

CS205 C/ C++ Programming - Lab Assignment 4

Name: 唐千栋(Qiandong Tang)

SID: 11612730

Part 1 - Analysis

The problem is to display on the standard output the name of the block to which most characters belong. First, we define a struct `Block` to store the language of each Unicode block. Second, we write a function `search` to binary search which block the unicode value belongs to. Then for each unicode value we read from the text, we find the corresponding block of it and count it. Finally, we would print the max count of block.

Part 2 - Code

```
1  #include <algorithm>
2  #include <cstdio>
3  #include <cstring>
4  #include <fstream>
5  #include <iostream>
6  #include <sstream>
7  #include <string>
8  #include <vector>
9  #include "utf8.h"
10 #define pb push_back
11 using namespace std;
12
13 const int maxn = 3e3 + 10;
14
15 string trim(string s) {
16     s.erase(s.find_last_not_of(" \t\r\n") + 1);
17     s.erase(0, s.find_first_not_of(" \t\r\n"));
18     return s;
19 }
20
21 vector<string> Split(const string& s, const string& split_str) {
22     string::size_type pos1, pos2;
23     pos1 = 0;
24     pos2 = s.find(split_str);
25     vector<string> v;
26     while (pos2 != string::npos) {
27         v.pb(s.substr(pos1, pos2 - pos1));
28         pos1 = pos2 + split_str.size();
29         pos2 = s.find(split_str, pos1);
30     }
```

```

31     if (pos1 != s.size()) v.pb(s.substr(pos1));
32     return v;
33 }
34
35 int hex_string_to_int(const string& s) {
36     int res = stoi(s.c_str(), NULL, 16);
37     return res;
38 }
39
40 struct Block {
41     int start, end;
42     string lang;
43 };
44
45 ostream& operator<<(ostream& out, const Block& block) {
46     out << "Language: " << block.lang << ", Range: " << block.start << "-"
47         << block.end;
48     return out;
49 }
50
51 int sz;
52 Block blocks[maxn];
53
54 void read_utf(string path) {
55     ifstream ifile(path);
56     string line;
57     sz = 0;
58     while (getline(ifile, line)) {
59         if (line[0] == '#' || line.empty()) continue;
60         istringstream sin(line);
61         string field;
62         vector<string> v;
63         while (getline(sin, field, ';')) {
64             field = trim(field);
65             v.pb(field);
66         }
67         vector<string> vs = Split(v[0], "..");
68         blocks[sz].start = hex_string_to_int(vs[0]);
69         blocks[sz].end = hex_string_to_int(vs[1]);
70         blocks[sz].lang = v[1];
71         sz++;
72     }
73 }
74
75 int search(unsigned int codepoint) {
76     int l = 0;
77     int r = sz - 1;
78     int mid;
79     if (codepoint > blocks[r].end) return -1;

```

```

80     while (l <= r) {
81         mid = (l + r) >> 1;
82         if (blocks[mid].start <= codepoint && codepoint <= blocks[mid].end) {
83             break;
84         }
85         if (codepoint > blocks[mid].end) {
86             l = mid + 1;
87         } else {
88             r = mid - 1;
89         }
90     }
91     return mid;
92 }
93
94 int cnts[maxn];
95
96 int main() {
97     read_utf("Blocks.txt");
98     memset(cnts, 0, sizeof(int));
99     string s;
100     unsigned char* p;
101     int bytes_in_char;
102     unsigned int codepoint;
103     int x;
104
105     while (!cin.eof()) {
106         cin >> s;
107         p = (unsigned char*)s.c_str();
108         while (*p) {
109             codepoint = utf8_to_codepoint(p, &bytes_in_char);
110             if (codepoint) {
111                 // printf("%c %u (%X) %d byte character\n", *p, codepoint,
112                 //     codepoint, bytes_in_char);
113                 // p += bytes_in_char; // Same as the line that follows
114                 _utf8_incr(p);
115             } else {
116                 printf("%c Invalid UTF-8\n", *p);
117                 p++;
118             }
119             x = search(codepoint);
120             if (x >= 0) cnts[x]++;
121         }
122     }
123     cout << blocks[max_element(cnts, cnts + sz) - cnts].lang << endl;
124     return 0;
125 }
126

```

Part 3 - Result & Verification

Test case:

```
→ Assignment4 git:(master) X ./main < ./TEST\ DATA\ FOR\ LAB\ 4/sample.txt
Armenian
→ Assignment4 git:(master) X ./main < ./TEST\ DATA\ FOR\ LAB\ 4/sample2.txt
Georgian
→ Assignment4 git:(master) X ./main < ./TEST\ DATA\ FOR\ LAB\ 4/sample3.txt
Lao
→ Assignment4 git:(master) X ./main < ./TEST\ DATA\ FOR\ LAB\ 4/sample4.txt
Malayalam
→ Assignment4 git:(master) X ./main < ./TEST\ DATA\ FOR\ LAB\ 4/sample5.txt
Devanagari
→ Assignment4 git:(master) X ./main < ./TEST\ DATA\ FOR\ LAB\ 4/sample6.txt
Georgian
```

Part 4 - Difficulties & Solutions

1. How to convert character to unicode value

using `codepoint = utf8_to_codepoint(p, &bytes_in_char);`

2. How to split string

```
1  vector<string> Split(const string& s, const string& split_str) {
2      string::size_type pos1, pos2;
3      pos1 = 0;
4      pos2 = s.find(split_str);
5      vector<string> v;
6      while (pos2 != string::npos) {
7          v.pb(s.substr(pos1, pos2 - pos1));
8          pos1 = pos2 + split_str.size();
9          pos2 = s.find(split_str, pos1);
10     }
11     if (pos1 != s.size()) v.pb(s.substr(pos1));
12     return v;
13 }
```