



华南师范大学

本科学生实验（实践）报告

院 系：计算机学院

实验课程：	编译原理
实验项目：	XLEX-词法自动生成器
指导老师：	黄煜廉
开课时间：	2023~2024年度第一学期
专 业：	计算机科学与技术
班 级：	计算机科学与技术3班

华南师范大学教务处

华南师范大学实验报告

学生姓名	卢泓钢	学号	20212131096
专业	计算机科学与技术	年级、班级	2021级计科3班
课程名称	编译原理	实验项目	XLEX-词法自动生成器
实验类型	综合	实验时间	2023 年 9 月 28 日
实验指导老师	黄煜廉	实验评分	

华南师范大学实验报告

1 实验内容

2 实验目的

3 实验文档

3.1 实验文档：基于Qt的C++ XLEX-词法自动生成器

3.2 引言

3.3 数据结构与设计思路

3.4 实现细节

3.4.1 正则表达式到NFA

3.4.2 NFA到DFA

3.4.3 DFA到最小化DFA

3.4.4 DFA到词法分析程序

3.5 测试

4 实验总结

5 参考文献

1 实验内容

设计一个应用软件，以实现将正则表达式-->NFA--->DFA-->DFA最小化-->词法分析程序

必做内容

- (1) 正则表达式应该支持单个字符，运算符有： 连接、选择 (|)、闭包 (*)、括号 ()、可选 (?)
- (2) 要提供一个源程序编辑界面，让用户输入一行（一个）或多行（多个）正则表达式（可保存、打开正则表达式文件）
- (3) 需要提供窗口以使用户可以查看转换得到的**NFA**（用状态转换表呈现即可）
- (4) 需要提供窗口以使用户可以查看转换得到的**DFA**（用状态转换表呈现即可）
- (5) 需要提供窗口以使用户可以查看转换得到的最小化**DFA**（用状态转换表呈现即可）
- (6) 需要提供窗口以使用户可以查看转换得到的词法分析程序（该分析程序需要用C/C++语言描述）
- (7) 用户界面应该是**windows**界面
- (9) 应该书写完善的软件文档

选做实验

- (1) 扩充正则表达式的运算符，如 []、正闭包 (+) 等。

2 实验目的

设计一个应用软件，以实现将正则表达式-->NFA--->DFA-->DFA最小化-->词法分析程序

3 实验文档

3.1 实验文档：基于Qt的C++ XLEX-词法自动生成器

3.2 引言

本实验旨在开发一个基于Qt的XLEX-词法自动生成器，该软件能够将用户输入的正则表达式转化为词法分析程序。具体步骤包括从正则表达式构建非确定性有限自动机（NFA），将NFA转换为确定性有限自动机（DFA），并对DFA进行最小化处理，最后生成相应的词法分析程序。文件解析器在软件工程规范下开发，具有高度的可维护性和扩展性。

3.3 数据结构与设计思路

核心过程由 xlexgenerator.h 中的 XLEXGenerator 类来实现，简要声明如下：

```
1 class XLEXGenerator {
2 public:
3     XLEXGenerator() {}
4     XLEXGenerator(string reg) {regExp = reg;}
5     void init();    // 初始化
```

```

6      void reg2NFA(); // Thompson构造法
7      void NFA2DFA(); // 子集构造法
8      void minimizeDFA(); // 最小化DFA
9      void DFA2XLEX(); // DFA转换为词法分析程序
10
11 public:
12     /* get成员变量 / set成员变量 (略) */
13
14 private:
15     /* 成员变量 */
16     /* reg to NFA */
17     int idNFA = 1; // 节点编号
18     int start; // 起始状态
19     int end; // 终止状态
20     string regExp; // 正则表达式
21     set<char> charset; // 字符集
22     map<int, map<char, set<int>>> NFA; // NFA
23
24     /* NFA to DFA */
25     int idDFA = 1; // DFA节点编号
26     map<int, map<char, int>> DFA; // DFA
27     map<set<int>, int> stateMap; // 状态集合映射
28     map<int, set<int>> stateMapReverse; // 状态映射集合
29     set<int> DFAendSet; // DFA终止状态集合
30
31     /* 最小化 DFA */
32     int idMinDFA = 1; // 最小化DFA节点编号
33     map<int, map<char, int>> minDFA; // 最小化DFA
34     set<int> minDFAendSet; // minDFA终止状态集合
35
36     /* DFA to 词法分析程序 */
37     string xlex; // 词法分析程序
38
39 private:
40     /* 工具函数 (略) */
41 };

```

public成员函数 `void reg2NFA()` 使用Thompson构造法，实现正则表达式到NFA的转换； `void NFA2DFA()` 使用子集构造法，实现NFA到DFA的转换； `void minimizeDFA()` 使用逐一合并状态的方法，实现DFA到最小化DFA的转换； `void DFA2XLEX()` 实现DFA到词法分析程序的转换。而函数 `void init()` 用于初始化参数，并根据传入的正则表达式，依次调用上述函数，实现正则表达式到词法分析程序的转换。

主要的数据结构

- 使用 `map<int, map<char, set<int>>>` 的数据结构存储NFA，表示NFA状态 u 通过字符 c 能到达的状态集合 s ，即 `NFA[u][c] = s`。
- 使用 `map<int, map<char, int>>` 的数据结构存储DFA，表示DFA状态 u 通过字符 c 能到达的状态 v ，即 `DFA[u][c] = v`。
- 使用 `map<int, map<char, int>>` 的数据结构存储最小化DFA，表示最小化DFA状态 u 通过字符 c 能到达的状态 v ，即 `minDFA[u][c] = v`。

- 使用 `map<set<int>, int>` 的数据结构存储NFA状态集合与DFA状态的映射，即 `stateMap[u] = v` 表示NFA状态集合u对应的DFA状态为v。
- 使用 `map<int, set<int>>` 的数据结构存储DFA状态与NFA状态集合的映射，即 `stateMapReverse[v] = u` 表示DFA状态v对应的NFA状态集合为u。
- 使用 `set<int>` 的数据结构存储DFA终止状态集合，即 `DFAendSet`。
- 使用 `set<int>` 的数据结构存储最小化DFA终止状态集合，即 `minDFAendSet`。

软件的总体结构可以分为以下模块：

- 用户界面模块：提供正则表达式输入、编辑、打开和保存的功能。
- 编译模块：包括正则表达式到 NFA 转换、NFA 到 DFA 转换、DFA 最小化以及词法分析程序生成。
- 输出模块：展示 NFA、DFA、最小化 DFA 的状态转换表，以及生成的词法分析程序。

3.4 实现细节

3.4.1 正则表达式到NFA

函数 `void reg2NFA()` 使用Thompson构造法，实现正则表达式到NFA的转换。Thompson构造法的核心思想是通过双栈操作，将正则表达式转换为NFA。一个是操作符栈，一个是操作数栈。其中操作数栈存储的是当前已经构造好的NFA的初态与终态。

为了简化代码，我实现了工具函数 `pii binOp(char op, pii a, pii b)` 和 `pii monOp(char op, pii a)`，用于处理双目运算符和单目运算符：

```

1  pii XLEXGenerator::binOp(char op, pii a, pii b) {
2      switch (op) {
3          case '.': return opConnect(a, b);
4          case '|': return opUnion(a, b);
5      } return {};
6  }
7
8  pii XLEXGenerator::monOp(char op, pii a) {
9      switch (op) {
10         case '*': return opClosure(a);
11         case '?': return opOption(a);
12         case '+': return opPosClosure(a);
13     } return {};
14 }

```

这两个函数根据不同的操作符分别调用不同的工具函数，以函数 `pii opUnion(pii a, pii b)` 为例，其实现如下：

```

1 // 并操作, 创建两个新的节点, 将其连接到a和b的起始状态, 将a和b的终止状态连接到另一个新的节点
2 pii XLEXGenerator::opUnion(pii a, pii b) {
3     NFA[idNFA][ '# '].insert(a.first);
4     NFA[idNFA++][ '# '].insert(b.first);
5     NFA[a.second][ '# '].insert(idNFA);
6     NFA[b.second][ '# '].insert(idNFA++);
7     return {idNFA - 2, idNFA - 1};
8 }

```

我对正则表达式扩充了正闭包 (+) 运算, 添加函数 `pii opPosClosure(pii a)`, 其实现如下:

```

1 // 正闭包操作, 创建两个新的节点, 将其连接到a的起始状态, 将a的终止状态连接到另一个新的节点, 将a
  的终止状态连接到a的起始状态
2 pii XLEXGenerator::opPosClosure(pii a) {
3     NFA[idNFA][ '# '].insert(a.first);
4     NFA[a.second][ '# '].insert(a.first);
5     NFA[a.second][ '# '].insert(idNFA + 1);
6     return {idNFA++, idNFA++};
7 }

```

对于正则表达式中的连接运算, 因为在原始正则表达式中没有专门的连接运算符, 所以需要在正则表达式中添加连接符, 即将正则表达式中的 `ab` 转换为 `a.b`。函数 `void addConnector()` 用于为正则表达式添加连接符, 具体实现如下:

```

1 void XLEXGenerator::addConnector() {
2     for (int i = 0; i < regExp.length() - 1; i++) {
3         char it = regExp[i], next = regExp[i + 1];
4         if (it == '(' || it == '|')
5             continue;
6         if (next == ')' || next == '|' || next == '*' || next == '+' || next == '?')
7             continue;
8         regExp.insert(i + 1, ".");
9         i++;
10    }
11 }

```

综上所述, 函数 `void reg2NFA()` 具体实现如下:

```

1 // 使用Thompson构造法, 通过双栈操作, 将正则表达式转换为NFA
2 void XLEXGenerator::reg2NFA() {
3     addConnector(); // 为正则表达式添加连接符
4     stack<pii> nodeStack; // 操作数栈 <起始状态, 终止状态>
5     stack<char> opStack; // 操作符栈
6     for (int i = 0; i < regExp.length(); i++) {
7         char it = regExp[i];
8         if (it == '(') {
9             opStack.push(it);
10        } else if (it == ')') { // 遇到右括号, 弹出操作符栈中的操作符, 直到遇到左括号
11            while (opStack.top() != '(') {
12                /* 将运算符与操作数取出来运算 */
13            }
14            opStack.pop(); // 弹出左括号
15        } else if (it == '*' || it == '.' || it == '|' || it == '?' || it == '+') {
16            // 遇到操作符, 弹出操作符栈中优先级大于等于当前操作符的操作符
17        }
18    }
19 }

```

```

16         while (!opStack.empty() && priority(opStack.top()) >= priority(it)) {
17             // 弹出优先级大于等于当前操作符的操作符
18             /* 将运算符与操作数取出来运算 */
19             }
20             opStack.push(it); // 将当前操作符压入操作符栈
21         } else { // 遇到操作数, 创建一个新的节点, 将其压入操作数栈
22             NFA[idNFA][it].insert(idNFA + 1);
23             nodeStack.push({idNFA++, idNFA++});
24             charset.insert(it);
25         }
26     } while (!opStack.empty()) { // 将操作符栈中剩余的操作符弹出
27         /* 将运算符与操作数取出来运算 */
28     }
29     // 获取起始状态和终止状态
30     start = nodeStack.top().first;
31     end = nodeStack.top().second;
32 }

```

3.4.2 NFA到DFA

函数 `void NFA2DFA()` 使用子集构造法, 从 **NFA** 的起始状态出发, 构建 **DFA** 状态和状态转移表。具体步骤包括:

- 获取**NFA**中起始状态的闭包作为**DFA**的起始状态。
- 使用广度优先搜索 (**BFS**) 来构建 **DFA** 状态和状态转移表。
- 遍历字符集, 获取状态集合通过字符转换后的状态集合的闭包。
- 构建 **DFA** 终止状态集合。

具体实现如下:

```

1 // 使用子集构造法, 将NFA转换为DFA
2 void XLEXGenerator::NFA2DFA() {
3     set<int> startSet = getClosure({start}); // 获取起始状态的闭包
4     stateMap[startSet] = idDFA;
5     stateMapReverse[idDFA++] = startSet;
6     queue<set<int>> q; // BFS队列
7     q.push(startSet);
8     while (!q.empty()) { // BFS
9         set<int> u = q.front(); q.pop();
10        for (char c : charset) {
11            auto v = getClosure(move(u, c)); // 获取状态集合u通过字符c转换后的状态集
12            合的闭包
13            if (v.empty()) continue; // 空集不处理
14            if (!stateMap.count(v)) {
15                stateMap[v] = idDFA; // 添加新状态
16                stateMapReverse[idDFA++] = v;
17                q.push(v);
18            }
19            DFA[stateMap[u]][c] = stateMap[v]; // 添加状态转换
20        }
21    }
22 }

```

```

21     for (auto& [k, v]:stateMapReverse) {    // 获取终止状态集合
22         if (v.count(end)) {
23             DFAendSet.insert(k);
24         }
25     }
26 }

```

3.4.3 DFA到最小化DFA

函数 `void minimizeDFA()` 使用逐一合并状态的方法，实现DFA到最小化DFA的转换，通过并查集识别等价状态并合并它们。具体步骤包括：

- 使用并查集识别等价状态。
- 合并等价状态，构建最小化 DFA 状态和状态转移表。
- 获取最小化 DFA 终止状态集合。

使用二重循环遍历DFA中的每一对状态，如果两个状态通过字符转移后的状态在同一个集合中，表明这两个状态是等价的，可以进行合并。使用并查集来维护等价关系，如果DFA中两个状态等价，则将两个状态在并查集中merge起来，最后将并查集中的每个集合作为最小化DFA的一个状态。

并查集的实现如下：

```

1  struct DSU {
2      vector<int> fa, siz;
3      DSU(int n) : fa(n), siz(n, 1) { iota(fa.begin(), fa.end(), 0); }
4      int find(int x) {
5          return fa[x] == x? x:(fa[x] = find(fa[x]));
6      }
7      bool same(int x, int y) { return find(x) == find(y); }
8      bool merge(int x, int y) {
9          x = find(x); y = find(y);
10         if (x == y) return false;
11         if (x > y) swap(x, y);
12         siz[x] += siz[y];
13         fa[y] = x;
14         return true;
15     }
16     int size(int x) { return siz[find(x)]; }
17 };

```

具体实现如下：

```

1  // 使用合并的方法，将DFA最小化
2  void XLEXGenerator::minimizeDFA() {
3      DSU dsu(idDFA);
4      map<int, int> DFA2minDFA;
5      for (int i=1;i<idDFA;i++) {
6          if (i != dsu.find(i))    // 如果已经合并则不处理
7              continue;
8          for (int j=1;j<idDFA;j++) {
9              if (dsu.same(i, j)) continue;
10             // 终态只能与终态合并，非终态只能与非终态合并
11             if (DFAendSet.count(i) ^ DFAendSet.count(j))

```



```

12         continue;
13     bool flag = true;
14     for (char c : charset) {
15         if (!DFA[i].count(c) && !DFA[j].count(c)) // 如果两个状态都不存在转
移则不处理
16             continue;
17         if (DFA[i].count(c) ^ DFA[j].count(c)) { // 如果两个状态只有一个存在转
移则不合并
18             flag = false;
19             break;
20         }
21         if (!dsu.same(DFA[i][c], DFA[j][c])) { // 如果两个状态转移后的状态不在同
一个集合则不合并
22             flag = false;
23             break;
24         }
25     }
26     if (flag) { // 合并两个状态
27         dsu.merge(i, j);
28     }
29 }
30 }
31 auto move = [&](set<int> states, char c) {
32     set<int> res;
33     for (auto state : states) {
34         if (!DFA.count(state)) continue;
35         if (!DFA[state].count(c)) continue;
36         res.insert(DFA[state][c]);
37     }
38     return res;
39 };
40 map<int, int> DS2minS;
41 // 划分出来的集合
42 map<int, set<int>> partitionSet;
43 for (int i=1; i<idDFA; i++) {
44     if (i == dsu.find(i)) { // 根节点
45         partitionSet[idMinDFA] = {i};
46         DS2minS[i] = idMinDFA++;
47     } else {
48         partitionSet[DS2minS[dsu.find(i)]].insert(i);
49     }
50 }
51 for (auto& [node, states]:partitionSet) {
52     for (char c : charset) {
53         set<int> next = move(states, c);
54         if (next.empty()) continue;
55         int nextNode = DS2minS[dsu.find(*next.begin())];
56         minDFA[node][c] = nextNode;
57     }
58     // 记录终态
59     for (auto state : states) {
60         if (DFAendSet.count(state)) {
61             minDFAendSet.insert(node);
62             break;
63         }

```

```

64     }
65 }
66 }

```

3.4.4 DFA到词法分析程序

函数 `void DFA2XLEX()` 实现minDFA到词法分析程序的转换，具体步骤包括：

- 使用状态机设计，根据当前状态和输入字符进行状态转移。
- 生成 C++ 代码，根据最小化 DFA 设计状态转移。
- 输出匹配成功或失败。

使用一个变量保持当前的状态，并将转换写成一个双层嵌套的case语句，其中第1个case语句测试当前的状态，嵌套着的第2层测试输入字符及所给状态。

具体实现如下：

```

1  // 根据最小化DFA生成c++词法分析程序
2  void XLEXGenerator::DFA2XLEX() {
3      stringstream ss;
4      ss << "#include<iostream>\n";
5      ss << "#include<string>\n";
6      ss << "using namespace std;\n";
7      ss << "int main() {\n";
8      ss << "    string s;\n";
9      ss << "    cin>>s;\n";
10     ss << "    int state = 1;\n";
11     ss << "    for (char c : s) {\n";
12     ss << "        switch (state) {\n";
13     for (auto& [node, ma]:minDFA) {
14         ss << "            case " << node << ":\n";
15         ss << "                switch (c) {\n";
16         for (auto& [k, v]:ma) {
17             ss << "                    case '" << k << "': state = " << v << ";
18             break;\n";
19         }
20         ss << "                default: state = 0;\n";
21         ss << "            }\n";
22         ss << "        break;\n";
23     }
24     ss << "        default: state = 0;\n";
25     ss << "    }\n";
26     ss << "    switch (state) {\n";
27     for (auto it : minDFAendSet) {
28         ss << "        case " << it << ": cout<<\\"匹配成功\ "<<endl; break;\n";
29     }
30     ss << "        default: cout<<\\"匹配失败\ "<<endl;\n";
31     ss << "    }\n";
32     ss << "}\n";
33     xlex = ss.str();
34 }

```

3.5 测试

选择正则表达式 `a(a|b)*` 作为测试用例测试结果如下：

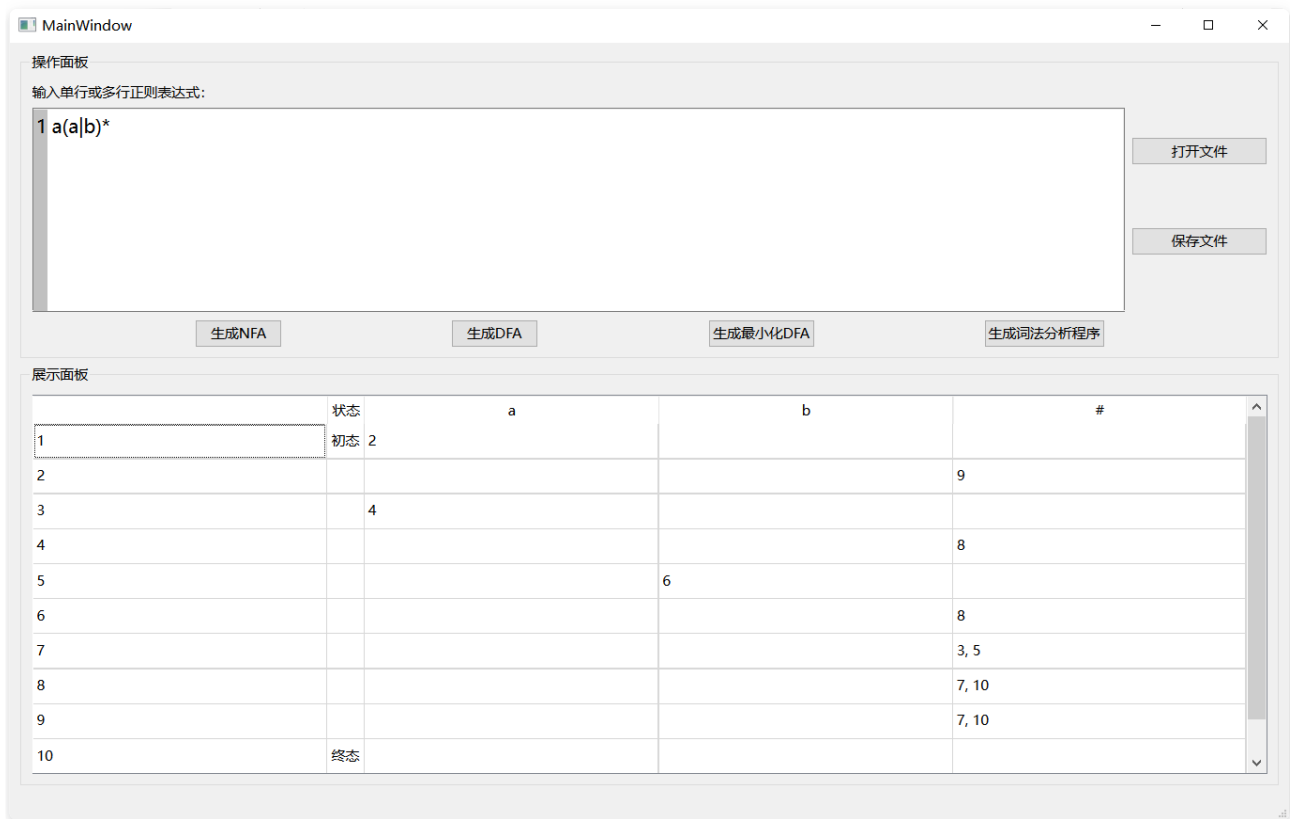


图1. 正则表达式转NFA

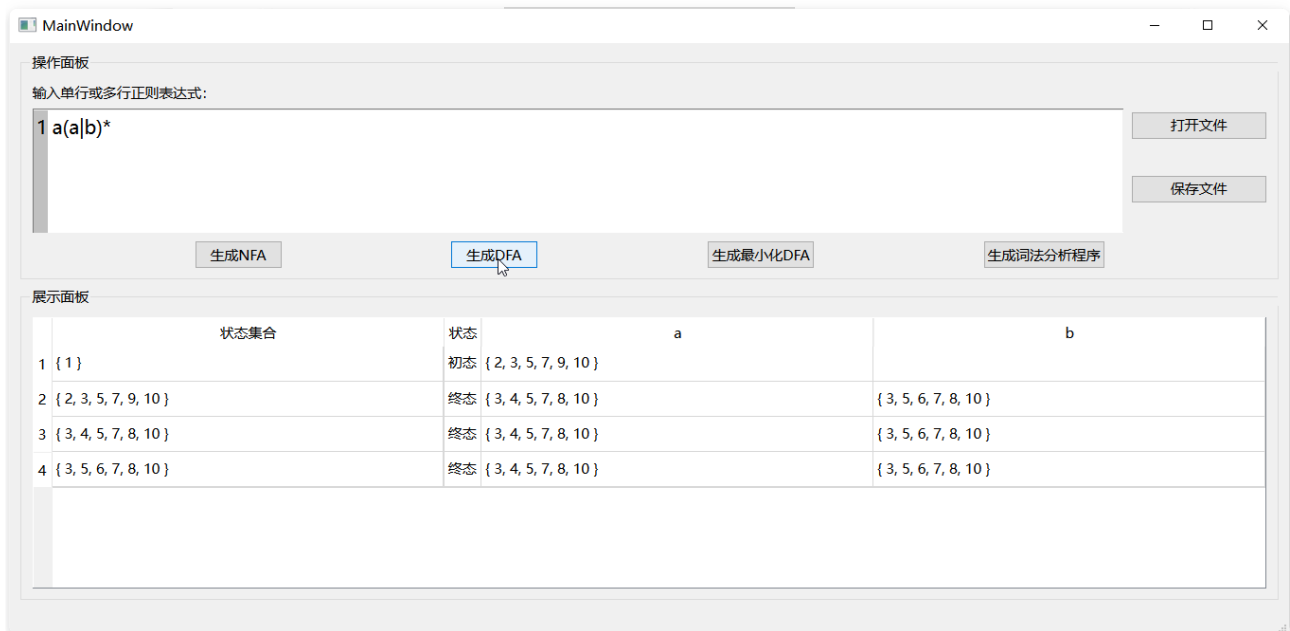


图2. NFA转DFA

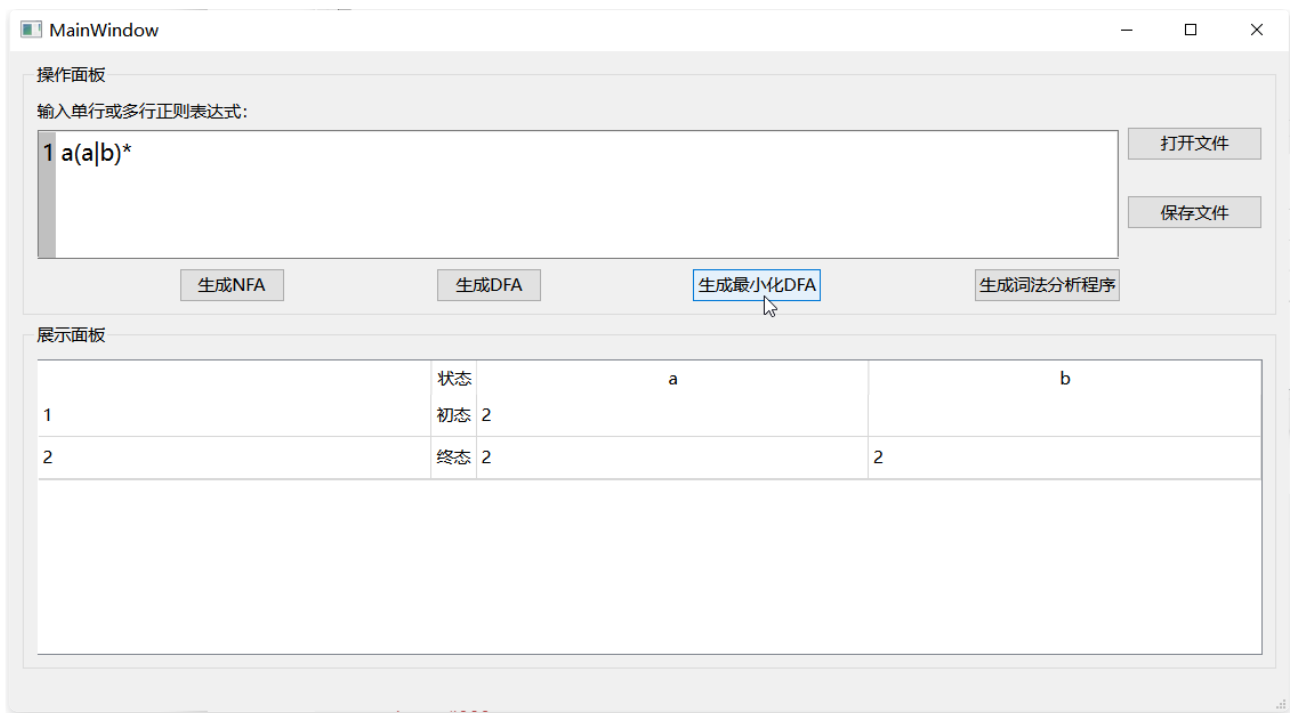


图3. 最小化DFA

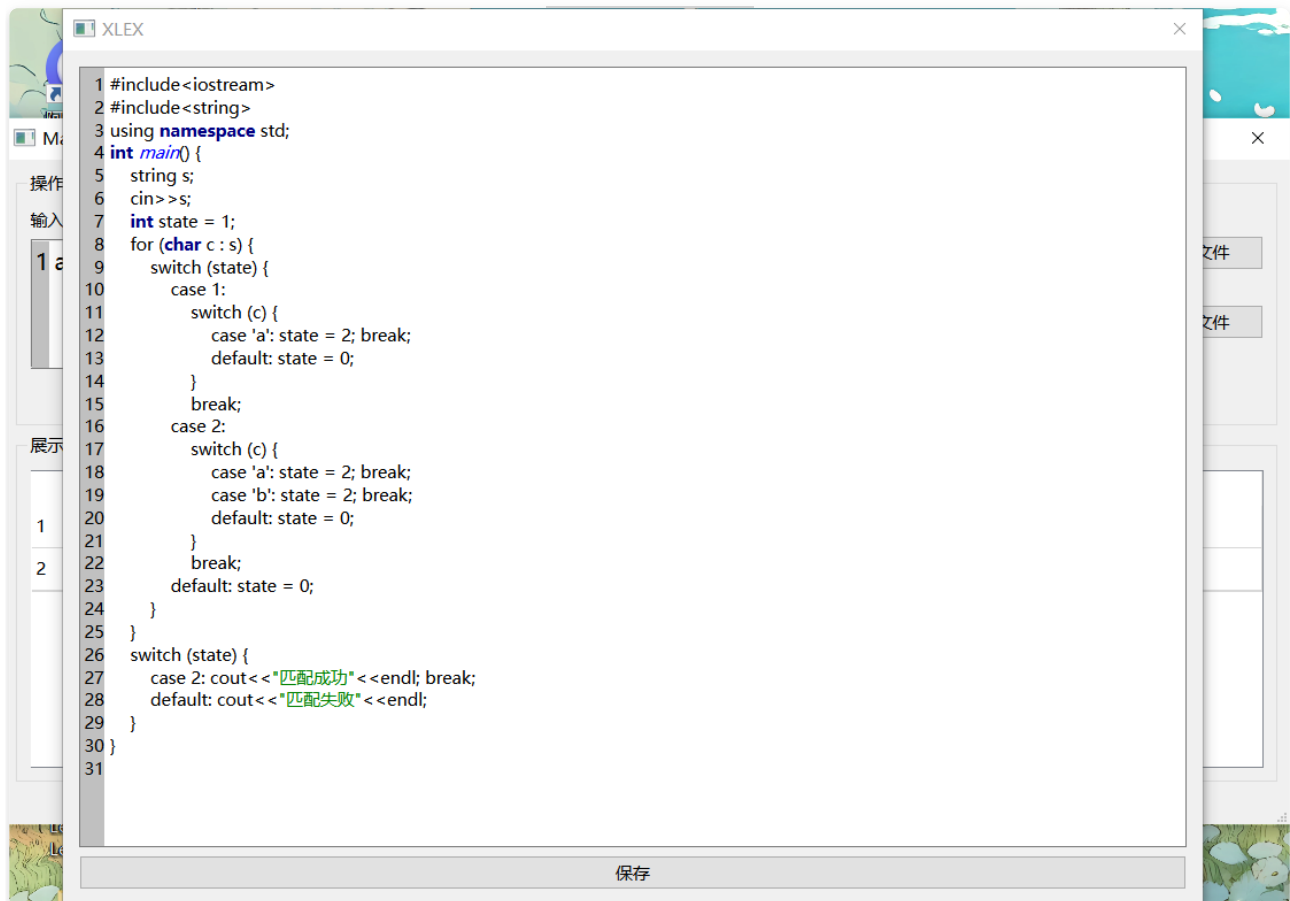


图4. 生成词法分析程序

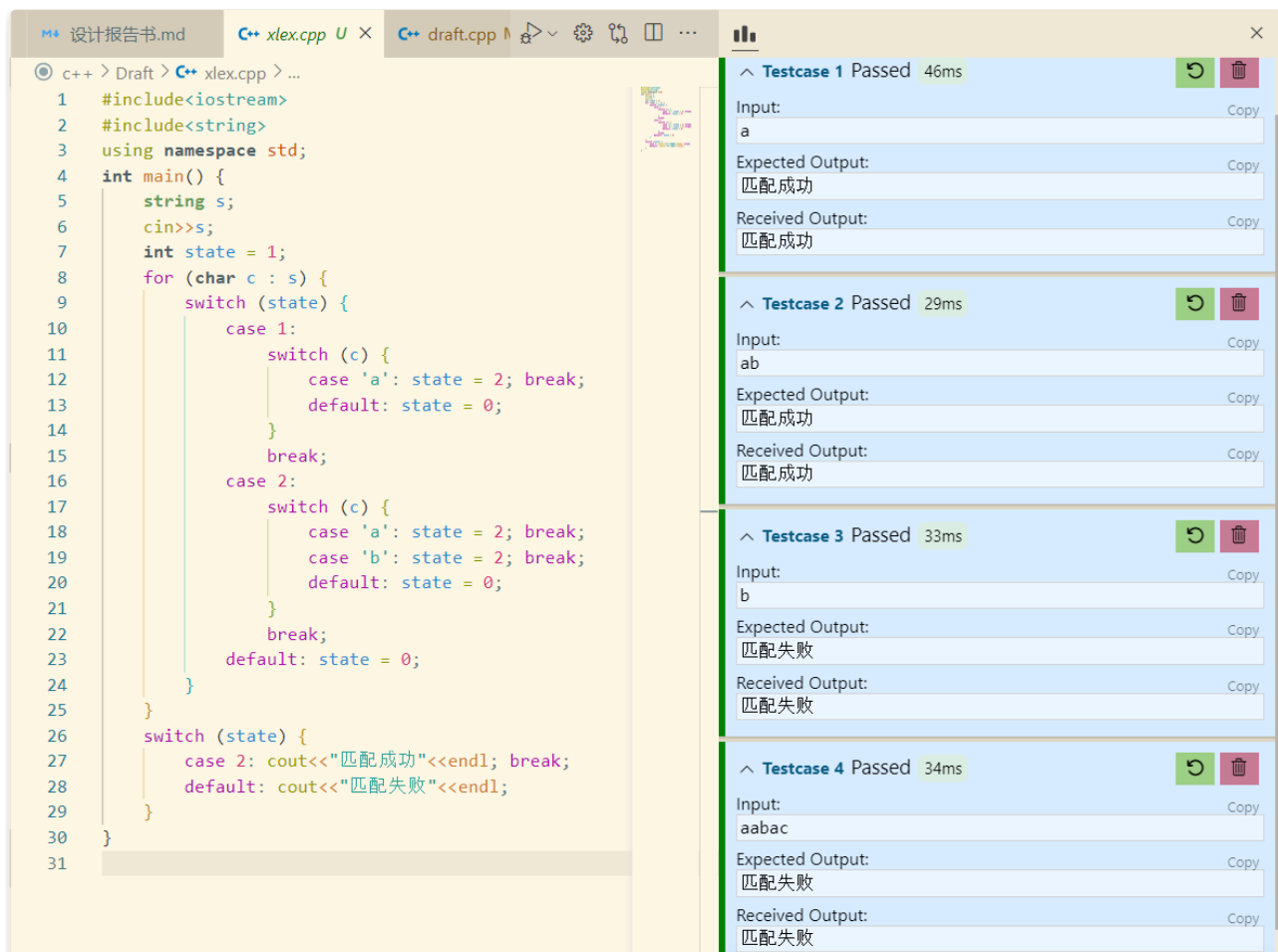


图5. 测试生成的词法分析程序

通过该测试，可以看出该程序能完成将正则表达式-->NFA-->DFA-->DFA最小化-->词法分析程序的全过程。

更多测试样例请看Testfile文件夹内的测试文档。

4 实验总结

这次实验让我深入理解了从正则表达式到词法分析程序的转换过程。首先，通过Thompson构造法，使用双栈操作，将正则表达式转换为NFA。这个步骤中，最关键的是需要对每个正则表达式里的每个操作符所对应的操作构建NFA中的状态和状态之间的转移关系。接着，使用子集构造法将NFA转换为DFA，去除NFA中的空转移与一对多映射。

同时我学习了如何对DFA进行最小化处理，通过状态的等价性来合并状态，减少了DFA中的状态数量，能使词法分析程序更为简洁和高效。最后，需要将最小化DFA转换为词法分析程序，通过状态机设计，根据当前状态和输入字符进行状态转移，生成C++代码，输出匹配成功或失败。

在项目中，我还强化了使用Qt框架来开发GUI应用程序的技能，能更轻松地开发基于qt的GUI程序。

5 参考文献

[《Qt 学习之路 2》目录 - DevBean Tech World](#)

[【精选】NFA到DFA的转化](#)

黄煜廉老师的ppt