

Course Project: An Interpreter for SSQL

Design Document

Kuang Yuanyuan

January 2, 2015

CHAPTER 1

INFORMATION

1.1 Team Information

Name	Student ID	E-mail
Sun Jiacheng	12330285	291624707@qq.com
Kuang Yuanyuan	12330153	596755905@qq.com
Qiu Zhilin	12330268	847300960@qq.com
Sun Dongliang	12330284	1255993541@qq.com
Wang Kaibin	12330305	wkbpluto@qq.com

1.2 Tarball Structure

/doc -- documents of project

/src -- source code

1.3 Compile

cd ./src

make -- compile SimpleDB

make test -- compile and run sample test

make tar -- tar the project.

CHAPTER 2

PROJECT

2.1 Project Website

<https://github.com/thomasking0529/SimpleDB>

2.2 Code Style

K&R(CDT*s default style)

2.3 APIs

2.3.1 Lexer

enum TokenType:

KEYWORD -- select from where delete table create int values primary
key default into insert case insensitive

ID -- id (identifier) is a sequence of digits, underline and letters. All
identifiers should start with a letter or an underline. The
maximum length of an identifier is 64.

NUM -- num (number) is a sequence of digits. (of 32-bits) only integers.

OP -- Arithmetical operators: +, -, *, /, unary -, unary +

Relational operators: <, >, <=>, ==, >=, <=

Logical operators: &&, ||, !

Assignment operator: =.

Basic punctuation("(", ")", ",", ";")

struct Token:

TokenType type -- token type.

std::string value -- token value, store the original string of token.

class Lexer:

std::list<Token> GetTokens(const std::string& s) -- get a list of tokens
from input string

2.3.2 Parser

enum Action:

```
CREATE -- create table
INSERT -- insert one row to table
DELETE -- delete row(s) from table
SELECT -- query row(s) from table
INVALID -- if unknown keywords detected
```

enum Op:

```
PLUS, // +, both unary and binary
MINUS, // -, both unary and binary
MULTIPLY, // *
DIVIDE, // /
LT, // <
GT, // >
NE, // <>
E, // ==
GTE, // >=
LTE, // <=
AND, // &&
OR, // ||
NOT, // !
EQ, // =
LB, // (
RB, // )
COMMA, // ,
```

struct Condition:

```
Condition* lop-- left subtree
Condition* rop-- right subtree
Op op -- operator
std::string opd - operand
```

struct Property:

```
std::string id -- property id
int default_value -- default value
operator== -- operator overloading
```

struct Statement:

Action act -- action.

std::string table -- table to operate on.

std::list<Property> prop_list --property to return or add, for create and select

std::vector<std::string> key_idx -- list of keys

std::vector<int> value_list -- list of values

Condition* cond -- where clauses, for select and delete

class Polish:

Polish() -- Constructor, initialize private variables

void neglect() -- set neglect flag, for neglect numbers

void Insert(const std::string& item) -- inset a token

int Calculate() -- calculate the expression

Condition* buildTree -- build a condition tree

class Parser:

Parser() -- Constructor, initialize private variables

Statement Parse(const std::string& s) -- main parser

2.3.3 Core

struct Row:

int key -- key value, for identify row

std::vector<int> -- key indices

std::vector<int> cols -- values

struct Table:

std::string id -- table name

std::vector<Property> -- properties of table

std::vector<int> -- key indices of table

unsigned long long count -- if no primary key, use this

std::set<Row> rows -- set of rows, use Row.key for key

void Insert(const std::list<int>& record) -- insert

void Delete(const Condition* cond) -- delete

std::vector<int> Query(const Condition* cond) -- query

class SimpleDB:

SimpleDB() -- Constructor, initialize private variables

void Execute(const std::string& stmt) -- execute raw input

2.4 Design

2.4.1 Lexer

Lexer reads a single string of statement and extracts tokens from it.

Finite State Automata:

read a character from input string;

check character for:

a-z *A-Z* *0-9*: concat to temp string

split characters: space, \t, \n, save temp string

symbols: save temp string

2.4.2 Parser

Design

Use LL(1) parser.

LL Parser:

An LL parser is called an LL(k) parser if it uses k tokens of look ahead when parsing a sentence. If such a parser exists for a certain grammar and it can parse sentences of this grammar without backtracking then it is called an LL(k) grammar. LL parser can only parse languages that have LL(k) grammars without ϵ -rules. LL(k) grammars without ϵ -rules can generate more languages the higher the number k of look ahead tokens. A corollary of this is that not all context-free languages can be recognized by an LL(k) parser. An LL parser is called an LL(*) parser (an LL-regular parser) if it is not restricted to a finite k tokens of look ahead, but can make parsing decisions by recognizing whether the following tokens belong to a regular language (for example by means of a Deterministic Finite Automaton).

Using the method introduced in the textbook, we construct the following First and Follow table:

	First	Follow
ssql_stmt	create insert delete select	\$
create_stmt	create	\$
decl_list	id primary)
decl_list*	, ϵ)
decl	id primary	,)
default_spec	default ϵ	,)
column_list	id) from
column_list*	, ϵ) from
insert_stmt	insert	\$
value_list	(- + num)

value_list*	, ϵ)
simple_expr	(-+ num	,)
simple_expr*	+ - ϵ	,)
simple_term	(-+ num	+ -,)
simple_term*	ϵ /*	+ -,)
simple_unary	(-+ num	/ + -,) *
delete_stmt	delete	\$
where_clause	where ϵ	;
disjunct	(!-+ id num	;)
disjunct*	ϵ	;)
conjunct	(! -+ id num	;)
conjunct*	&& ϵ	;)
bool	(! -+ id num	&& ;)
rop	== > < >= <= <>	(-+ id num
query_list	select	\$
select_list	id *	From
expr	(-+ id num	&& == > < >= <= <>;)
expr*	+ - ϵ	&& == > < >= <= <>;)
term	(+ -id num	;) && <> == > < >= <= + -
term*	ϵ /*	;) && <> == > < >= <= + -
unary	(-+ id num	;) && <> == > < >= <= + - * /
comp	(+ -id num	&& ;)

And with the help of the method in the textbook, we generate the following Prediction table:

M[ssql_stmt, create] = ssql_stmt create_stmt
 M[ssql_stmt, insert] = ssql_stmt insert_stmt
 M[ssql_stmt, delete] = ssql_stmt delete_stmt
 M[ssql_stmt, select] = ssql_stmt query_stmt
 M[create_stmt, create] = create_stmt create table id (decl_list);
 M[create_stmt, create] = create_stmt create table id (decl_list);
 M[decl_list, id] = decl_list decl decl_list*
 M[decl_list, primary] = decl_list decl decl_list*
 ...

The remainder is in the Appendix at the end of the document.

The table is stored and read for parsing in the file "Rule.txt" and terminals are stored in the file "Terminals".

Implementation

- Basic Algorithm

```
Initialization:
    List token = lexer.getTokens();
    Map table = readFile"Rule.txt"
    Map action;
    Stack procedure;
    Stack father;
    procedure.push($);
    procedure.push(ssql_stmt);
    father.push($);
    father.push(ssql_stmt);
    ip = token.first;
    X = procedure.top;
    F = father.top;
Loop
    father.pop;
    procedure.pop;
    If (X match ip) {
        ip = token.next;
        action[ip].translate(ip, F);
    } else if (ip is a terminals but no match X) {
        Error();
    } else if (X is  $\epsilon$ ) {
        do nothing
    } else if (table[X, ip] =  $Y_1 Y_2 \dots Y_k$ ) {
        procedure.push( $Y_k, Y_{k-1} \dots Y_1$ );
        father.push(X, X, ..., X);
    }
    X = procedure.top;
    F = father.top;
Until X = $
if token is not empty, Error;
```

- Implement detail

The implement of the parser is as the pseudocode above, which only use one more stack to store the Grammar Productions' left part, it's still an LL(1) grammar as we still only look one more token to decide which production to use. The father stack store the Grammar productions' left part of the corresponding symbol in the procedure stack, just to help with the translation. Then we can translate by tokens seeing it's father.

- Translation problem

We generate a structure Statement when parsing and the SimpleDB's database manager will use the Statement to do query create etc. As was mentioned before, with the help of the Grammar productions' left part which the token is matched on, we could get the meaning of the token and generate the statement. However, when it comes to the expression, we use the algorithm of reverse polish expression, which implements as class Polish. When a token is seem as part of a expression, we insert it to a object of Polish, and for the math expression, calculate the result, for bool expression, build a condition tree in the Statement.

- Detail of translation:

Use a object of Polish to reduce the where clause. Use reverse polish to check the priority of braces. Set a flag for neglect numbers, if unary plus encountered, discard it. If a unary minus encountered, if a number followed, flag = !flag, this help to reduce expression like ---1; if an id followed, expand the expression to (0-id); if a left brace followed, expand to " (0-" , count the left braces, count unary braces followed by left brace and insert the ")" in correct place. If unit operator ! found, make it become a binary one and always give another operand "!!", and when building the condition tree, recognize the "!!" and do something.

2.4.3 Core

Core part of SimpleDB is supposed to manage the database in memory. Get the Statement instance from parser to Execute order.

SimpleDB a set of Table, Table manage a set of Row. Do query, insert, delete on Table.

CHAPTER 3

PROJECT SCHEDULE

3.1 Stage 1

3.1.1 Project architecture and APIs design

Author Deadline

Sun Jiacheng 12.5.2014

3.1.2 Lexer implementation

Author Deadline

Sun Jiacheng 12.5.2014

3.1.3 Group Discussion

Contents:

Stage 2 and Stage 3 work distribution.

Date:

12.6.2014

3.2 Stage 2

3.2.1 Parser

Author Deadline

Qiu Zhilin 12.23.2014

3.2.2 Core

Author Deadline

Sun Dongliang 12.23.2014

3.3 Stage 3

3.3.1 Test Document

Author Deadline

Wang Kaibin 1.2.2014

3.3.2 Design Document

Author Deadline

Kuang Yuanyuan 1.2.2014

CHAPTER 4

APPENDIX

4.1 Prediction table

$M[\text{ssql_stmt}, \text{create}] = \text{ssql_stmt} \rightarrow \text{create_stmt}$
 $M[\text{ssql_stmt}, \text{insert}] = \text{ssql_stmt} \rightarrow \text{insert_stmt}$
 $M[\text{ssql_stmt}, \text{delete}] = \text{ssql_stmt} \rightarrow \text{delete_stmt}$
 $M[\text{ssql_stmt}, \text{select}] = \text{ssql_stmt} \rightarrow \text{query_stmt}$
 $M[\text{create_stmt}, \text{create}] = \text{create_stmt} \rightarrow \text{create table id (decl_list);}$
 $M[\text{decl_list}, \text{id}] = \text{decl_list} \rightarrow \text{decl decl_list}^*$
 $M[\text{decl_list}, \text{primary}] = \text{decl_list} \rightarrow \text{decl decl_list}^*$
 $M[\text{decl_list}^*, ,] = \text{decl_list}^* \rightarrow , \text{decl decl_list}^*$
 $M[\text{decl_list}^*,)] = \text{decl_list}^* \rightarrow \epsilon$
 $M[\text{decl}, \text{id}] = \text{decl} \rightarrow \text{id int default_spec}$
 $M[\text{decl}, \text{primary}] = \text{decl} \rightarrow \text{primary key (column_list);}$
 $M[\text{column_list}, \text{id}] = \text{column_list} \rightarrow \text{id column_list}^*$
 $M[\text{column_list}^*, ,] = \text{column_list}^* \rightarrow , \text{id column_list}^*$
 $M[\text{column_list}^*,)] = \text{column_list}^* \rightarrow \epsilon$
 $M[\text{column_list}^*, \text{from}] = \text{column_list}^* \rightarrow \epsilon$
 $M[\text{default_spec}, \text{default}] = \text{default_spec} \quad \text{default} = \text{simple_expr}$
 $M[\text{default_spec}, ,] = \text{default_spec} \rightarrow \epsilon$
 $M[\text{default_spec},)] = \text{default_spec} \rightarrow \epsilon$
 $M[\text{simple_expr}, (] = \text{simple_expr} \rightarrow \text{simple_term simple_expr}^*$
 $M[\text{simple_expr}, +] = \text{simple_expr} \rightarrow \text{simple_term simple_expr}^*$
 $M[\text{simple_expr}, -] = \text{simple_expr} \rightarrow \text{simple_term simple_expr}^*$
 $M[\text{simple_expr}, \text{num}] = \text{simple_expr} \rightarrow \text{simple_term simple_expr}^*$
 $M[\text{simple_expr}^*, +] = \text{simple_expr}^* \rightarrow + \text{simple_term simple_expr}^*$
 $M[\text{simple_expr}^*, -] = \text{simple_expr}^* \rightarrow - \text{simple_term simple_expr}^*$
 $M[\text{simple_expr}^*,)] = \text{simple_expr}^* \rightarrow \epsilon$
 $M[\text{simple_expr}^*, ,] = \text{simple_expr}^* \rightarrow \epsilon$
 $M[\text{simple_term}, (] = \text{simple_term} \rightarrow \text{simple_unary simple_term}^*$
 $M[\text{simple_term}, -] = \text{simple_term} \rightarrow \text{simple_unary simple_term}^*$
 $M[\text{simple_term}, +] = \text{simple_term} \rightarrow \text{simple_unary simple_term}^*$
 $M[\text{simple_term}, \text{num}] = \text{simple_term} \rightarrow \text{simple_unary simple_term}^*$
 $M[\text{simple_term}^*, *] = \text{simple_term}^* \rightarrow * \text{simple_unary simple_term}^*$
 $M[\text{simple_term}^*, /] = \text{simple_term}^* \rightarrow / \text{simple_unary simple_term}^*$
 $M[\text{simple_term}^*,)] = \text{simple_term}^* \rightarrow \epsilon$
 $M[\text{simple_term}^*, ,] = \text{simple_term}^* \rightarrow \epsilon$
 $M[\text{simple_term}^*, +] = \text{simple_term}^* \rightarrow \epsilon$
 $M[\text{simple_term}^*, -] = \text{simple_term}^* \rightarrow \epsilon$
 $M[\text{simple_unary}, (] = \text{simple_unary} \rightarrow (\text{simple_expr})$
 $M[\text{simple_unary}, -] = \text{simple_unary} \rightarrow - \text{simple_unary}$

$M[\text{simple_unary}, +] = \text{simple_unary} \rightarrow +\text{simple_unary}$
 $M[\text{simple_unary}, \text{num}] = \text{simple_unary} \rightarrow \text{num}$
 $M[\text{insert_stmt}, \text{insert}] = \text{insert_stmt} \rightarrow \text{insert into id (column_list) values (value_list);}$
 $M[\text{value_list}, ()] = \text{value_list} \rightarrow \text{simple_expr value_list}^*$
 $M[\text{value_list}, +] = \text{value_list} \rightarrow \text{simple_expr value_list}^*$
 $M[\text{value_list}, -] = \text{value_list} \rightarrow \text{simple_expr value_list}^*$
 $M[\text{value_list}, \text{num}] = \text{value_list} \rightarrow \text{simple_expr value_list}^*$
 $M[\text{value_list}^*, ,] = \text{value_list}^* \rightarrow , \text{simple_expr value_list}^*$
 $M[\text{value_list}^*,)] = \text{value_list}^* \rightarrow \epsilon$
 $M[\text{delete_stmt}, \text{delete}] = \text{delete_stmt} \rightarrow \text{delete from id where_clause;}$
 $M[\text{where_clause}, \text{where}] = \text{where_clause} \rightarrow \text{where disjunct}$
 $M[\text{where_clause}, ;] = \text{where_clause} \rightarrow \epsilon$
 $M[\text{disjunct}, ()] = \text{disjunct} \rightarrow \text{conjunct disjunct}^*$
 $M[\text{disjunct}, -] = \text{disjunct} \rightarrow \text{conjunct disjunct}^*$
 $M[\text{disjunct}, +] = \text{disjunct} \rightarrow \text{conjunct disjunct}^*$
 $M[\text{disjunct}, \text{id}] = \text{disjunct} \rightarrow \text{conjunct disjunct}^*$
 $M[\text{disjunct}, \text{num}] = \text{disjunct} \rightarrow \text{conjunct disjunct}^*$
 $M[\text{disjunct}, !] = \text{disjunct} \rightarrow \text{conjunct disjunct}^*$
 $M[\text{disjunct}^*, ||] = \text{disjunct}^* \rightarrow || \text{conjunct disjunct}^*$
 $M[\text{disjunct}^*,)] = \text{disjunct}^* \rightarrow \epsilon$
 $M[\text{disjunct}^*, ;] = \text{disjunct}^* \rightarrow \epsilon$
 $M[\text{conjunct}, -] = \text{conjunct} \rightarrow \text{bool conjunct}^*$
 $M[\text{conjunct}, +] = \text{conjunct} \rightarrow \text{bool conjunct}^*$
 $M[\text{conjunct}, \text{id}] = \text{conjunct} \rightarrow \text{bool conjunct}^*$
 $M[\text{conjunct}, \text{num}] = \text{conjunct} \rightarrow \text{bool conjunct}^*$
 $M[\text{conjunct}, ()] = \text{conjunct} \rightarrow \text{bool conjunct}^*$
 $M[\text{conjunct}, !] = \text{conjunct} \rightarrow \text{bool conjunct}^*$
 $M[\text{conjunct}^*, \&\&] = \text{conjunct}^* \rightarrow \&\& \text{bool conjunct}^*$
 $M[\text{conjunct}^*,)] = \text{conjunct}^* \rightarrow \epsilon$
 $M[\text{conjunct}^*, ||] = \text{conjunct}^* \rightarrow \epsilon$
 $M[\text{conjunct}^*, ;] = \text{conjunct}^* \rightarrow \epsilon$
 $M[\text{bool}, ()] = \text{bool} \rightarrow (\text{disjunct})$
 $M[\text{bool}, !] = \text{bool} \rightarrow !\text{bool}$
 $M[\text{bool}, -] = \text{bool} \rightarrow \text{comp}$
 $M[\text{bool}, +] = \text{bool} \rightarrow \text{comp}$
 $M[\text{bool}, \text{id}] = \text{bool} \rightarrow \text{comp}$
 $M[\text{bool}, \text{num}] = \text{bool} \rightarrow \text{comp}$
 $M[\text{comp}, -] = \text{comp} \rightarrow \text{expr rop expr}$
 $M[\text{comp}, +] = \text{comp} \rightarrow \text{expr rop expr}$
 $M[\text{comp}, \text{id}] = \text{comp} \rightarrow \text{expr rop expr}$
 $M[\text{comp}, \text{num}] = \text{comp} \rightarrow \text{expr rop expr}$
 $M[\text{rop}, <>] = \text{rop} \rightarrow <>$
 $M[\text{rop}, ==] = \text{rop} \rightarrow ==$
 $M[\text{rop}, <] = \text{rop} \rightarrow <$

$M[\text{rop}, >] = \text{rop} \rightarrow >$
 $M[\text{rop}, >=] = \text{rop} \rightarrow >=$
 $M[\text{rop}, <=] = \text{rop} \rightarrow <=$
 $M[\text{expr}, -] = \text{expr} \rightarrow \text{term expr}^*$
 $M[\text{expr}, +] = \text{expr} \rightarrow \text{term expr}^*$
 $M[\text{expr}, \text{id}] = \text{expr} \rightarrow \text{term expr}^*$
 $M[\text{expr}, \text{num}] = \text{expr} \rightarrow \text{term expr}^*$
 $M[\text{expr}^*, +] = \text{expr}^* \rightarrow + \text{term expr}^*$
 $M[\text{expr}^*, -] = \text{expr}^* \rightarrow - \text{term expr}^*$
 $M[\text{expr}^*,)] = \text{expr}^* \rightarrow \epsilon$
 $M[\text{expr}^*, ;] = \text{expr}^* \rightarrow \epsilon$
 $M[\text{expr}^*, ||] = \text{expr}^* \rightarrow \epsilon$
 $M[\text{expr}^*, \&\&] = \text{expr}^* \rightarrow \epsilon$
 $M[\text{expr}^*, <>] = \text{expr}^* \rightarrow \epsilon$
 $M[\text{expr}^*, ==] = \text{expr}^* \rightarrow \epsilon$
 $M[\text{expr}^*, >] = \text{expr}^* \rightarrow \epsilon$
 $M[\text{expr}^*, <] = \text{expr}^* \rightarrow \epsilon$
 $M[\text{expr}^*, >=] = \text{expr}^* \rightarrow \epsilon$
 $M[\text{expr}^*, <=] = \text{expr}^* \rightarrow \epsilon$
 $M[\text{term}, -] = \text{term} \rightarrow \text{unary term}^*$
 $M[\text{term}, +] = \text{term} \rightarrow \text{unary term}^*$
 $M[\text{term}, \text{id}] = \text{term} \rightarrow \text{unary term}^*$
 $M[\text{term}, \text{num}] = \text{term} \rightarrow \text{unary term}^*$
 $M[\text{term}^*,)] = \text{term}^* \rightarrow \epsilon$
 $M[\text{term}^*, ;] = \text{term}^* \rightarrow \epsilon$
 $M[\text{term}^*, ||] = \text{term}^* \rightarrow \epsilon$
 $M[\text{term}^*, \&\&] = \text{term}^* \rightarrow \epsilon$
 $M[\text{term}^*, <>] = \text{term}^* \rightarrow \epsilon$
 $M[\text{term}^*, ==] = \text{term}^* \rightarrow \epsilon$
 $M[\text{term}^*, >] = \text{term}^* \rightarrow \epsilon$
 $M[\text{term}^*, <] = \text{term}^* \rightarrow \epsilon$
 $M[\text{term}^*, >=] = \text{term}^* \rightarrow \epsilon$
 $M[\text{term}^*, <=] = \text{term}^* \rightarrow \epsilon$
 $M[\text{term}^*, +] = \text{term}^* \rightarrow \epsilon$
 $M[\text{term}^*, -] = \text{term}^* \rightarrow \epsilon$
 $M[\text{term}^*, *] = \text{term}^* \rightarrow * \text{unary term}^*$
 $M[\text{term}^*, /] = \text{term}^* \rightarrow / \text{unary term}^*$
 $M[\text{unary}, -] = \text{unary} \rightarrow \text{unary}$
 $M[\text{unary}, +] = \text{unary} \rightarrow + \text{unary}$
 $M[\text{unary}, \text{id}] = \text{unary} \rightarrow \text{id}$
 $M[\text{unary}, \text{num}] = \text{unary} \rightarrow \text{num}$
 $M[\text{query_list}, \text{select}] = \text{query_list} \rightarrow \text{select select_list from id where_clause;}$
 $M[\text{select_list}, *] = \text{select_list} \rightarrow *$
 $M[\text{select_list}, \text{id}] = \text{select_list} \rightarrow \text{column_list}$