

Database Languages and their Compilers

Prof. Dr. Torsten Grust

Database Systems Research Group
U Tübingen

Winter 2010

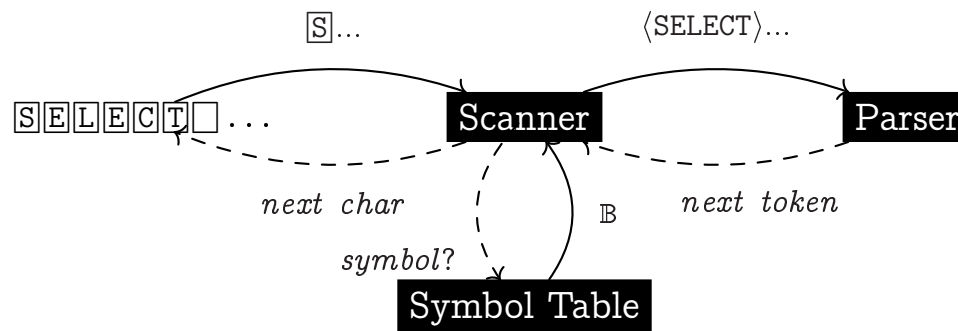


Copyright © 1995 United Feature Syndicate, Inc.
Redistribution in whole or in part prohibited

1 Query Parsing

1.1 The Scanner (Lexer)

- ▶ On the most basic level, any query input is nothing but a stream of characters, including **whitespace** (e.g., , `␣`, `→`);
- ▶ a scanner (or lexer) drops whitespace, and maps sequences of characters to tokens (`→` returns, `-->` calls):



- ▶ The parser drives the whole scanning and parsing process.

► Example tokenization:

SELECT *	≡	S E L E C T * ↵	scan →	⟨SELECT⟩ ⟨STAR⟩
FROM R	≡	F R O M R ↵		⟨FROM⟩ ⟨ID, "R"⟩
WHERE A = 42	≡	W H E R E A = 4 2 ↵		⟨WHERE⟩ ⟨ID, "A"⟩ ⟨EQ⟩ ⟨NUM, 42⟩

► The SQL scanner will generate a number of different tokens which can be classified as follows (with $b :: \mathbb{B}, n :: \mathbb{N}, s :: \mathbb{S}$)³:

Reserved words	Symbols	Constants	Identifiers
⟨SELECT⟩, ⟨SUM⟩, ⟨AND⟩	⟨STAR⟩, ⟨DOT⟩, ⟨EQ⟩	⟨BOOL, b ⟩, ⟨NUM, n ⟩, ⟨STRING, s ⟩	⟨ID, s ⟩

► Reserved words (and symbols) make the parsing of a language a lot simpler:

- SELECT FROM, SELECT FROM WHERE, SELECT, FROM WHERE SELECT = FROM
- The scanner uses a **symbol table** to recognize reserved words in the input:

symbol? ("SELECT") \dashrightarrow true

symbol? ("Mail") \dashrightarrow false .

³ $b, n,$ and s are known as **semantic token values**.

A Hand-Crafted Scanner

```
1 WHILE true DO  
2    $c \leftarrow \text{nextchar}()$ ;  
3   IF  $c \in \{\square, \blacksquare, \boxtimes\}$  THEN  
4      $\lfloor$ ;  
5   IF  $c \in \{0 \dots 9\}$  THEN  
6      $n \leftarrow \text{val}(c)$ ;  
7      $c \leftarrow \text{nextchar}()$ ;  
8     WHILE  $c \in \{0 \dots 9\}$  DO  
9        $n \leftarrow 10 * n + \text{val}(c)$ ;  
10       $c \leftarrow \text{nextchar}()$ ;  
11       $\text{pushback}(c)$ ;  
12      RETURN  $\langle \text{NUM}, n \rangle$ ;  
13   IF  $c \in \{\square, \square, \equiv, \ltimes, \triangleright, \oplus, \ominus, \ast, \surd\}$  THEN  
14      $s \leftarrow c$ ;  
15     IF  $c \in \{\ltimes, \triangleright\}$  THEN  
16        $c \leftarrow \text{nextchar}()$ ;  
17       IF  $c \in \{\equiv, \triangleright\}$  THEN  
18          $s \leftarrow s@c$ ;  
19       ELSE  
20          $\lfloor \text{pushback}(c)$ ;  
21       RETURN  $\text{lookup}(s)$ ;  
22   IF  $c \in \{a \dots z, A \dots Z, \square\}$  THEN  
23      $s \leftarrow c$ ;  
24      $c \leftarrow \text{nextchar}$ ;  
25     WHILE  $c \in \{a \dots z, A \dots Z, \square, \square, 0 \dots 9\}$  DO  
26        $s \leftarrow s@c$ ;  
27        $c \leftarrow \text{nextchar}()$ ;  
28      $\text{pushback}(c)$ ;  
29     IF  $\text{symbol?}(s)$  THEN  
30       RETURN  $\text{lookup}(s)$ ;  
31     ELSE  
32       RETURN  $\langle \text{ID}, s \rangle$ ;
```



Some remarks on

nextchar(), *pushback(c)*, *lookup(s)*

► This scanner implementation makes use of recurring patterns of code:

■ Check for **presence** of a specific character c (or a set of characters):

```
IF  $c \in \{\square, \uparrow, \rightarrow\}$  THEN  
└ ...;
```

■ Check for a **sequence** of specific characters:

```
IF  $c \in \{\langle, \rangle\}$  THEN  
└  $c \leftarrow nextchar()$ ;  
  IF  $c \in \{=, >\}$  THEN  
  └ ...;
```

■ Check for **repetitions** of the same character (or set of characters):

```
WHILE  $c \in \{a \dots z, A \dots Z, \_ , 0 \dots 9\}$  DO  
└ ...;  
   $c \leftarrow nextchar()$ ;
```

1.1.1 Regular Expressions



A Regular Expression Pattern (Reg.Ex.)

... matches a well-defined set of character strings. Any reg.ex. is built by using the building blocks listed below.

► Let Σ be an ordered set of characters (e.g., the ASCII code set). A **reg.ex. over Σ** takes one of the following forms:

- ① $[c]$ $c \in \Sigma$, matches the character $[c]$ only;
- ② $[c-c']$ $c, c' \in \Sigma$, matches any character in $\{c \dots c'\}$;
- ③ $[\^c-c']$ $c, c' \in \Sigma$, matches any character *not* in $\{c \dots c'\}$;
- ④ $re\ re'$ re, re' reg.ex., matches any **sequence** of strings matched by re and re' ;
- ⑤ re^* re reg.ex., matches 0 or more **repetitions** of strings matched by re ;
- ⑥ $re \mid re'$ re, re' reg.ex., matches any string matched by re **or** re' ;
- ⑦ ε matches the empty string "" only;
- ⑧ (re) re reg.ex., matches exactly the strings matched by re .

- Some handy abbreviations are in common use, e.g.:

$$\begin{aligned}
 re^+ &\equiv re\,re^* \\
 re? &\equiv re \mid \varepsilon \\
 . &\equiv [c-c'] \quad \text{with } \{c \dots c'\} = \Sigma \\
 [c-c'd-d'] &\equiv [c-c'] \mid [d-d']
 \end{aligned}$$

Example 1.1

- A crude reg.ex. to match e-mail addresses like "foo@bar.org" (with $alpha \equiv [a-zA-Z_]$, $alphanumeric \equiv alpha \mid [0-9]$):

$$alphanumeric^+ @ (alphanumeric^+ [.])^+ alpha^+$$

- Matching a floating point literal like "123.45E-10" (with $digit \equiv [0-9]$):

$$([+ | -])? ((digit^* [.] digit^+) \mid (digit^+ ([.] digit^*)?)) (E ([+ | -])? digit^+)?$$

Matching a C-style comment

Our SQL processor will accept (and ignore) C-style `/*...*/` comments. How could a suitable reg.ex. look like?

1.1.2 Regular Expressions for SQL

- ▶ The set of reg.ex. needed to identify SQL tokens is actually quite simple:

Reserved words (and symbols)

`[S][E][L][E][C][T]` `[G][R][O][U][P]` *ws*⁺ `[B][Y]` `[<][=]`

Constants

Numbers

`[0-9]`⁺

Booleans

`([T][R][U][E])` | `([F][A][L][S][E])`

Strings

`'` `[^']`* `'`

Identifiers

`[a-zA-Z_]` `[a-zA-Z_0-9]`*

Negative number literals?

Why don't we match negative number constants like `-42` here?

Rob Pike's Regular Expression Matcher

```
1  /* match: search for regexp anywhere in the text */
2  int match (char *regexp, char *text)
3  {
4      if (regexp[0] == '^')
5          return matchhere(regexp + 1, text);
6      do { /* must look even if string is empty */
7          if (matchhere(regexp, text))
8              return 1;
9      } while (*text++ != '\0');
10     return 0;
11 }
12
13 /* matchhere: search for regexp at the beginning of text */
14 int matchhere(char *regexp, char *text)
15 {
16     if (regexp[0] == '\0')
17         return 1;
18     if (regexp[1] == '*')
19         return matchstar(regexp[0], regexp + 2, text);
20     if (regexp[0] == '$' && regexp[1] == '\0')
21         return *text == '\0';
22     if (*text != '\0' && (regexp[0] == '.' || regexp[0] == *text))
23         return matchhere(regexp + 1, text + 1);
24     return 0;
25 }
```

```
1
2  /* matchstar: leftmost longest search for c*regexp */
3  int matchstar(int c, char *regexp, char *text)
4  {
5      char *t;
6
7      for (t = text; *t != '\0' && (*t == c || c == '.' ); t++)
8          ;
9      do { /* * matches zero or more */
10         if (matchhere(regexp, t))
11             return 1;
12     } while (t-- > text);
13     return 0;
14 }
```

- Obviously, reg.ex. provide a flexible tool to concisely specify scanners. A possible *nexttoken()* routine:

Build a table with entries $re_i \mapsto \langle t_i, n_i \rangle, \quad i = 1 \dots ;$

WHILE *nextchar()* \neq EOF DO

pushback();

 Call *nextchar()* repeatedly and match the input against all re_i ;

 IF input s matches re_j THEN

 RETURN $\langle t_j, n_j \rangle$;

 ELSE

 Call some default action;



- **N.B.** *nexttoken()* repeatedly scans the input for the **longest** match it can find (consider the SQL reg.exs. `[OR]` and `[OR][D][E][R] ws+ [B][Y]`).
- In general, the value of n_j will depend on the actual text matched against re_j (consider entry $[0-9]^+ \mapsto \langle \text{NUM}, n \rangle$):

IF input s matches re_j THEN

 RETURN $\langle t_j, n_j(s) \rangle$

ELSE

 ...;

1.1.3 `lex` – A Scanner Generator

- ▶ The scanner generator `lex` implements this approach:
 - **input:** a table of entries $re_i \mapsto \langle t_i, n_i \rangle$;
 - **output:** a C routine `yylex()` that implements a scanner for the specified language.
- ▶ A `lex` input file consists of a series of definitions and the actual table entries (separated by `%%`).
- ▶ In an action, the actual matched input s is available in the variable `char *yytext`, it's length in `int yyleng`.
- ▶ A semantic token value has to be assigned to the global variable `int yylval`. The entry $re_j \mapsto \langle t_j, n_j(s) \rangle$, is thus implemented by:

```
re_j          { yylval = n_j(yytext); return t_j; }
```

```

1  %{
2  #include <tokens.h>
3  %}
4
5  digit                [0-9]
6  ws                   [ \t\n]
7  identifier           [a-zA-Z_] [a-zA-Z_0-9]*
8  other                .
9
10 %%
11
12 "'\"([^\n]|\\.|\n)*'"  { *(yytext+yytext-1) = '\0';
13                        yylval = (int)yytext+1; return STRING; }
14 {digit}+             yylval = atoi(yytext);    return NUM;
15 "TRUE" | "FALSE"    yylval = *yytext == 'T'; return BOOL;
16 "SELECT"            return SELECT;
17 "GROUP"{ws}+"BY"    return GROUPBY;
18 {identifier}        yylval = (int)yytext;    return ID;
19 "/*"([^\n]|[\n])*"*/" | {ws}+ ;
20 {other}              return *yytext;

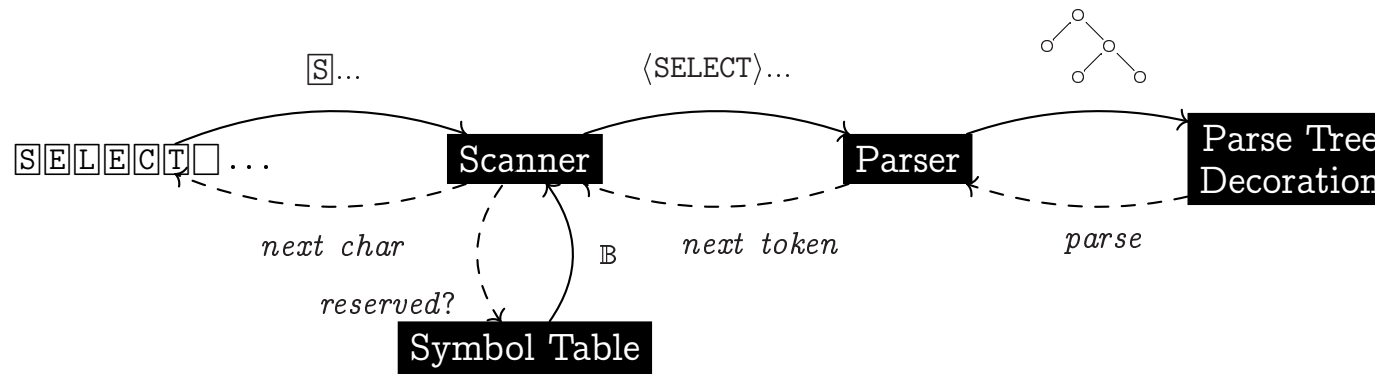
```

lex Input File

A few details about the lex input snippet printed on the previous slide and a live demo of the same.

1.2 Syntax Analysis (Parsing)

- ▶ We will use **context-free grammars**, a set of rules precisely describing the valid queries, to implement a **parser** for SQL.
- ▶ The parser will
 - implement the **syntax check**;
 - generate a **parse tree** representing the input query structure:

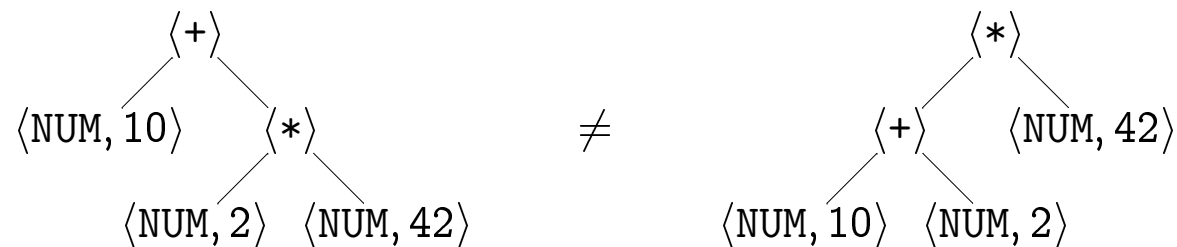


- ▶ Such parsers are easy to specify, extensible, and may be derived from the associated grammar in a fully automated fashion.

- ▶ Although the scanner emits a **linear token sequence**, the parser has to reveal the implicit **tree-shaped structure**:⁴

■ $10 + 2 * 42 \xrightarrow{\text{scan}} \langle \text{NUM}, 10 \rangle \langle + \rangle \langle \text{NUM}, 2 \rangle \langle * \rangle \langle \text{NUM}, 42 \rangle$

The well-known language of mathematics imposes the following structure on the above token sequence:



Dangling Tokens

Will there be a problem in determining the tree-shape for the following query?

```
SELECT  *
FROM    R, SELECT  *
                FROM  S
WHERE   A = 42
```

⁴The input language's **semantics** determines the tree shape for a token sequence.

1.2.1 Context-Free Grammars

- ▶ A list of production rules facilitates query syntax checking as well as parse tree construction.

- Production rule i :

$$S_{i_0} \rightarrow S_{i_1} S_{i_2} \dots S_{i_n}$$

- Symbols S_{i_j} ($j \geq 0$) are either **terminal** (i.e. tokens) or **nonterminal**;
 - the left-hand symbol S_{i_0} is always nonterminal;
 - S_{i_0} may appear on the right-hand side of its production rule (**recursion**).
- ▶ A list of production rules forms a **context-free grammar** iff
 - there is at least one rule of the form $S_i \rightarrow \dots$ for all nonterminal symbols S_i occurring on the right-hand side of a rule;
 - one nonterminal symbol is marked as the **start symbol** (conventionally, S_{0_0} is the start symbol).

- The following context-free grammar describes the syntax of (simplified) SQL SELECT-FROM blocks:

$$sf \rightarrow \langle \text{SELECT} \rangle \textit{proj_list} \langle \text{FROM} \rangle \textit{from_list}$$
$$\textit{proj_list} \rightarrow \textit{expr}$$
$$\textit{proj_list} \rightarrow \textit{proj_list} \langle , \rangle \textit{proj_list}$$
$$\textit{from_list} \rightarrow \langle \text{ID}, _ \rangle$$
$$\textit{from_list} \rightarrow \textit{from_list} \langle , \rangle \textit{from_list}$$
$$\textit{expr} \rightarrow \langle \text{ID}, _ \rangle$$
$$\textit{expr} \rightarrow \langle \text{NUM}, _ \rangle$$
$$\textit{expr} \rightarrow \textit{expr} \langle + \rangle \textit{expr}$$
$$\textit{expr} \rightarrow \textit{expr} \langle * \rangle \textit{expr}$$
$$\textit{expr} \rightarrow \langle (\rangle \textit{expr} \langle) \rangle$$

- sf , $proj_list$, $from_list$, $expr$ are the nonterminal symbols of this grammar; sf is the start symbol.

N.B. The grammar is indifferent to the semantic token values ($\langle \text{NUM}, _ \rangle$ vs. $\langle \text{NUM}, n \rangle$).

1.2.2 Derivations and Parse Trees

- ▶ Using this context-free grammar, we can perform a **syntax check** for a given query (token sequence) q as follows:

```
 $d \leftarrow S_{0_0};$   
WHILE  $\exists$  nonterminal symbol  $S_{i_0}$  in  $d$  DO  
⊙  $\left[ \begin{array}{l} \text{Choose any rule } S_{i_0} \rightarrow S_{i_1} S_{i_2} \dots S_{i_n}; \\ \text{Replace } S_{i_0} \text{ by } S_{i_1} S_{i_2} \dots S_{i_n} \text{ in } d; \end{array} \right.$   
  
RETURN  $d = q;$ 
```

- If this procedure returns *true*, q is a syntactically valid SELECT-FROM block;
- if we list d 's values before each replacement step we obtain a **derivation** for q .

⊙ Choice of Rule

For one grammar and one query q there may be many (or none) possible derivations.

- One possible derivation for the query below is shown here (nonterminal symbol chosen for next replacement):

```
SELECT  A, A + 1
FROM    R, S
```

sf

↳ ⟨SELECT⟩ proj_list ⟨FROM⟩ from_list

↳ ⟨SELECT⟩ proj_list ⟨FROM⟩ from_list ⟨,⟩ from_list

↳ ⟨SELECT⟩ proj_list ⟨FROM⟩ from_list ⟨,⟩ ⟨ID, "S"⟩

↳ ⟨SELECT⟩ proj_list ⟨FROM⟩ ⟨ID, "R"⟩ ⟨,⟩ ⟨ID, "S"⟩

↳ ⟨SELECT⟩ proj_list ⟨,⟩ proj_list ⟨FROM⟩ ⟨ID, "R"⟩ ⟨,⟩ ⟨ID, "S"⟩

↳ ⟨SELECT⟩ proj_list ⟨,⟩ expr ⟨FROM⟩ ⟨ID, "R"⟩ ⟨,⟩ ⟨ID, "S"⟩

↳ ⟨SELECT⟩ proj_list ⟨,⟩ expr ⟨+⟩ expr ⟨FROM⟩ ⟨ID, "R"⟩ ⟨,⟩ ⟨ID, "S"⟩

↳ ⟨SELECT⟩ proj_list ⟨,⟩ expr ⟨+⟩ ⟨NUM, 1⟩ ⟨FROM⟩ ⟨ID, "R"⟩ ⟨,⟩ ⟨ID, "S"⟩

↳ ⟨SELECT⟩ proj_list ⟨,⟩ ⟨ID, "A"⟩ ⟨+⟩ ⟨NUM, 1⟩ ⟨FROM⟩ ⟨ID, "R"⟩ ⟨,⟩ ⟨ID, "S"⟩

↳ ⟨SELECT⟩ expr ⟨,⟩ ⟨ID, "A"⟩ ⟨+⟩ ⟨NUM, 1⟩ ⟨FROM⟩ ⟨ID, "R"⟩ ⟨,⟩ ⟨ID, "S"⟩

↳ ⟨SELECT⟩ ⟨ID, "A"⟩ ⟨,⟩ ⟨ID, "A"⟩ ⟨+⟩ ⟨NUM, 1⟩ ⟨FROM⟩ ⟨ID, "R"⟩ ⟨,⟩ ⟨ID, "S"⟩

- We can modify our syntax check procedure for a query q to automatically produce a **parse tree** as a side-effect during derivation:

$d \leftarrow S_{0_0};$

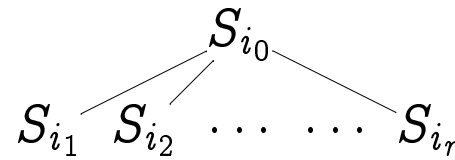
$t \leftarrow S_{0_0};$

WHILE \exists *nonterminal symbol* S_{i_0} in d DO

 Choose any rule $S_{i_0} \rightarrow S_{i_1} S_{i_2} \dots S_{i_n};$

 Replace S_{i_0} by $S_{i_1} S_{i_2} \dots S_{i_n}$ in $d;$

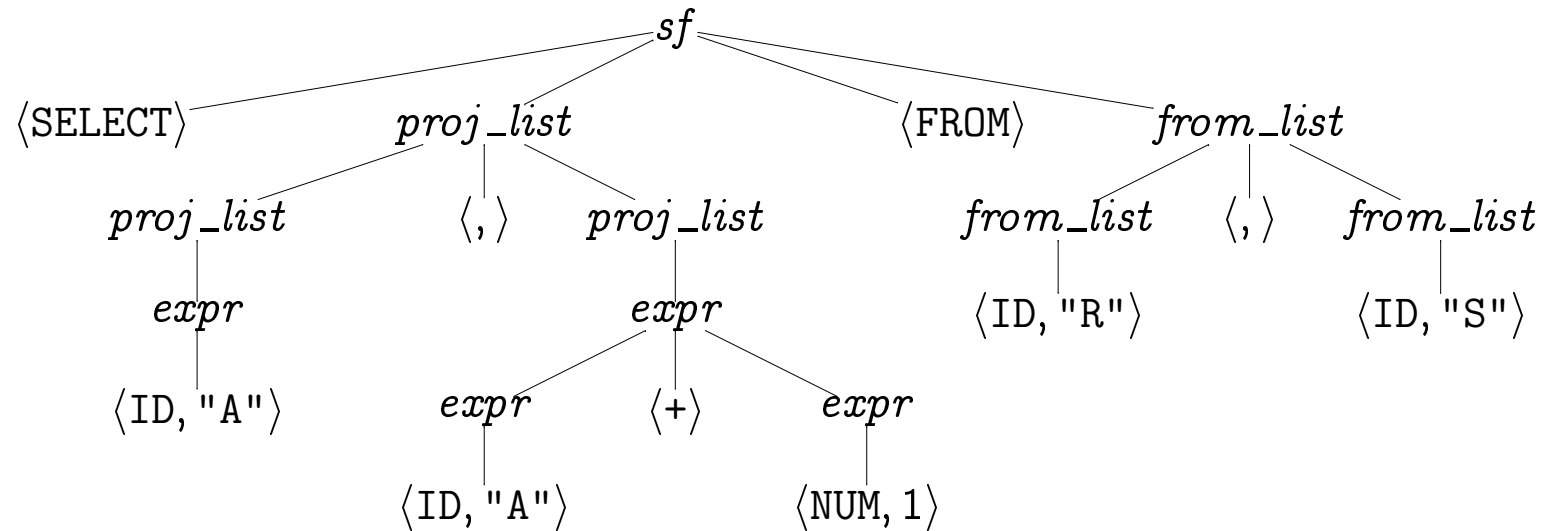
 Replace S_{i_0} by  in $t;$



RETURN $(d = q, t);$

- These parse trees provide exactly the tree-shaped query structure we were aiming for in this section.

- The derivation for q shown earlier corresponds to the following parse tree:

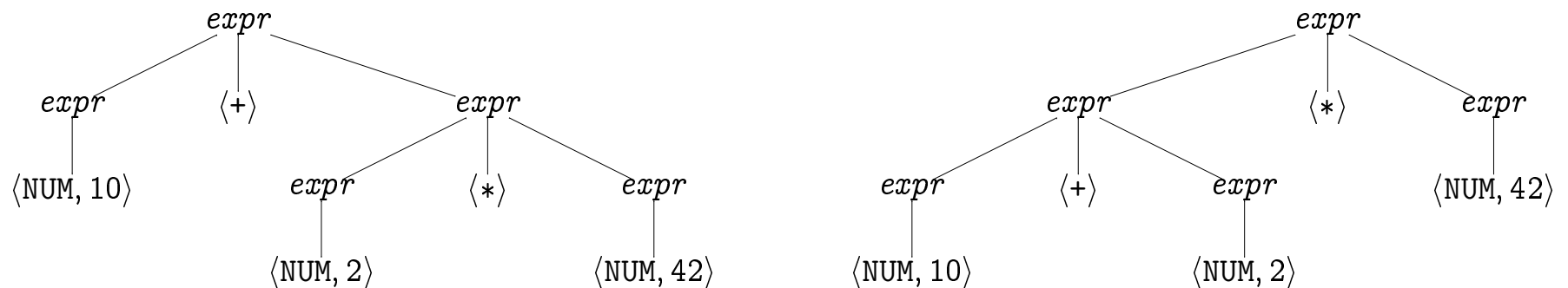


- For a successful syntax check, a **pre-order traversal** of the parse tree's leaves⁵ equals q .

⁵Also known as the tree's **front**.

1.2.3 Ambiguous Grammars

- ▶ There are grammars for which we can derive **different** parse trees for one input;
- ▶ obviously, this is a **bad thing**: parsers use parse trees to derive meaning (see the beginning of our discussion of parsers).
- Given the simplified SELECT-FROM block grammar (using *expr* as the start symbol), we can derive two parse trees for the input $10 + 2 * 42$:



- SQL's notion of the usual **operator precedence** is not reflected in the grammar.

✎ Comma-Separated Lists

Have a closer look at the (sub-)parse tree (start symbol *from_list*) for the relation list `R, S, T` in the SQL query

```
SELECT  A
        FROM  R, S, T .
```

Any ambiguity here? If there is, do such list-like constructs deserve special attention?

► We can encode operator precedence in a context-free grammar.

① Identify the precedence levels and the operators in each level (operators on higher levels bind tighter):

level	operator(s)
⋮	⋮
6	* / ...
5	+ - ...
⋮	⋮

② Introduce new separate production rules for each precedence level:

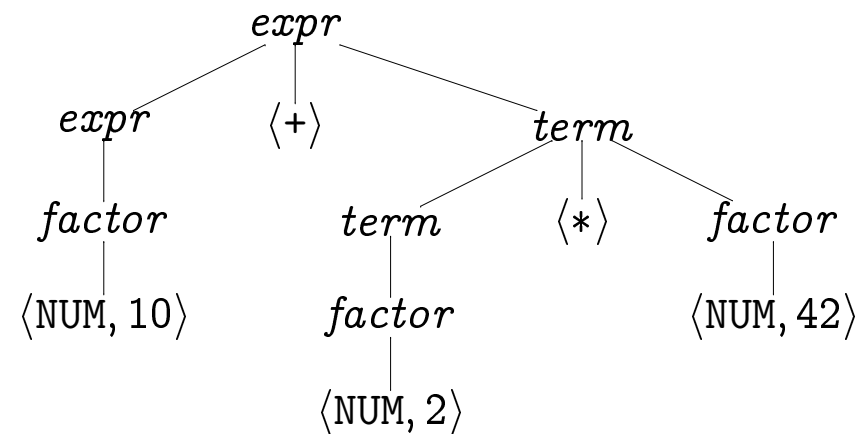
$$\begin{array}{lll}
 \textit{expr} & \rightarrow & \textit{expr} \langle + \rangle \textit{term} \\
 \textit{expr} & \rightarrow & \textit{expr} \langle - \rangle \textit{term} \\
 \textit{expr} & \rightarrow & \textit{term}
 \end{array}
 \quad \left. \vphantom{\begin{array}{l} \textit{expr} \\ \textit{expr} \\ \textit{expr} \end{array}} \right\} \text{level 5}$$

$$\begin{array}{lll}
 \textit{term} & \rightarrow & \textit{term} \langle * \rangle \textit{factor} \\
 \textit{term} & \rightarrow & \textit{term} \langle / \rangle \textit{factor} \\
 \textit{term} & \rightarrow & \textit{factor}
 \end{array}
 \quad \left. \vphantom{\begin{array}{l} \textit{term} \\ \textit{term} \\ \textit{term} \end{array}} \right\} \text{level 6}$$

$$\begin{array}{ll}
 \textit{factor} & \rightarrow \langle \text{ID}, - \rangle \\
 \textit{factor} & \rightarrow \langle \text{NUM}, - \rangle \\
 \textit{factor} & \rightarrow \langle (\rangle \textit{expr} \langle) \rangle
 \end{array}$$

- ▶ Leftmost derivation and parse tree for the expression $10 + 2 * 42$ (start symbol *expr*):

expr
 ↙
expr $\langle + \rangle$ *term*
 ↙
factor $\langle + \rangle$ *term*
 ↙
 $\langle \text{NUM}, 10 \rangle$ $\langle + \rangle$ *term*
 ↙
 $\langle \text{NUM}, 10 \rangle$ $\langle + \rangle$ *term* $\langle * \rangle$ *factor*
 ↙
 $\langle \text{NUM}, 10 \rangle$ $\langle + \rangle$ *factor* $\langle * \rangle$ *factor*
 ↙
 $\langle \text{NUM}, 10 \rangle$ $\langle + \rangle$ $\langle \text{NUM}, 2 \rangle$ $\langle * \rangle$ *factor*
 ↙
 $\langle \text{NUM}, 10 \rangle$ $\langle + \rangle$ $\langle \text{NUM}, 2 \rangle$ $\langle * \rangle$ $\langle \text{NUM}, 42 \rangle$



Precedence captured?

How does this grammar transformation capture operator precedence?
 How about associativity of operators (remember $10 - 20 - 30$)?

1.2.4 Predictive Parsing and Recursive Descent

- Once we have constructed a context-free grammar for a language, we can **derive a parser program** from it using a simple strategy.

In essence, the **production rules for nonterminal symbol S_{i_0}** are turned into a **single procedure $S_{i_0}()$** . $S_{i_0}()$ is recursive whenever the production rules for S are:

$$S_{i_0} \rightarrow S_{i_1} S_{i_2} \dots S_{i_n}$$

- ① S_{i_j} nonterminal: call procedure $S_{i_j}()$;
- ② S_{i_j} terminal $\langle t \rangle$: call $eat(\langle t \rangle)$.

$eat(\langle t \rangle)$:

```
IF  $lookahead = \langle t \rangle$  THEN
|  $lookahead \leftarrow nexttoken()$ ;
ELSE
|  $error()$ ;
```

- ③ Parsing starts by executing: $lookahead \leftarrow nexttoken()$; $S_{0_0}()$;

- Translate the following excerpt of the SQL grammar into a parser program:

```

      :
  quantifier → ⟨FORALL⟩ var ⟨IN⟩ coll ⟨:⟩ condition
  quantifier → ⟨EXISTS⟩ var ⟨IN⟩ coll ⟨:⟩ condition

      var → ⟨ID, _⟩
      :
```

Procedure *quantifier*():

```
SWITCH lookahead DO
  CASE ⟨FORALL⟩
  | eat(⟨FORALL⟩); var(); eat(⟨IN⟩); coll(); eat(⟨:⟩); condition();
  CASE ⟨EXISTS⟩
  | eat(⟨EXISTS⟩); var(); eat(⟨IN⟩); coll(); eat(⟨:⟩); condition();
  OTHERWISE
  | error();
```

Procedure *var*():

```
eat(⟨ID, _⟩);
```

✍ **Excerpt: lex-based predictive recursive descent SQL parser**

A recursive descent parser (using a lex-built scanner) for a subset of SQL clauses:

quantifier → $\langle \text{FORALL} \rangle$ *var* $\langle \text{IN} \rangle$ *coll* $\langle : \rangle$ *condition*

quantifier → $\langle \text{EXISTS} \rangle$ *var* $\langle \text{IN} \rangle$ *coll* $\langle : \rangle$ *condition*

coll → *table_ref*

condition → *var rel var*

rel → $\langle = \rangle$

rel → $\langle < \rangle$

rel → $\langle \leq \rangle$

rel → $\langle \geq \rangle$

rel → $\langle < \rangle$

rel → $\langle > \rangle$

var → $\langle \text{ID}, - \rangle$

table_ref → $\langle \text{ID}, - \rangle$

- Constructing a recursive descent parser for a given grammar seems simple enough. Let us add the VALUES SQL clause:

$$\begin{aligned} coll &\rightarrow table_ref \\ coll &\rightarrow coll_const \\ \\ coll_const &\rightarrow \langle \text{VALUES} \rangle coll_elems \\ \\ coll_elems &\rightarrow const \\ coll_elems &\rightarrow coll_elems \langle , \rangle const \\ \\ const &\rightarrow \langle \text{NUM} \rangle \end{aligned}$$

Procedure *coll_elems()*:

```
SWITCH lookahead DO
  CASE  $\langle \text{NUM} \rangle$  ①
  | const();
  CASE  $\langle \text{NUM} \rangle$  ②
  | coll_elems(); eat( $\langle , \rangle$ ); const();
  OTHERWISE
  | error();
```

- Procedure *coll_elems()* faces a dilemma: given only **one lookahead token** (marked below) how can it determine the proper case?

Input	Case
... VALUES <u>4</u> 2	①
... VALUES <u>2</u> , 10, 42	②



A single lookahead token must suffice

Predictive parsing will only work if knowledge of the **next terminal symbol** (the lookahead) is sufficient to determine which production rule to use next.

- ▶ Given (the right-hand side of) a production rule $S_{i_0} \rightarrow S_{i_1} \dots S_{i_n}$, $FIRST(S_{i_1} \dots S_{i_n})$ is the set of all tokens that can occur as the first token in any token sequence derivable from $S_{i_1} \dots S_{i_n}$.

■ Given this grammar

$$term \rightarrow term \langle * \rangle factor$$

$$term \rightarrow term \langle / \rangle factor$$

$$term \rightarrow factor$$

$$factor \rightarrow \langle ID, - \rangle$$

$$factor \rightarrow \langle NUM, - \rangle$$

$$factor \rightarrow \langle (\rangle term \langle) \rangle$$

we have

$$FIRST(term \langle * \rangle factor) = \{ \langle ID, - \rangle, \langle NUM, - \rangle, \langle (\rangle \} .$$

⚡ Computing *FIRST*

Determining the *FIRST* set looks very simple. While computing $FIRST(S_{i_1} S_{i_2} S_{i_3})$ it seems as if S_{i_2} and S_{i_3} can be ignored and $FIRST(S_{i_1})$ is all that matters. This may not be true. (Why?)

► Clearly, if we are given production rules

$$\begin{array}{l} S_{i_0} \rightarrow \sigma_1 \\ S_{i_0} \rightarrow \sigma_2 \end{array}$$

with $\langle t \rangle \in FIRST(\sigma_1) \cap FIRST(\sigma_2)$, a predictive parser will not know what to do if the lookahead is $\langle t \rangle$.

■ Production rules for *coll_elems*:

$$\begin{array}{l} coll_elems \rightarrow const \\ coll_elems \rightarrow coll_elems \langle , \rangle const \end{array}$$

$$const \rightarrow \langle NUM \rangle$$

$$FIRST(const) \cap FIRST(coll_elems \langle , \rangle const) = \{ \langle NUM \rangle \}.$$

► Observe that the core of our problem of common elements in *FIRST* sets is the **left-recursion** in the production rule for *coll_elems*.

■ A production rule is said to be **left-recursive** if it takes the following general form:

$$\begin{aligned} S &\rightarrow S \sigma_1 \\ S &\rightarrow \sigma_2 \end{aligned}$$

where σ_2 does not start with S .

Eliminating left-recursion

Fortunately, we can get rid of left-recursion in rules all together by a simple rule transformation:

$$\begin{aligned} S &\rightarrow S \sigma_1 \\ S &\rightarrow \sigma_2 \end{aligned} \quad \dashrightarrow$$

■ The modified rules pose no problem for a predictive parser.

- ▶ A similar problem arises when **two or more production rules** for the same nonterminal S **start with the same symbols**:

$$\begin{array}{l}
 sfw \rightarrow \underbrace{\text{SELECT } proj_list \text{ FROM } from_list} \\
 sfw \rightarrow \underbrace{\text{SELECT } proj_list \text{ FROM } from_list}_{=} \text{ WHERE } condition
 \end{array}$$

- Only after the parser has seen the last token of the expansion of $from_list$ it can decide which rule to choose (*lookahead* = $\langle \text{WHERE} \rangle ?$).
This does not go together with predictive one symbol lookahead parsing.

Left-factorization

Once more we can transform the production rules and make predictive parsing possible again:

$$\begin{array}{l}
 S \rightarrow \sigma_1 \sigma_2 \\
 S \rightarrow \sigma_1 \sigma_3
 \end{array}
 \quad \dashrightarrow$$

1.2.5 yacc – A Parser Generator

- ▶ Once we have come up with a grammar for a language, it is simple enough to **automatically derive a parser** from it. This is exactly what yacc (Yet Another Compiler Compiler) provides.
- ▶ yacc-generated parsers are **not** predictive but use a parsing technique called LALR(1)⁶.
- ▶ yacc-generated parsers
 - can deal with left-recursion in rules;
 - do not need their rules left-factorized;
 - are able to use operator precedence and associativity tables to resolve conflicts during rule selection.

⁶1-symbol Lookahead, Left-to-Right Parse, Rightmost Derivation, see [1, 2].

- ▶ A yacc input file comes in the general form

```

parser declarations           (tokens, operator precedence/associativity)
%%
grammar rules
%%
C program code

```

- ▶ Parser declarations:

```

1 %token OR AND EQ ...
2 %token BOOL ID SELECT FROM WHERE ...
3
4 %left      OR
5 %left      AND
6 %nonassoc  EQ NEQ LEQ GEQ LT GT
7 %left      '+'
8 %left      '*'
9 %nonassoc  NOT

```

- Implements this operator precedence table:

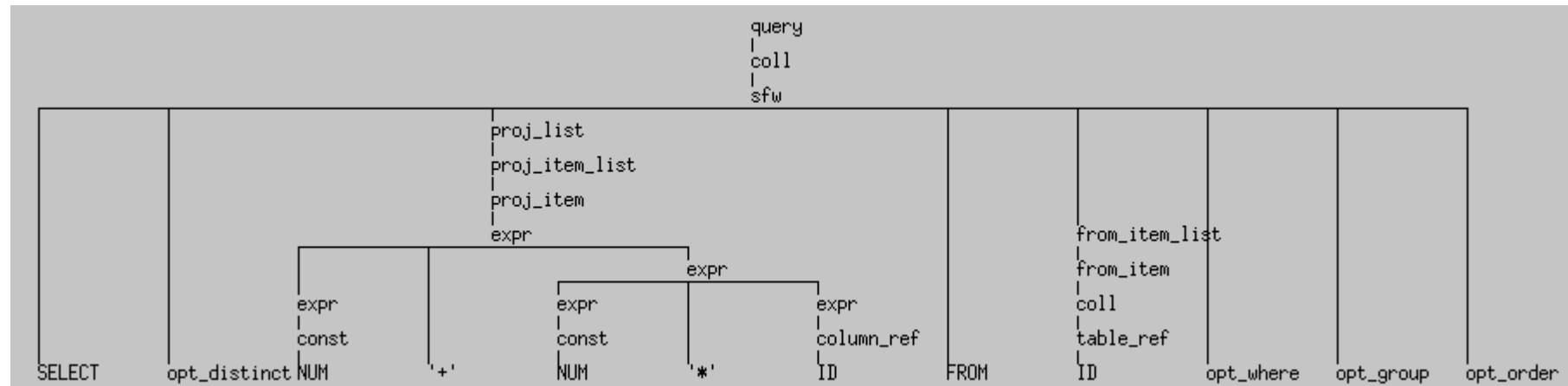
level	operator(s)
1	OR
2	AND
3	= <> <= >= < >
4	+
5	*
6	NOT

✍ yacc Input File

Take a closer look at the yacc input file for our SQL dialect and double check the corresponding parse trees.

► Parse tree for SQL query

```
SELECT  2 + 3 * A
FROM    R
```



1.3 Recommended Reading

- [1] A. V. AHO, R. SETHI, AND J. D. ULLMAN, *Compilers: Principles, Techniques, and Tools*, Addison Wesley, 1986. ISBN 0-201-10088-6. The “Dragon” Book.
- [2] A. W. APPEL, *Modern Compiler Implementation in C: Basic Techniques*, Cambridge University Press, 1997. ISBN 0-521-58653-4.

Contents

0	Introduction to Query Compilation	1
0.1	Welcome!	1
0.2	Administrativa	3
0.3	Some Remarks on these Slides	4
0.4	Relational Databases and SQL	6
0.5	A Guided Tour through an SQL Query Processor	14
0.5.1	Character Streams and Tokens	15
0.5.2	Identify Valid SQL Syntax	16
0.5.3	Resolve the Meaning of Variables and Identifiers	17
0.5.4	Check Types and Schemas	18

0.5.5	How Could Query Evaluation Look Like?	19
0.5.6	One Query, Many Programs	20
0.5.7	Query Operators	21
0.5.8	Query Operator Trees and Rewriting	22
0.5.9	The Cheaper, the Better: Query Cost Models	23
0.5.10	Executing Operator Trees on a Machine	24
0.6	Recommended Reading	27
1	Query Parsing	28
1.1	The Scanner (Lexer)	28
1.1.1	Regular Expressions	33
1.1.2	Regular Expressions for SQL	35
1.1.3	<u>lex – A Scanner Generator</u>	<u>39</u>

1.2	Syntax Analysis (Parsing)	42
1.2.1	Context-Free Grammars	44
1.2.2	Derivations and Parse Trees	46
1.2.3	Ambiguous Grammars	50
1.2.4	Predictive Parsing and Recursive Descent	54
1.2.5	yacc – A Parser Generator	63
1.3	Recommended Reading	67