



# 编译实验实验报告

姓名: 施鹏飞

学号: 18040100077

2021 年 6 月 30 日

# 目录

<b>1 实验内容</b>	<b>1</b>
<b>2 实验环境</b>	<b>1</b>
<b>3 内容分析</b>	<b>1</b>
<b>4 具体实现</b>	<b>2</b>
4.1 Lex 词法分析 . . . . .	2
4.1.1 解析 SQL 关键字 . . . . .	2
4.1.2 解析自定义的字串 . . . . .	2
4.2 解析标点符号 . . . . .	3
4.3 数据库系统的结构设计 . . . . .	4
4.3.1 物理结构 . . . . .	4
4.3.2 逻辑结构 . . . . .	5
4.3.2.1 基本链表 . . . . .	5
4.3.2.2 数据库系统结构 . . . . .	6
4.3.2.3 数据库结构 . . . . .	6
4.3.2.4 表单结构 . . . . .	6
4.3.2.5 属性域结构 . . . . .	7
4.3.2.6 条件结构 . . . . .	8
4.4 Yacc 的语法分析并与接口的交互 . . . . .	8
4.4.1 总体语句框架 . . . . .	9
4.4.2 创建数据库 . . . . .	10
4.4.3 使用数据库 . . . . .	10
4.4.4 显示所有数据库 . . . . .	11
4.4.5 删除数据库 . . . . .	11

4.4.6	创建表单 . . . . .	11
4.4.7	删除表单 . . . . .	13
4.4.8	显示所有表单 . . . . .	13
4.4.9	展示表单所有属性 . . . . .	13
4.4.10	插入元组 . . . . .	13
4.4.11	查询元组 . . . . .	15
4.4.12	退出系统 . . . . .	17
4.4.13	删除元组与更新元组 . . . . .	18
<b>5</b>	<b>结果展示</b>	<b>18</b>
<b>6</b>	<b>局限与展望</b>	<b>20</b>
<b>7</b>	<b>个人体会</b>	<b>21</b>
<b>8</b>	<b>程序清单</b>	<b>21</b>

## 1 实验内容

本次编译原理实验中，我们借助 Lex/Yacc 等分析工具，实现了一个较为简单的数据库系统。该数据库系统为磁盘数据库，将数据存储在本机，程序结束后数据不会消失。用户可以在终端输入 SQL 语句操作数据库，执行相应的 SQL 功能。

## 2 实验环境

Ubuntu 20.04.2 LTS, flex 2.6.4, bison (GNU Bison) 3.5.1 以及 C++ 开发环境。

## 3 内容分析

首先，我们考虑希望实现的所有数据库功能及对应的 SQL 语句格式，罗列如下：

操作描述	SQL 语句
创建数据库	create database {数据库名}
删除数据库	drop database {数据库名}
使用数据库	use {数据库名}
显示所有数据库	show databases
创建表单	create table {列名 列类型, ...}
删除表单	delete table {表单名}
展示所有表单	show tables
展示某表单的属性	desc {表单名}
插入	insert into {表单名} ... [where ...]
查找	select ... from {表单名, ...}
删除	delete from {表单名} [where ...]
更新	update {表单名} set ... [where ...]
删除表单	drop table {表单名}
退出	exit

综合考虑后，我使用 Lex 生成文法记号，Yacc 对语法进行分析，当识别出不同的命令时，再分不同的情况进行功能实现：

- 对于仅涉及到一、两个字符串的命令（如创建数据库、删除表单、退出等），Yacc 直接将该字符串传给功能实现接口进行处理。

- 对于涉及到以列表形式给出指令（如创建数据表、插入元组、查找元组等）的情况时，Yacc 在对语法树归结时，直接将 SQL 命令中的列表信息以 `{内容所在的指针 + 内容有效字节}` 的格式一条条

插入预先给出的全局链表中；当整条语句解析完成后，Yacc 通知功能实现模块，功能实现模块访问该全局链表，按照预先商定好的顺序对其进行解析，进而完成 SQL 功能。

## 4 具体实现

考虑到在本套系统中，Lex 与其余部分的联系稍微弱一些，因而在本章节我们主要分三部分进行阐述。第一部分介绍 Lex 处理单词的方式，第二部分介绍数据库物理结构与逻辑结构，第三部分介绍 Yacc 对记号流的解析以及对应 C++ 数据库操作接口的设计。

### 4.1 Lex 词法分析

Lex 负责解析输入的字符串，并将其转换为记号流传达给 Yacc。

#### 4.1.1 解析 SQL 关键字

用户可能输入 SQL 关键字，此时 Lex 不区分大小写的解析出他们，并输出对应的宏定义值。以几个 SQL 关键字为例，给出代码详细说明 Lex 如何处理关键字：

##### 定义部分

```
1 create [cC][rR][eE][aA][tT][eE]
2 database [dD][aA][tT][aA][bB][aA][sS][eE]
3 databases [dD][aA][tT][aA][bB][aA][sS][eE][sS]
4 use [uU][sS][eE]
5 table [tT][aA][bB][lL][eE]
```

##### 识别规则部分

全部大写的返回值为 Yacc 中预先宏定义的常量值。

```
1 {create} {return CREATE;}
2 {database} {return DATABASE;}
3 {databases} {return DATABASES;}
4 {use} {return USE;}
5 {table} {return TABLE;}
```

#### 4.1.2 解析自定义的字串

用户可能输入如下三种自定义的类型：

- 整数 (由 0-9 构成的数, 允许前导 0)
- 字符串 (在本实验中我们认为字符串必须以**双引号**为首尾, 暂不与标准 SQL 对应)
- 变量名 (包含数字、52 个大小写字母与下划线, 且不以数字开头)。

### 定义部分

两个基本的类型 `charType` 和 `digitType` 分别表示可用的字符以及 10 个数码。 `integer` 是由若干个基本数码组成的一个合法的允许前导 0 的整数。 `id` 以一个字符起头, 随后跟任意数量 (可以为 0) 的数码或字符。 `string` 以双引号为首尾, 中间有任意数量 (可以为 0) 的数码或字符。

a.l

```
1 charType [a-zA-Z_]
2 digitType [0-9]
3
4 integer {digitType}+
5 id {charType}({charType}|{digitType})*
6 string \"({charType}|{digitType})*\"
```

### 识别规则部分

当识别到了这三种记号时, Lex 将其以字符串指针的形式存储在 `yylval` 中以供 Yacc 使用。**特别地**, 整数类型也以字符串存储, 解析的工作交给后人进行。

a.l

```
1 {id} {
2     yyval.str = strdup(yytext);
3     return ID;
4 }
5 {integer} {
6     yyval.str = strdup(yytext);
7     return INTEGER;
8 }
9
10 {string} {
11     yyval.str = strdup(yytext);
12     return STRING;
13 }
```

## 4.2 解析标点符号

下面给出所用到的全部标点信息。

a.l

```
1 ; {return FIN;}
2 \{ {return LB;}
3 \} {return RB;}
4 , {return COMMA;}
5 \" {return QM;}
6 \< {return BELOW;}
7 = {return EQU;}
8 \* {return STAR;}
```

## 4.3 数据库系统的结构设计

### 4.3.1 物理结构

在本地磁盘上, 数据库系统文件全部存放在 `DBMS/` 文件夹下。在 `DBMS` 下有一份元数据 `.dbinfo`, 用于说明该数据库系统有哪些数据库。其中记载了明文形式的所有数据库名字, 以供系统访问对应的次级文件夹。而在某一个数据库下, 以 `Database1` 为例, 其所有的表单以及元数据全部存储在文件夹 `Database1/` 下。在该文件夹下首先有一份元数据 `.tableinfo`, 用于说明该数据库下有哪些表单。其中记载了明文形式的所有表单名, 以供系统进一步访问表单文件。每一份表单文件的文件名即为表单名, 在其内部, 数据按照如下的字节顺序排列:

- 前 32 位为一个 32 位整数  $n$ , 表明该表单有多少个属性域。
- 接下来跟  $n$  块数据代表属性域的详细属性。每一块数据前 32 位为一个 32 位整数  $l_i$ , 表示该属性域的名称长度, 随后跟  $l_i$  字节的信息, 代表该属性域的名称, 再后跟 2 个 32 位整数, 分别为该属性域的占用字节长度和属性类型。
- 然后 32 位为一个 32 位的整数  $m$ , 表明该表单有多少个元组。
- 接下来跟  $m$  块数据代表表单的具体内容。每一块数据又分成  $n$  小块, 这每一小块就代表了该元组在对应属性域上的分量。在每一小块中, 首先给出 32 位的一个整数  $k_i$ , 代表该元组在该属性上拥有  $k_i$  字节的信息, 随后跟出  $k_i$  字节的数据, 表明其具体的取值。

尽管在上述具体阐述中, 说明字节块的逻辑意义看起来可能较为费力, 但是实际在对物理结构读写时只需要 `循环 + 字节读写` 即可, 比较方便。

数据库系统的层次物理结构, 如图1所示。

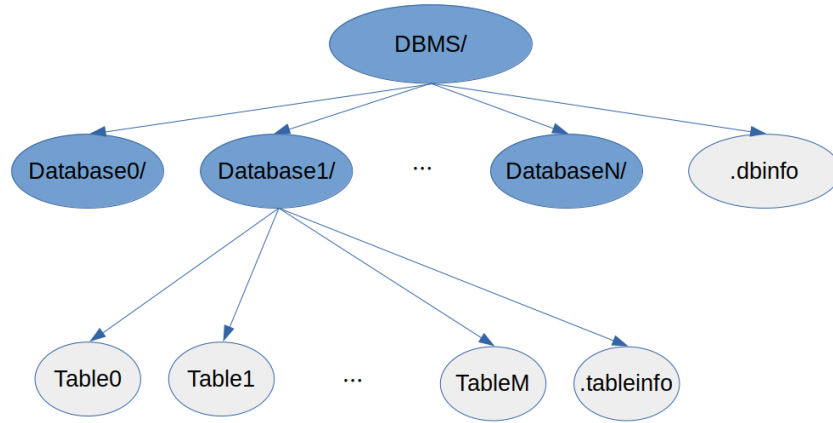


图 1: SQL 系统文件层次结构

### 4.3.2 逻辑结构

对应于4.3.1给出的物理结构，我们对该数据库系统的逻辑存储结构进行说明。

**4.3.2.1 基本链表** 本数据库系统使用自定义的链表结构，该链表的值域为一个空指针与该空指针指向内容的有效字节长度。有了这样两条信息，我们就可以使用该链表存储任何信息，十分方便。数据结构与函数接口如下所示。

下面给出链表的数据接口与函数接口信息。

list.h

```

1 struct ListElement{
2     int bytes;
3     void *val;
4     ListElement *pre, *next;
5 };
6
7 struct List{
8     ListElement *front, *back;
9     int length;
10
11     List();
12     void push_back(void *val, int L);
13     void erase(ListElement*);
14     void copy(List *from); //将from链表拥有所有元素信息拷贝到当前链表
15     void pop_front();
16     void pop_back();
17     void clear();
18     void save(const char *file, int append); //向文件file写入字节流，append非0则以附加形式缀到文件尾
19 };
  
```



**4.3.2.2 数据库系统结构** 数据库系统的结构较为简单，是一个链表元素为数据库的链表。

数据库系统提供如下的全局函数接口：创建数据库、使用数据库、删除数据库、显示所有的数据库、从本地读所有数据库、将所有数据库写入本地。

数据结构与函数接口如下所示。

sql.h

```
1 extern List dbs;
2 void read();
3 void write();
4 int create_database(const char* databaseName, int bytes);
5 int use_database(const char* databaseName, int bytes);
6 int drop_database(const char* databaseName, int bytes);
7 void show_databases();
```

**4.3.2.3 数据库结构** 数据库结构中有该数据库的名称，还有该数据库下所有表单形成的链表。

对某个数据库而言，其提供如下的成员函数接口：查找某表单、显示所有表单、展示某表单的属性域、删除某表单、删除所有表单、从本地读数据库、将数据库写入本地。

数据结构与函数接口如下所示。

sql.h

```
1 struct DATABASE{
2     char *name;
3     List table;
4
5     ListElement* find_table(const char* tableName, int bytes);
6     void show_tables();
7     int desc_table(const char *tableName, int bytes);
8     int drop_table(const char *tableName, int bytes);
9     void drop_all_table();
10
11     void read();
12     void write();
13 };
```

**4.3.2.4 表单结构** 表单结构中有该表单属于哪个数据库、该表单的名称、该表单拥有的属性数量、该表单的元组数量，以及所有属性构成的链表和所有元组构成的链表。

**注意**，在这里，一条元组本身就是一个链表，因而在该结构中，存储的元组链表其实是链表的链表。

对于某个表单而言，其提供如下的成员函数接口：检查元组信息合法性（在 where 子句中用到，检查 where 子句约束的属性是否都确实是存在的一个合法属性）、展示表单所有属性域、从本地读表单、将表单写入本地。

数据结构与函数接口如下所示。

sql.h

```
1 struct TABLE{
2     char *belong_to_db;
3     char *name;
4     int fieldNumber, turpleNumber;
5     List field, turpleEntry;
6
7     List* check_all_turple();
8     int check_turple(List *turple, ListElement *conElem);
9     void desc();
10
11     void read();
12     void write();
13 };
```

**4.3.2.5 属性域结构** 属性域结构记录了某属性的名称、占用的字节数以及其类型。属性域的类型是我们提前宏定义的整数类型，不同的整数代表不同的数据类型，目前暂时只支持INT和CHAR()。

在属性域中有一个特别的值域selected，该属性仅在进行 select 操作时使用。当调用 select 语句时，首先检查所有的属性，如果 select 语句要求展示该属性则将selected设置为真，不然则设置为假，最终在展示时我们就可以根据该值域的真假决定是否输出。

数据结构与函数接口如下所示。

sql.h

```
1 struct FIELD{
2     char *name;
3     int bytes;
4     int type;
5
6     // 仅仅在执行select操作时使用
7     int selected;
8 };
```

**4.3.2.6 条件结构** `条件`结构应用于 where 子句，用于记录某一个原子条件的信息。**注意**，该结构体只能表示二元比较符，并且要求左操作数必须为数据属性，右操作数必须是属性值（如 Age<20, Name < “ZhangSan”）。其内部有左操作数的指针及字节长度、右操作数的指针及字节长度、比较符的类型，比较符的类型预先进行了宏定义，是一个整数类型。

数据结构如下所示。

sql.h

```
1 struct CONDITION{
2     void *lval, *rval;
3     int lbytes, rbytes;
4     int cmp;
5 };
```

## 4.4 Yacc 的语法分析并与接口的交互

考虑到 Yacc 在对语法树归结时会相应的进行语义翻译，二者联系十分密切，所以这里我们将两部分放在一起进行说明，以期能有更好的效果。在逐一说明各个功能的 Yacc 语法分析与 C++ 接口使用方式之前，我们先给出 Yacc 用到的一些设置选项。

a.y

```
1 %union{
2     char *str;
3     int num;
4     struct {
5         char *p;
6         int bytes;
7     } M;
8 }
9
10 %token CREATE DATABASE DATABASES USE TABLE TABLES SHOW DESC INSERT INTO VALUES SELECT UPDATE DELETE
    DROP EXIT INT DOUBLE CHAR AND OR FROM WHERE SET
11 %token FIN LB RB COMMA BELOW EQU STAR QM
12 %token <str> ID STRING INTEGER
13 %type <str> create_database use_database desc_table drop_table drop_database value
14 %type <M> datatype
15 %type <num> create_table insert select
16
17 %left OR
18 %left AND
```

结合 Yacc 文件的部分相关代码，我们做如下说明：

- union 联合体有三种可能的情况：字符串指针、32 位整数、指针 + 字节数构成的结构体，在之后的程序中会根据不同的情况分配不同的数据类型。

- %token 给出了可能 Lex 可能用到的宏定义，自动生成相应的存放在 `y.tab.h` 中。
- 由于版本不同，在本实验环境中我们使用 %type 关键字来为文法中所有的终结符、非终结符分配类型。
- 使用 %left 指定左结合的优先级，解决多个逻辑表达式之间的组合问题。

#### 4.4.1 总体语句框架

`statements` 由一条或多条的 `statement` 组成。而每一个 `statement` 则分别可能产生各种不同的非终结符，对应不同的功能；根据其返回值打印一定的提示信息。

**a.y**

```
1 statements:
2     statement
3     | statements statement
4     ;
5
6 statement:
7     create_database {
8         if(create_database($1, strlen($1))) puts("create_database ok.");
9         else puts("create_database fail.");
10    }
11    | use_database {
12        if(use_database($1, strlen($1))) puts("use_database ok.");
13        else puts("use_database fail.");
14    }
15    | create_table {
16        if($1) puts("create_table ok.");
17        else puts("create_table fail.");
18    }
19    | show_databases {
20        show_databases();
21    }
22    | show_tables {
23        if(current_database == NULL) puts("No database selected!");
24        else current_database->show_tables();
25    }
26    | desc_table {
27        if(current_database == NULL) puts("No database selected!");
28        else current_database->desc_table($1, strlen($1));
29    }
30    | insert {
31        if($1) puts("insert_into ok.");
```

```

32     else puts("insert_into fail.");
33 }
34 |delete {printf("Delete something ok\n");}
35 |update {printf("Update ok\n");}
36 |select {
37     if($1 == 1) puts("Select ok");
38     else puts("Select fail");
39 }
40 |drop_table {
41     if(current_database == NULL) puts("No database selected!");
42     else if(current_database->drop_table($1, strlen($1)) == 0) puts("drop_table fail.");
43     else puts("drop_table ok.");
44 }
45 |drop_database {
46     if(drop_database($1, strlen($1))) puts("drop_database ok.");
47     else puts("drop_database fail.");
48 }
49 |exit {return 0;}
50 ;

```

#### 4.4.2 创建数据库

为了创建数据库,只需要将数据库的名称作为当前语句的内容向上传, `statement` 调用 `create_database` 进行创建。

`create_database` 函数接受一个字符串及其字节数,首先检查当前数据库系统下是否存在一个同名的数据库,如果不存在,那么先在文件系统中创建一个对应的数据库文件夹,并在文件夹内创建一个数据表元数据(初始为空),再更新内存中的信息,创建一个新的数据库元素,将对应的值初始化后插入到数据库系统链表尾。至此全部操作完成。

**a.y**

```

1 create_database:
2     CREATE DATABASE ID FIN {$$ = $3;}

```

#### 4.4.3 使用数据库

使用数据库也类似地只需要将数据库的名称作为当前语句的内容向上传, `statement` 调用 `use_database` 进行选择。

`use_database` 函数接受一个字符串及其字节数，检查该字符串对应的数据库是否存在，存在则操作成功。

**a.y**

```
1 use_database:
2     USE ID FIN {$$ = $2;}
```

#### 4.4.4 显示所有数据库

显示数据库不需要任何用户自行输入的参数，只要归结成功 `statement` 就调用 `show_databases` 展示所有的数据库。

`show_databases` 函数顺序遍历数据库系统链表中所有的数据表，并美观地将所有表名显示。

**a.y**

```
1 show_databases: SHOW DATABASES FIN
```

#### 4.4.5 删除数据库

删除数据库也类似地将数据库名称向上传，`statement` 调用 `statement` 调用 `drop_database` 进行删除。

`drop_database` 函数接受一个字符串及其字节数，首先检查当前数据库系统下是否存在一个同名的数据库，如果存在，那么先在文件系统中删除对应的数据库文件夹，再更新内存中的信息，将该数据库对应的链表元素删去。至此全部操作完成。

**a.y**

```
1 drop_database:
2     DROP DATABASE ID FIN {$$ = $3;}
```

#### 4.4.6 创建表单

创建表单不仅需要表单名，还需要该表单各种属性的名称、类型、长度等信息，这些额外信息的数量未知，因而我们使用一个全局的**链表**来存储所有的这些信息。

归结进行时首先肯定先读取到表单名时，将该字符串加入到链表，成为链表的第一个元素。

接下来当归结到某一个属性域时，首先插入一个链表元素表示该属性的名称；接下来插入一个链表元素该元素的第一个字节表示该属性的类型，是一个提前约定好的整数，如果该类型是字符串数组，则该链表元素后面的所有字节则给出该数组的长度，以字符串形式给出，交由后人解析成整数。

按照 Yacc 归结的顺序，可以保证该链表中所有的信息均是保持原有顺序的。

当创建表单整条语句全部读取完成后，接口函数 `temp_create_table` 就对该全局链表进行解析，创建相应的表单，这里就不予赘述了。

**注意！** 我们定义了一个全局的数据库指针 `current_database` 表明当前使用的数据库，在使用该创建表单功能前请先使用 `use` 指令选定数据库，不然该函数将运行失败并给出相应的报错。在本报告之后也有一些 SQL 功能需要提前使用 `use` 指令，之后便不再提及了，但是每一个之后对应的接口函数都对此做了错误处理，会产生相应的报错。

**a.y**

```
1 create_table:
2     CREATE TABLE ID LB {
3         temp.clear();
4         temp.push_back($3, strlen($3));
5     }
6     fields
7     RB FIN {
8         if(temp_create_table() == 0) $$ = 0;
9         else $$ = 1;
10    }
11 fields:
12     field
13     |fields COMMA field
14 field:
15     ID{
16         temp.push_back($1, strlen($1));
17     }
18     datatype {
19         temp.push_back($3.p, $3.bytes);
20     }
21 datatype:
22     INT {
23         $$ .bytes = 1;
24         $$ .p = (char*)malloc($$ .bytes);
25         $$ .p[0] = 1;
26     }
27     |CHAR LB INTEGER RB {
28         $$ .bytes = 1 + strlen($3);
```

```

29     $$p = (char*)malloc($$.bytes);
30     $$p[0] = 2;
31     for(int i = 1; i < $$.bytes; i++) $$p[i] = $3 [i - 1];
32 }

```

#### 4.4.7 删除表单

删除表单跟之前的删除数据库很类似，上传表单名后 statement 使用 `drop_table` 指令。

**a.y**

```

1 drop_table:
2     DROP TABLE ID FIN {$$ = $3;}

```

#### 4.4.8 显示所有表单

显示所有表单不需要额外的字符串参数，直接调用 `show_tables` 即可。

**a.y**

```

1 show_tables: SHOW TABLES FIN

```

#### 4.4.9 展示表单所有属性

展示表单上传表单名称，调用 `desc` 成员函数即可。

**a.y**

```

1 desc_table: DESC ID FIN {$$ = $2;}

```

#### 4.4.10 插入元组

在 insert 功能下，我们支持两种插入模式。一种是不提供属性名，但是必须按顺序给全所有属性值；另一种提供属性名，相应地给出对应的属性值。

`insert_values` 对应解析出 `VALUE (属性值 1, 属性值 2, ...)` 的结构。



而 `insert_method` 则允许用户给出自定义的属性列，属性列之间用逗号分别，所有属性列用括号包围（哪怕是一个属性也需要括号）。

根据插入模式的不同，信息在全局链表中的组成方式也不同。

- 如果插入一条拥有全属性的元组，那么链表第一个元素记录了该表单的名称，用于函数索引出将要操作的表单；随后就给出所有属性域的内容以及长度直到结尾。接口函数会检查插入的属性值的数量是否与该表单拥有的属性属性匹配，但是暂时无法检测属性值的类型也匹配。

- 如果插入一条只拥有部分属性的元组，那么链表第一个元素同样记录了该表单的名称，用于函数索引出将要操作的表单；随后跟一个有效字节长度为 1、指针内容第一个字节为 0 的特征链表元素，表示即将开始属性名的插入。接下来跟着若干个链表元素，存储了将要插入元素拥有的属性名，以一个相同的特征 0 元素结束。当读取到第二个特征 0 元素后，就类似全属性元组结构一样，跟着所有对应的属性值直到文件结尾。

接口函数 `temp_insert` 解析该链表，并进行相应的插入工作，并会对一些常见错误进行汇报。当遇到一个不匹配的属性名后，该函数并不会报错，但也不会进行任何插入工作（此处没有与标准 DBMS 接轨，日后改进），那些没有初始值的属性值均被设置为空，展示时显示为 `(null)`。

a.y

```
1 insert:
2   INSERT INTO ID {
3     temp.clear();
4     temp.push_back($3, strlen($3));
5   }
6   insert_method
7 FIN {
8   if(temp_insert()) $$ = 1;
9   else $$ = 0;
10 }
11
12 insert_method:
13   LB {
14     char *p = (char*)malloc(1); p[0] = 0;
15     temp.push_back(p, 1);
16   }
17   cols RB {
18     char *p = (char*)malloc(1); p[0] = 0;
19     temp.push_back(p, 1);
20   }
21   insert_values
22   |insert_values
```

```

23 cols:
24     ID {
25         temp.push_back($1, strlen($1));
26     }
27     |cols COMMA ID {
28         temp.push_back($3, strlen($3));
29     }
30 insert_values:
31     VALUES LB values RB
32 values:
33     value {
34         temp.push_back($1, strlen($1));
35     }
36     |values COMMA value {
37         temp.push_back($3, strlen($3));
38     }
39 value:
40     INTEGER {
41         $$ = $1;
42     }
43     |STRING {
44         $$ = $1;
45     }

```

#### 4.4.11 查询元组

查询元组涉及到 where 条件子句，这里我们先对该类子句做一个说明。

对于任意涉及到 where 字句的其余 SQL 语句，where 都是可选的，因而在 Yacc 中 `conditions` 的解析也分为了有条件 `with_conditions` 和无条件。

那么有条件可以归约成若干条原子条件（即不包含逻辑与、逻辑或等，只有单独的一个逻辑表达式）的逻辑组合，正如 `with_conditions` 所示。

在本实验中，暂只支持 `OR, AND` 作为逻辑组合方式，`<=` 作为逻辑比较符；尽管 `>, <=, >=` 等其余运算符都能较快的通过 `<` 来转换得到，但是在测试时仅采用了小于、等于这两个逻辑比较符，之后由于时间紧促，便没有来得及新加功能。

在对条件进行归约的时候，我们使用**后缀表达式**解决优先级问题。使用后缀表达式不会产生优先级问题，使用一个栈线性地扫描就可以无歧义地完成解析工作。可以看到，在归约时，归约到 `value` 时会向全局链表中加入数值/属性名，但归约到比较符时不进行任何操作，当整条逻辑比较语句全部读取完

成后才将比较符加入到链表中。通过这样的归约方式，我们就可以获得所有条件的后缀表达式。

因为字节长度的不同，所以接口函数可以通过分析字节长度来判断当前链表元素对应的是一个原子条件，还是一个逻辑组合符号。

接口函数根据该全局链表解析条件，遍历查询的表单下的所有元组，如果当前元组符合所有的条件，则将该元组加入到一个候选元组链表中，该候选元组链表用于随后的输出。

**需要注意**，尽管这套 DBMS 的实现十分朴素，许多地方的时间复杂度、空间复杂度都没有仔细优化，但是这里由于涉及到了元组的复制，我们还是给出了一个较好的方法，以免出现复杂度爆炸：候选元组链表中不照搬复制原有的元组内容，而是存原有元组的指针信息。由于候选元组链表存的只是一个指针，那么其空间消耗至多也就是一个指针的大小加一个 32 位的整数，不会因为元组数据内容的巨大而消耗过多的时间空间。

a.y

```
1 conditions:
2     WHERE with_conditions
3     |
4 with_conditions:
5     with_conditions AND with_conditions {
6         char *p = (char*)malloc(1); p[0] = 1;
7         cons.push_back(p, 1);
8     }
9     |with_conditions OR with_conditions {
10        char *p = (char*)malloc(1); p[0] = 2;
11        cons.push_back(p, 1);
12    }
13    |LB with_conditions RB
14    |ID BELOW value {
15        CONDITION *now = new CONDITION;
16        now->lval = $1; now->lbytes = strlen($1);
17        now->rval = $3; now->rbytes = strlen($3);
18        now->cmp = 1;
19        cons.push_back(now, sizeof(CONDITION));
20    }
21    |ID EQU value {
22        CONDITION *now = new CONDITION;
23        now->lval = $1; now->lbytes = strlen($1);
24        now->rval = $3; now->rbytes = strlen($3);
25        now->cmp = 2;
26        cons.push_back(now, sizeof(CONDITION));
27    }
```

说明完了 where 子句对应的模块，剩余的部分就很简洁了。select 语句要么使用 `*` 要求显示所有的属性，要么给出需要的属性列表，以逗号分隔。本系统暂时只支持单表查询，不支持多表查询、子表查

询等操作。

在数据解析方面，由于有区分数据块归属的困难，我们再使用另一个全局的链表，与 where 子句使用的链表不同。这里使用的全局链表第一个元素给出表单的名称用于索引，随后给出所有希望查询的属性名；而 where 字句的链表专职处理逻辑语句的后缀表达式。

接口函数 `temp_select` 会检查属性名是否正确，并对一些常见错误进行汇报，随后将美观地展示出所有元组的信息。

**a.y**

```
1 select:
2     SELECT{
3         clear_cons();
4         temp.clear();
5     }
6     select_what
7 FROM ID{
8     temp.push_back($5, strlen($5));
9 }
10 conditions FIN{
11     if(temp_select() == 1) $$ = 1;
12     else $$ = 0;
13 }
14 ids:
15     ID {
16         temp.push_back($1, strlen($1));
17     }
18     |ids COMMA ID {
19         temp.push_back($3, strlen($3));
20     }
21 select_what:
22     STAR
23     |ids
```

#### 4.4.12 退出系统

读取到退出指令，系统先妥善地将内存数据库存储到本地，再安全地退出。

**a.y**

```
1 exit:
2     EXIT
```

#### 4.4.13 删除元组与更新元组

由于时间紧张等原因，delete 和 update 操作并没有完成语义翻译，该系统暂时只能用 Yacc 解析出两种语句的结构，不能执行实际的操作，十分遗憾！

##### a.y delete

```
1 delete:
2     DELETE FROM ID conditions FIN;
```

##### a.y update

```
1 update:
2     UPDATE ID
3     SET set_values
4     conditions FIN;
5 set_values:
6     set_value
7     |set_values COMMA set_value
8 set_value:
9     ID EQU value
```

## 5 结果展示

使用如下的测试数据作为标准输入。

##### in

```
1 create database db;
2 show databases;
3 use db;
4
5 create table tab(
6     id char(100),
7     name char(100),
8     age int
9 );
10
11 show tables;
12 desc tab;
13
14 insert into tab values ("id0", "person0", 10);
15 insert into tab values ("id1", "person1", 20);
16 insert into tab values ("id2", "person2", 30);
17 insert into tab values ("id3", "person3", 40);
18
19 select * from tab;
```

```

20 select id, age from tab;
21
22 select * from tab where age < 20;
23 select * from tab where (age < 20);
24 select * from tab where age < 20 OR id < "id3";
25 select * from tab where (age < 20 OR id < "id3") AND name = "person2";
26
27 update tab set id = "idx";
28 update tab set id = "idx" where id < "id2";
29
30 exit

```

测试结果如下所示，该系统能正确地执行上述语句。

```

rqdmap@rqdmap-MacBookPro: ~/codes/SQL_Interpreter/src
rqdmap@rqdmap-MacBookPro:~/codes/SQL_Interpreter/src$ ./main < in
Hello
create_database ok.
12
+-----+
| Database |
+-----+
| db       |
+-----+
use_database ok.
create_table ok.
+-----+
| Tables_in_db |
+-----+
| tab          |
+-----+
+-----+-----+
| Field | Type      |
+-----+-----+
| id    | char(100) |
| name  | char(100) |
| age   | int       |
+-----+-----+
insert_into ok.
insert_into ok.
insert_into ok.
insert_into ok.
+-----+-----+-----+
| id    | name      | age |
+-----+-----+-----+
| "id0" | "person0" | 10  |
| "id1" | "person1" | 20  |
| "id2" | "person2" | 30  |
| "id3" | "person3" | 40  |
+-----+-----+-----+
Select ok
+-----+-----+

```

```
rqdmap@rqdmap-MacBookPro: ~/codes/SQL_Interpreter/src
+-----+
| id | age |
+-----+
| "id0" | 10 |
| "id1" | 20 |
| "id2" | 30 |
| "id3" | 40 |
+-----+
Select ok
+-----+-----+
| id | name | age |
+-----+-----+
| "id0" | "person0" | 10 |
+-----+-----+
Select ok
+-----+-----+
| id | name | age |
+-----+-----+
| "id0" | "person0" | 10 |
+-----+-----+
Select ok
+-----+-----+
| id | name | age |
+-----+-----+
| "id0" | "person0" | 10 |
| "id1" | "person1" | 20 |
| "id2" | "person2" | 30 |
+-----+-----+
Select ok
+-----+-----+
| id | name | age |
+-----+-----+
| "id2" | "person2" | 30 |
+-----+-----+
Select ok
Bye
rqdmap@rqdmap-MacBookPro:~/codes/SQL_Interpreter/src$
```

## 6 局限与展望

- 该系统功能尚十分不完善，没有实现更新与删除操作，支持的数据类型、操作符、比较符不全面，不支持多表查询、子表查询，不支持主码、外码，不支持 check 语句，不支持数值表达式运算等。
- 该系统某些实现并没有与标准 DBMS 匹配，不能执行所有的标准 SQL 语句。
- 支持的条件类型十分单薄，不支持左值和右值的任意选择，可通过在条件结构体中添加数据类型进一步进行划分。
- 涉及到指定属性的函数操作并没有对该属性对应的数据类型进行合法性检查，默认用户遵循规范。
- 时空复杂度仍可以大力优化，并对代码和程序结构进一步美化。

## 7 个人体会

本次实验我受益匪浅，熟悉了 Lex/Yacc 分析工具的使用，扎实了 C++ 编程的功底，也提高了 Latex 书写文档的能力。半年来这是第三份千行代码以上的程序项目，从第一份 RQD-FTP 开始甚至不会写 C++ 多文件编译、乱在 h 文件中写实现，到第二份二进制表示的高精度整数类，再到现在这份垃圾的 DBMS 全套系统，果然大工程使人提高，在神秘的方面提高了退役 ACM 选手的代码能力。正如小兵老师所说，这次实验项目令人难忘，自己实现 DBMS 十分有趣，但是实现出来的东西却差强人意。尽管有一千行左右但也就是我 2 天日夜肝出来的结果，日后如果有时间还希望进一步优化，能真正媲美一个标准的 DBMS。当然，时空复杂度方面就不妄想与之比肩了，功能完备即心满意足。

## 8 程序清单

程序名	描述
MakeFile	
rqdmap.h	个人编程惯用的一些宏定义
list.h, list.cpp	自定义链表结构接口及定义
a.l	Lex 源文件
lex.yy.c	Lex 编译结果
a.y	Yacc 文件
y.tab.h, y.tab.c	Yacc 编译结果
sql.h, sql.cpp	SQL 函数接口及定义
main	可执行程序，数据库系统
in	测试数据
init	bash 脚本，初始化数据库
DBMS	数据库本地文件