in$ FIRST$(\beta)$

         add FOLLOW$(A)$ to FOLLOW$(B)$

   **until** no further progress

---

**Note:** FIRST$(\beta)$ is calculated in a similar way as the inner loop of FIRST.

# Table-driven predictive parser

Input | | | | | a | + | b | $

```
X
Y
Z
$
```
Stack

Predictive
Parsing
Program

Output

Parsing
Table
M

# Parsing table

For each production $A \rightarrow \alpha$ of the grammar do:

1. For each terminal $a \in \text{FIRST}(\alpha)$, add $A \rightarrow \alpha$ to $M[A, a]$.

2. If $\varepsilon \in \text{FIRST}(\alpha)$, then for each terminal $b$ in $\text{FOLLOW}(A)$, add $A \rightarrow \alpha$ to $M[A, b]$. If $\varepsilon \in \text{FIRST}(\alpha)$ and $\$ \in \text{FOLLOW}(A)$, add $A \rightarrow \alpha$ to $M[A, \$]$ as well.

All empty entries must be set to **error**.

# Parsing table

| Production | First | Follow |
|---|---|---|
| $E \rightarrow T E'$ | $\{(, \textbf{id}\}$ | $\{), \$\}$ |
| $E' \rightarrow +T E' \mid \varepsilon$ | $\{+, \varepsilon\}$ | $\{), \$\}$ |
| $T \rightarrow F T'$ | $\{(, \textbf{id}\}$ | $\{+, ), \$\}$ |
| $T' \rightarrow *F T' \mid \varepsilon$ | $\{*, \varepsilon\}$ | $\{+, ), \$\}$ |
| $F \rightarrow (E) \mid \textbf{id}$ | $\{(, \textbf{id}\}$ | $\{+, *, ), \$\}$ |

| | **id** | + | * | ( | ) , $ |
|---|---|---|---|---|---|
| $E$ | $E \rightarrow T E'$ | | | $E \rightarrow T E'$ | |
| $E'$ | | $E' \rightarrow +T E'$ | | | $E' \rightarrow \varepsilon$ |
| $T$ | $T \rightarrow F T'$ | | | $T \rightarrow F T'$ | |
| $T'$ | | $T' \rightarrow \varepsilon$ | $T' \rightarrow *F T'$ | | $T' \rightarrow \varepsilon$ |
| $F$ | $F \rightarrow \textbf{id}$ | | | $F \rightarrow (E)$ | |

# LL(1) grammars

A grammar is LL(1) if

- For every pair of rules $A \to \alpha_1$ and $A \to \alpha_2$, then $\text{FIRST}(\alpha_1) \cap \text{FIRST}(\alpha_2) = \emptyset$.

- If $A \to \alpha$ and $A$ is nullable, then $\text{FIRST}(\alpha) \cap \text{FOLLOW}(A) = \emptyset$.

If a grammar is LL(1), then no conflicts appear in the parsing table.

# Predicting parsing algorithm

*ip* is the pointer to the input (initially pointing at the first token);
The stack initially contains $S\$$, with $S$ on the top;
Set $X$ to the top of the stack symbol;
**while** $(X \neq \$)$ **do**
   $a = $ INPUT[*ip*];
   **if** $(X = a)$ pop the stack and advance *ip*;
   **else if** $(X$ is a terminal) error();
   **else if** $(M[X, a] = $ **error**$)$ error();
   **else if** $(M[X, a] = X \to Y_1 Y_2 \cdots Y_k)$
      Output the production $X \to Y_1 Y_2 \cdots Y_k$;
      Pop the stack;
      Push $Y_k, Y_{k-1}, \ldots, Y_1$ onto the stack, with $Y_1$ on top;
   Set $X$ to the top of the stack symbol;

---

**Exercise:** Execute the parsing algorithm with the input $z * (x + y)\$$.

# Generation of top-down parsers
# (ANTLR style)

# A simple example

```
instruction_list: ( instruction ) *

instruction: IDENT ":=" expr |
             IF expr THEN instruction_list

expr: (IDENT | NUM) (PLUS (IDENT | NUM))*
```

Parser

```
void instruction_list () {
  while (token==IDENT || token==IF) {
    instruction();
  }
}
```

# A simple example

```
instruction: IDENT ":=" expr |
             IF expr THEN instruction_list

void instruction () {
  if (token==IDENT) {
    nexttoken();
    if (token== COLON_EQUAL) {
      nexttoken(); expr();
    } else SYNTAXERROR();
  } else if (token==IF) {
    nexttoken(); expr();
    if (token==THEN) {
      nextoken(); instruction_list();
    } else SYNTAXERROR();
  } else SYNTAXERROR();
}
```

# A simple example

```
expr: (IDENT | NUM) (PLUS (IDENT | NUM))*
```

---

```
void expr() {
  if (token==IDENT || token==NUM) {
    nexttoken();
    while (token==PLUS) {
      nexttoken();
      if (token==IDENT || token==NUM){
        nexttoken();
      } else SYNTAXERROR();
    }
  } else SYNTAXERROR();
}
```

# ANTLR

- ANTLR uses the EBFN notation

- NULLABLE, FIRST and FOLLOW must be extended to the EBNF notation.

- A recursive top-down parser is generated.

# Nullable

- NULLABLE($\varepsilon$) = *true*.

- NULLABLE($E*$) = *true*.

- If NULLABLE($E$), then NULLABLE($E+$) = *true*.

- If $V \rightarrow E \in G$ and NULLABLE($E$), then NULLABLE($V$) = *true*.

- If NULLABLE($E_1$) and NULLABLE($E_2$), then NULLABLE($E_1\,E_2$) = *true*.

- If NULLABLE($E_1$) or NULLABLE($E_2$), then NULLABLE($E_1\,|\,E_2$) = *true*.

- Nothing else is NULLABLE.

# First

- If $c$ is a terminal, $\text{FIRST}(c) = \{c\}$.

- $\text{FIRST}(E*)$ and $\text{FIRST}(E+)$ contain $\text{FIRST}(E)$.

- If $V \to E \in G$, then $\text{FIRST}(V)$ contains $\text{FIRST}(E)$.

- $\text{FIRST}(E_1 E_2)$ contains $\text{FIRST}(E_1)$.

- If $\text{NULLABLE}(E_1)$, then $\text{FIRST}(E_1 E_2)$ contains $\text{FIRST}(E_2)$.

- $\text{FIRST}(E_1 | E_2)$ contains $\text{FIRST}(E_1)$ and $\text{FIRST}(E_2)$.

- Nothing else belongs to $\text{FIRST}$.

# Follow

- If $E_1\,E_2$ is and expression of the grammar, then $\textsc{Follow}(E_2)$ contains $\textsc{Follow}(E_1\,E_2)$, $\textsc{Follow}(E_1)$ contains $\textsc{First}(E_2)$ and, if $\textsc{Nullable}(E_2)$ then $\textsc{Follow}(E_1)$ contains $\textsc{Follow}(E_1\,E_2)$.

- If $E_1\,|\,E_2$ is and expression of the grammar, then $\textsc{Follow}(E_1)$ and $\textsc{Follow}(E_2)$ contain $\textsc{Follow}(E_1\,|\,E_2)$.

- If $E*$ or $E+$ are expressions of the grammar, then their $\textsc{Follow}$ is contained in $\textsc{Follow}(E)$.

- If $V \rightarrow E \in G$, then $\textsc{Follow}(E)$ contains $\textsc{Follow}(V)$.

- Nothing else belongs to $\textsc{Follow}$.

# Generating an LL(1) recursive-descent predictive parser

For every rule $A \to E$, a function is generated:

```
void A() {
    Parse(E, Follow(A))
}
```

where $\textbf{Parse}(E, F)$ is the code generated to recognize the EBNF expression $E$ followed by the tokens in $F$.

**Token** is a variable that represents the current token.

```
match(T) ≡
    if (Token == T) nexttoken();
    else SyntaxError();
```

# Parse(E,F)

```
Parse(E₁ | E₂ | ... | Eₙ, F) ≡
    if (Token ∈ First(E₁)) Parse(E₁,F);
    else if (Token ∈ First(E₂)) Parse(E₂,F);
    ...
    else if (Token ∈ First(Eₙ)) Parse(Eₙ,F);
    else if (no Eᵢ is nullable) SyntaxError();
    else if (Token ∉ F) SyntaxError();
```

$$\textbf{Parse(}E_1 \mid E_2 \mid \ldots \mid E_n, F\textbf{)} \equiv$$
$$\texttt{if (Token} \in \textbf{\textit{First}}(E_1)\texttt{)} \textbf{Parse(}E_1, F\textbf{);}$$
$$\texttt{else if (Token} \in \textbf{\textit{First}}(E_2)\texttt{)} \textbf{Parse(}E_2, F\textbf{);}$$
$$\ldots$$
$$\texttt{else if (Token} \in \textbf{\textit{First}}(E_n)\texttt{)} \textbf{Parse(}E_n, F\textbf{);}$$
$$\texttt{else if (no } E_i \texttt{ is nullable) SyntaxError();}$$
$$\texttt{else if (Token} \notin F\texttt{) SyntaxError();}$$

If the BNF version of the grammar is LL(1) then there are no conflicts between the different branches (including the case of a nullable $E_i$).

# Parse(E,F)

**Parse(**$E_1 E_2 \ldots E_n$**,** $F$**)** $\equiv$

   **Parse(**$E_1$**,** **$First$**$(E_2 \ldots E_n \cdot F)$**)**;

   **Parse(**$E_2$**,** **$First$**$(E_3 \ldots E_n \cdot F)$**)**;

   $\ldots$

   **Parse(**$E_n$**,** $F$**)**;


**where $First$$(E \cdot F)$ is**

- **$First$**$(E)$ **if** $E$ **is not nullable**

- **$First$**$(E) \cup F$**, otherwise**

# Parse(E,F)

```
Parse(E*,F) ≡
    while (Token ∈ First(E)) Parse(E,F);


Parse(E+,F) ≡
    do Parse(E,F); while (Token ∈ First(E,F));


Parse(E?,F) ≡
    if (Token ∈ First(E)) Parse(E,F);

Parse(a,F) ≡ Match(a);  // a is a terminal symbol


Parse(A,F) ≡ A();  // A is a non-terminal symbol
```

# LL(1) example

Grammar:

$$E \rightarrow T \, ('\mathtt{+}' \, T \mid '\mathtt{-}' \, T)*$$

$$T \rightarrow F \, ('\mathtt{*}' \, F \mid '\mathtt{/}' \, F)*$$

$$F \rightarrow \mathbf{ID} \mid \mathbf{NUM} \mid '\mathtt{(}' \, E \, '\mathtt{)}'$$

---

```
void E() {
    T();
    while (Token == '+' or Token == '-') {
        if (Token == '+') { Match('+'); T(); }
        else if (Token == '-') { Match('-'); T(); }
        else SyntaxError();      // redundant
    }
}
```

# LL(1) example

```
void T() {
    F();
    while (Token == '*' or Token == '/') {
        if (Token == '*') { Match('*'); F(); }
        else if (Token == '/') { Match('/'); F(); }
        else SyntaxError();        // redundant
    }
}

void F() {
    if (Token == ID) Match(ID);
    else if (Token == NUM) Match(NUM);
    else if (Token == '(') {
        Match('('); E(); Match(')');
    } else SyntaxError();
}
```

# Bottom-up Parsing

# Rightmost Derivation

$1:$    $e \rightarrow t + e$

$2:$    $e \rightarrow t$

$3:$    $t \rightarrow \textbf{Id} * t$

$4:$    $t \rightarrow \textbf{Id}$

A rightmost derivation for $\textbf{Id} * \textbf{Id} + \textbf{Id}$:

$e$

$t + e$

$t + t$

$t + \textbf{Id}$

$\textbf{Id} * t + \textbf{Id}$

$\textbf{Id} * \textbf{Id} + \textbf{Id}$

Basic idea of bottom-up parsing: construct this rightmost derivation backward.

# Handles

1 :  $e \rightarrow t + e$

2 :  $e \rightarrow t$

3 :  $t \rightarrow \mathbf{Id} * t$

4 :  $t \rightarrow \mathbf{Id}$

$\mathbf{Id} * \boxed{\mathbf{Id}} + \mathbf{Id}$

$\boxed{\mathbf{Id} * t} + \mathbf{Id}$

$t + \boxed{\mathbf{Id}}$

$t + \boxed{t}$

$\boxed{t + e}$

$e$



This is a reverse rightmost derivation for $\mathbf{Id} * \mathbf{Id} + \mathbf{Id}$.

Each highlighted section is a <span style="color:red">handle</span>.

Taken in order, the handles build the tree from the leaves to the root.

# Shift-reduce Parsing

| | | |
|---|---|---|
| $1:\quad e{\rightarrow}t+e$ | | |
| $2:\quad e{\rightarrow}t$ | | |
| $3:\quad t{\rightarrow}\mathbf{Id} * t$ | | |
| $4:\quad t{\rightarrow}\mathbf{Id}$ | | |

| stack | input | action |
|---|---|---|
| | $\mathbf{Id} * \mathbf{Id} + \mathbf{Id}$ | shift |
| $\mathbf{Id}$ | $* \mathbf{Id} + \mathbf{Id}$ | shift |
| $\mathbf{Id}*$ | $\mathbf{Id} + \mathbf{Id}$ | shift |
| $\mathbf{Id} * \mathbf{Id}$ | $+ \mathbf{Id}$ | reduce (4) |
| $\mathbf{Id} * t$ | $+ \mathbf{Id}$ | reduce (3) |
| $t$ | $+ \mathbf{Id}$ | shift |
| $t+$ | $\mathbf{Id}$ | shift |
| $t + \mathbf{Id}$ | | reduce (4) |
| $t + t$ | | reduce (2) |
| $t + e$ | | reduce (1) |
| $e$ | | accept |

Scan input left-to-right, looking for handles.

An oracle tells what to do

# LR Parsing

$1:$     $e \rightarrow t + e$

$2:$     $e \rightarrow t$

$3:$     $t \rightarrow \mathbf{Id} * t$

$4:$     $t \rightarrow \mathbf{Id}$

| stack | input | action |
|---|---|---|
| $\boxed{0}$ | **Id** * **Id** + **Id** $ | shift, goto 1 |

|   | action |   |   |   | goto |   |
|---|---|---|---|---|---|---|
|   | **Id** | $+$ | $*$ | \$ | $e$ | $t$ |
| 0 | s1 |    |    |    | 7 | 2 |
| 1 |    | r4 | s3 | r4 |   |   |
| 2 |    |    | s4 | r2 |   |   |
| 3 | s1 |    |    |    |   | 5 |
| 4 | s1 |    |    |    | 6 | 2 |
| 5 |    | r3 |    | r3 |   |   |
| 6 |    |    |    | r1 |   |   |
| 7 |    |    |    | acc |   |   |

1. Look at state on top of stack

2. and the next input token

3. to find the next action

4. In this case, shift the token onto the stack and go to state 1.

# LR Parsing

$1:\quad e{\rightarrow}t + e$

$2:\quad e{\rightarrow}t$

$3:\quad t{\rightarrow}\textbf{Id} \ast t$

$4:\quad t{\rightarrow}\textbf{Id}$

| | action | | | | goto | |
|---|---|---|---|---|---|---|
| | **Id** | $+$ | $\ast$ | $\$$ | $e$ | $t$ |
| 0 | s1 | | | | 7 | 2 |
| 1 | | r4 | s3 | r4 | | |
| 2 | | | s4 | | r2 | |
| 3 | s1 | | | | | 5 |
| 4 | s1 | | | | 6 | 2 |
| 5 | | | r3 | | r3 | |
| 6 | | | | r1 | | |
| 7 | | | | acc | | |

|stack | input | action |
|---|---|---|
| $\boxed{0}$ | **Id** \* **Id** + **Id** $\$$ | shift, goto 1 |
| $\boxed{0}\ \boxed{\textbf{Id}\,1}$ | \* **Id** + **Id** $\$$ | shift, goto 3 |
| $\boxed{0}\ \boxed{\textbf{Id}\,1}\ \boxed{\ast\,3}$ | **Id** + **Id** $\$$ | shift, goto 1 |
| $\boxed{0}\ \boxed{\textbf{Id}\,1}\ \boxed{\ast\,3}\ \boxed{\textbf{Id}\,1}$ | + **Id** $\$$ | reduce w/ 4 |

Action is reduce with rule 4 ($t \rightarrow$ **Id**). The right side is removed from the stack to reveal state 3. The goto table in state 3 tells us to go to state 5 when we reduce a $t$:

| stack | input | action |
|---|---|---|
| $\boxed{0}\ \boxed{\textbf{Id}\,1}\ \boxed{\ast\,3}\ \boxed{t\,5}$ | + **Id** $\$$ | |

# LR Parsing

$$1: \quad e \rightarrow t + e$$

$$2: \quad e \rightarrow t$$

$$3: \quad t \rightarrow \mathbf{Id} * t$$

$$4: \quad t \rightarrow \mathbf{Id}$$

|   | action | | | | goto | |
|---|---|---|---|---|---|---|
|   | **Id** | $+$ | $*$ | $\$$ | $e$ | $t$ |
| 0 | s1 |   |   |   | 7 | 2 |
| 1 |   | r4 | s3 | r4 |   |   |
| 2 |   | s4 |   | r2 |   |   |
| 3 | s1 |   |   |   |   | 5 |
| 4 | s1 |   |   |   | 6 | 2 |
| 5 |   | r3 |   | r3 |   |   |
| 6 |   |   |   | r1 |   |   |
| 7 |   |   |   | acc |   |   |

| stack | input | action |
|---|---|---|
| 0 | **Id** * **Id** + **Id** $\$$ | shift, goto 1 |
| 0 Id/1 | * **Id** + **Id** $\$$ | shift, goto 3 |
| 0 Id/1 */3 | **Id** + **Id** $\$$ | shift, goto 1 |
| 0 Id/1 */3 Id/1 | + **Id** $\$$ | reduce w/ 4 |
| 0 Id/1 */3 t/5 | + **Id** $\$$ | reduce w/ 3 |
| 0 t/2 | + **Id** $\$$ | shift, goto 4 |
| 0 t/2 +/4 | **Id** $\$$ | shift, goto 1 |
| 0 t/2 +/4 Id/1 | $\$$ | reduce w/ 4 |
| 0 t/2 +/4 t/2 | $\$$ | reduce w/ 2 |
| 0 t/2 +/4 e/6 | $\$$ | reduce w/ 1 |
| 0 e/7 | $\$$ | accept |

# Constructing the SLR Parse Table

The states are places we could be in a reverse-rightmost derivation. Let's represent such a place with a dot.

$1:\quad e{\to}t+e$

$2:\quad e{\to}t$

$3:\quad t{\to}\mathbf{Id}\ *\ t$

$4:\quad t{\to}\mathbf{Id}$

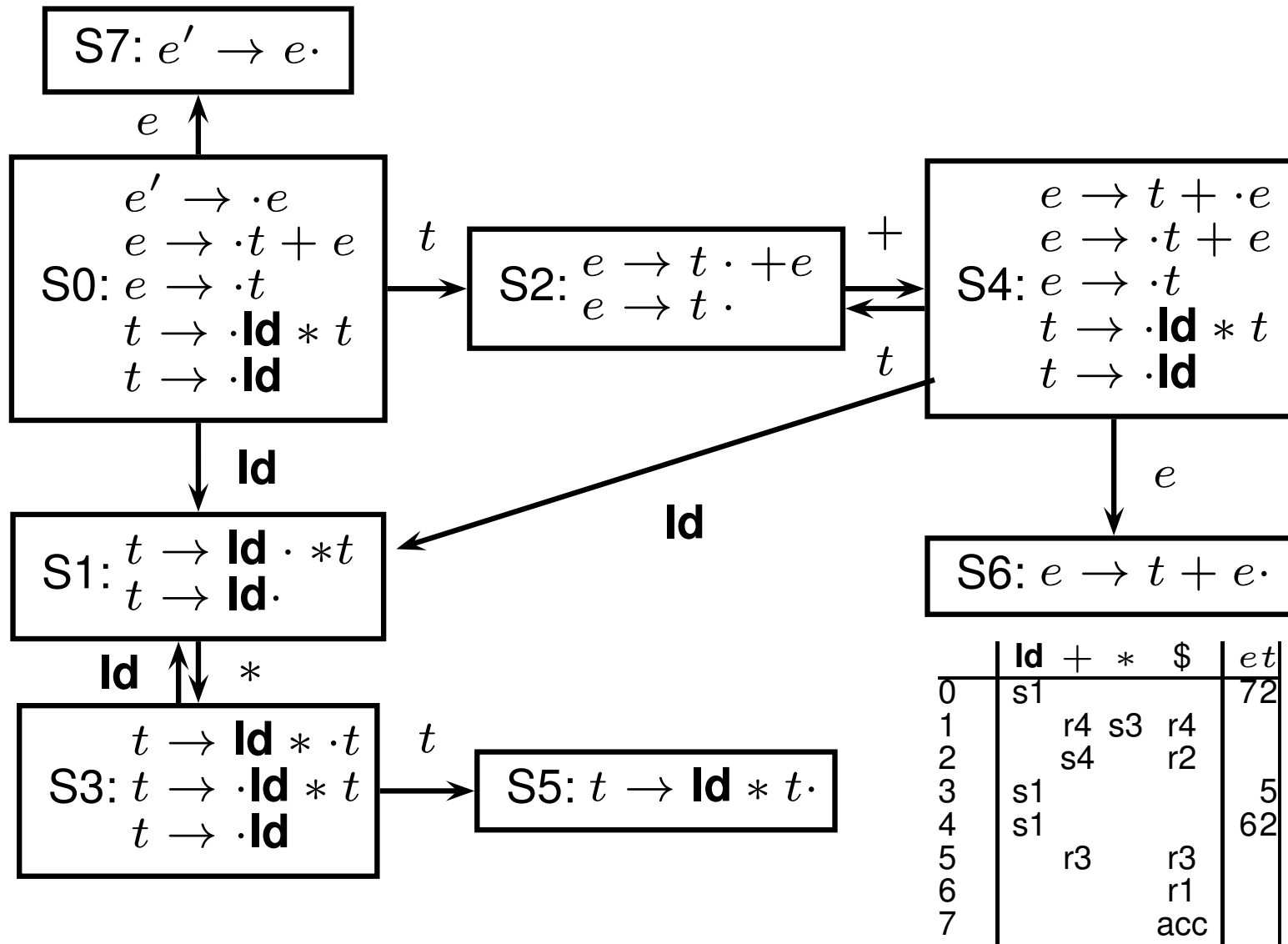Say we were at the beginning $(\cdot e)$. This corresponds to

$e' \to \cdot e$
$e \to \cdot t + e$
$e \to \cdot t$
$t \to \cdot \mathbf{Id} * t$
$t \to \cdot \mathbf{Id}$

The first is a placeholder. The second are the two possibilities when we're just before $e$. The last two are the two possibilities when we're just before $t$.

# Constructing the SLR Parsing Table

S7: $e' \to e\cdot$

$e' \to \cdot e$
$e \to \cdot t + e$
S0: $e \to \cdot t$
$t \to \cdot \mathbf{Id} * t$
$t \to \cdot \mathbf{Id}$

$t$

S2: $e \to t \cdot +e$
$e \to t \cdot$

$+$

$e \to t + \cdot e$
$e \to \cdot t + e$
S4: $e \to \cdot t$
$t \to \cdot \mathbf{Id} * t$
$t \to \cdot \mathbf{Id}$

$e$

$t$

$\mathbf{Id}$

S1: $t \to \mathbf{Id} \cdot *t$
$t \to \mathbf{Id}\cdot$

$e$

S6: $e \to t + e\cdot$

$\mathbf{Id}$

$*$

$t \to \mathbf{Id} * \cdot t$
S3: $t \to \cdot \mathbf{Id} * t$
$t \to \cdot \mathbf{Id}$

$t$

S5: $t \to \mathbf{Id} * t\cdot$

|   | $\mathbf{Id}$ | $+$ | $*$ | $\$$ | $e$ | $t$ |
|---|---|---|---|---|---|---|
| 0 | s1 |  |  |  | 7 | 2 |
| 1 |  | r4 | s3 | r4 |  |  |
| 2 |  | s4 |  | r2 |  |  |
| 3 | s1 |  |  |  |  | 5 |
| 4 | s1 |  |  |  | 6 | 2 |
| 5 |  | r3 |  | r3 |  |  |
| 6 |  |  |  | r1 |  |  |
| 7 |  |  |  | acc |  |  |

# The Punchline

This is a tricky, but mechanical procedure. The parser generators YACC, Bison, Cup, and others (but not ANTLR) use a modified version of this technique to generate fast bottom-up parsers.

You need to understand it to comprehend error messages:

Shift/reduce conflicts are caused by a state like

$t \to \textbf{Id} \cdot *t$

$t \to \textbf{Id} * t\cdot$

Reduce/reduce conflicts are caused by a state like

$t \to \textbf{Id} * t\cdot$

$e \to t + e\cdot$

# Exercises (grammars)

1. Write unambiguous grammars for the following languages:

   - The set of all strings of $a$'s and $b$'s that are palindromes.

   - Strings that match the pattern $a * b*$ and have more $a$'s than $b$'s.

   - Strings with balanced parenthesis and square braces. Example:

     `( [ [ ] ( ( ) [ ( ) ] [ ] ) ] )`

   - The set of all strings of $a$'s and $b$'s such that every $a$ is immediatelly followed by at least one $b$.

   - The set of all strings of $a$'s and $b$'s with an equal number of $a$'s and $b$'s.

   - The set of all strings of $a$'s and $b$'s with an different number of $a$'s and $b$'s.

   - Blocks of statements in Pascal or MH, where the semicolons ('$;$') separate the statements:

     `( statement;  ( statement ;  statement  ) ; statement )`

   - Blocks of statemens in C, where the semicolons ('$;$') follow each statement:

     `{ statement;  { statement;  statement; }  statement; }`

2. Specify the previous grammars in ANTLR notation, modifying the grammar when necessary.

# Exercises (parsing)

1. Calculate NULLABLE, FIRST and FOLLOW of the non-terminal symbols in the following grammar:

$$
\begin{aligned}
A &\rightarrow B \mid a \\
B &\rightarrow b \mid \varepsilon \\
C &\rightarrow c \mid A B C
\end{aligned}
$$

2. Consider the following grammar:

$$
S \rightarrow S\, S\, + \ \mid\ S\, S\, * \ \mid\ a
$$

and the string $a\, a\, +\, a *$.

- Give a leftmost derivation for the string.
- Give a rightmost derivation for the string.
- Give a parse tree for the string.
- Is the grammar ambiguous or unambiguous? Justify your answer.
- Describe the language generated by this grammar.

3. Consider the following grammar:

$$
\begin{aligned}
S &\rightarrow cABc \\
A &\rightarrow aAa \mid c \\
B &\rightarrow bBb \mid c
\end{aligned}
$$

- Calculate FIRST and FOLLOW for the non-terminal symbols.
- Construct the LL(1) parsing table and check whether it is an LL(1) grammar.

4. Calculate NULLABLE, FIRST and FOLLOW for the following grammar:

$$
\begin{aligned}
S &\rightarrow uBDz \\
B &\rightarrow Bv \mid w \\
D &\rightarrow EF \\
E &\rightarrow y \mid \varepsilon \\
F &\rightarrow x \mid \varepsilon
\end{aligned}
$$

Construct the LL(1) parsing table and give evidence that this grammar is not LL(1). Modify the grammar as little as possible to make an LL(1) grammar that accepts the same language.

5. Design a table-driven top-down parser for the following grammar:

$$
\begin{aligned}
S &\rightarrow E \,\$ \\
E &\rightarrow T + E \mid T \\
T &\rightarrow \text{num} * T \mid \text{num}
\end{aligned}
$$

6. Design a recursive-descent parser (ANTLR style) for the following grammar:

$$E \rightarrow T \ (' +' \ T) \ *$$

$$T \rightarrow F \ (' *' \ F) \ *$$

$$F \rightarrow ' (' \ E \ ')' \ | \ \mathbf{id}$$

7. Consider the following grammar:

$$G \rightarrow S \ \$$$

$$S \rightarrow A \ M$$

$$M \rightarrow S \ | \ \varepsilon$$

$$A \rightarrow a \ E \ | \ b \ A \ A$$

$$E \rightarrow a \ B \ | \ b \ A \ | \ \varepsilon$$

$$B \rightarrow b \ E \ | \ a \ B \ B$$

(a) Describe the language generated by the grammar.

(b) Give a parse treee for the string *abaa*$.

(c) Is it an LL(1) grammar ? Build the parsing table and identify the conflicts.

8. Design an SLR(1) parser for the following grammar:

$$
\begin{aligned}
S' &\rightarrow & S \; \$ \\
S &\rightarrow & V ; S \mid \varepsilon \\
V &\rightarrow & \texttt{int id}
\end{aligned}
$$

9. Design an SLR(1) and LL(1) parsers for the following grammar:

$$
\begin{aligned}
P &\rightarrow & E \; \$ \\
E &\rightarrow & \texttt{atom} \mid {}' E \mid ( E \, E_S ) \\
E_S &\rightarrow & E \, E_S \mid \varepsilon
\end{aligned}
$$

Give the leftmost derivation and the parse tree for the string `(cdr '(a b c))$`.