

# 项目文件结构说明

```
1 bin
2 |---xqz.exe          压缩工具
3 include              头文件
4 |---CodeTree.hpp     用于解压的编码树
5 |---cxxopts.hpp      github上找的命令选项解析器
6 |---fnode.hpp        用于压缩文件内文件结构的类
7 |---HuffmanTree.hpp  关于哈夫曼树和哈夫曼编码
8 |---xq_debug.hpp     debug用
9 |---xqzio.hpp        关于输入输出的一些函数
10 src                 源文件
11 |---CodeTree.cpp
12 |---fnode.cpp
13 |---HuffmanTree.cpp
14 |---xqzio.cpp
15 |---main.cpp
16 test
17 |---performance.xlsx 性能测试结果
18 .gitignore
19 Makefile
20 Project.pdf
21 README.md           此开发文档
```

## 项目设计、实现的具体思路

### 1. 建哈夫曼树&生成哈夫曼编码

- 建哈夫曼树
  - 所有节点放在一个优先队列里
  - 每次取出头部的两个，合并，放回队列，直到队列元素个数为1
- 生成哈夫曼编码
  - 从叶节点往上走，生成编码

```
1 #define uc unsigned char
2 typedef std::map<uc, std::string> Word2Code_T;
3 typedef std::map<std::string, uc> Code2Word_T;
4 extern int nodes_num; //节点个数
5 extern Word2Code_T word2code;
6 extern Code2Word_T code2word;
7
8 class HuffmanTreeNode
9 {
10 private:
11     int value = 0;
12     HuffmanTreeNode *father = nullptr;
13     char lr = 0; // 0-left_son, 1-right_son
14 public:
15     int idx = 0;
```

```

16 HuffmanTreeNode();
17 HuffmanTreeNode(int v);
18 HuffmanTreeNode(int v, HuffmanTreeNode *real_addr);
19 bool operator<(const HuffmanTreeNode &x) const;
20 void SetValue(int v);
21 void SetFather(HuffmanTreeNode *f);
22 void SetLR(char _lr);
23 int GetValue();
24 char GetLR();
25 // GenerateCode - 从此节点向上生成01序列, 仅用于叶节点
26 std::string GenerateCode();
27 };
28
29 extern std::map<int, HuffmanTreeNode *> idx2nodes;
30 extern HuffmanTreeNode *word2nodes_addr[256];
31
32 struct HuffmanForest_Queue_Cmp{bool operator()(int a, int b)};
33 typedef std::priority_queue<int, std::vector<int>, HuffmanForest_Queue_Cmp>
HuffmanForest_Queue;
34
35 class HuffmanForest
36 {
37 private:
38     HuffmanForest_Queue nodes;
39
40 public:
41     HuffmanForest();
42     // is_Tree - 是否是一棵树, 即是否生成完毕
43     bool is_Tree();
44     // AddNode - 用于添加初始节点(叶节点)
45     void AddNode(HuffmanTreeNode *x, uc word);
46     void Pop();
47     HuffmanTreeNode *Top();
48     void Push(HuffmanTreeNode *x);
49     word2Code_T GenerateWord2Code();
50     Code2Word_T GenerateCode2Word();
51 };
52
53 // MergeHuffmanTree - 合并两个节点
54 HuffmanTreeNode *MergeHuffmanTree(HuffmanTreeNode *x, HuffmanTreeNode *y);
55 // GenerateHuffmanTree - 生成哈夫曼树
56 void GenerateHuffmanTree(HuffmanForest *x);

```

## 2. 压缩

### 先判断是不是文件

- 文件
  - 文件是否为空
    - 是
      - 按空文件格式写入(格式见下)
    - 否
      - 读入并计数
      - 生成哈夫曼编码

- 评估压缩后是否有压缩效果
    - 若有, 按字节操作, 转换成哈夫曼编码写入压缩文件
    - 若无, 不压缩, 直接copy
- 文件夹
  - 对文件夹内所有文件进行压缩

## 压缩文件格式

```

1  [1 or 0]          # 1 - file; 0 - folder
2  ##### if file #####
3  .\[file_name]
4  [压缩文件长度(bit)]
5  # code2word
6  [code]:[word]
7  ...
8  *
9  [compressed file]
10 ##### if folder #####
11 .\[folder_name]
12 [number of empty folder]
13 \...              # 空文件夹相对根文件夹路径
14 ...
15 [number of files]
16 \[file_path]      # 相对根文件夹路径
17 [压缩文件长度]    -
18 # code2word        |
19 [code]:[word]      |
20 ...                |
21 *                  |
22 [compressed file]  -
23 ...
24 ##### fake compress #####
25 .\[file_name]
26 [文件长度(byte)]
27 **
28 ...
29 ##### empty file #####
30 .\[file_name]
31 0

```

```

1  extern ull word_freq[256];
2  extern const int MAX_IO_N;
3  extern std::vector<std::string> files;
4  extern std::vector<std::string> empty_folders;
5
6  bool isFile(const char* path);
7  // GetFileSize - 得到文件大小(按byte)
8  ull GetFileSize(const char* file_path);
9  // GetFoldersize - 得到文件夹大小(按byte)
10 ull GetFoldersize(const char* folder_path);
11 // GetRootPath - 得到根目录
12 std::string GetRootPath(std::string src);
13 // GetFolderFiles - 得到root下所有文件, 存在files里
14 void GetFolderFiles(std::string root);
15 // xqz_read_src_compress - 读文件, 并计数得到word_freq
16 void xqz_read_src_compress(const char *filename);

```

```
17 // cnt_freq - 对x中字符计数
18 void cnt_freq(uc *x, int l);
19 // compressed_length - 压缩文件长度(按bit)
20 ull compressed_length();
21 // compressed_code2word_length - 压缩文件加上code2word长度(按byte)
22 ull compressed_code2word_length();
```

### 3. 解压

先判断是文件夹还是文件

- 文件
  - 读入文件路径，并创建解压文件所在的文件夹（如果路径不存在）
  - 按照压缩文件格式依次读入各个部分
    - 若为fake compress，直接无需解压，copy到解压文件
    - 否则
      - 建立编码树，就是没有节点频率信息的哈夫曼树
      - 按位操作，0走左儿子，1走右儿子，直到叶节点，于是得到解压出来的word，写入解压文件
- 文件夹
  - 处理空文件夹
  - 处理文件

### 4. 展示压缩文件结构

- 按照压缩文件格式读入，抛弃除文件路径外的信息
- 按照文件路径创建fnode
- 用printDirTree()输出

## 性能测试结果

见 `./test/performance.xlslx`

思考：

- 哈夫曼压缩对字符出现频率差异较大的文件效果较好，对字符出现频率平均的文件效果很差
- 由于哈夫曼压缩是不定长编码，解压时需要按位操作，解压速度不好提升

## 遇到的问题和解决方案（项目重难点）

### 1. mkdir

默认的mkdir函数不能创建多级目录，用递归解决

### 2. 检查文件或文件夹是否存在

access()

### 3. 压缩和解压时读写问题

一次性读入0x1000个字符，然后处理

### 4. 用户输入的路径不规范

```
1 void parse_path(string &x)
2 {
3     for (int i = 0; i < x.size(); ++i) //所有的 '/' 换成 '\\'
4     {
5         if (x[i] == '/')
6         {
7             x[i] = '\\';
8         }
9     }
10    if (x[0] != '.' || x[0] != '\\') //若路径开头没有'.\\', 默认为当前目录下
11    {
12        x = ".\\" + x;
13    }
14    if (x[x.size() - 1] == '\\') //路径末尾若有 '\\', 去掉
15        x[x.size() - 1] = '\\0';
16    return;
17 }
```

## 优化技巧（项目亮点）

- 算法
  - 哈夫曼树
    - 利用优先队列生成哈夫曼树
  - 压缩
    - 压缩前评估效果，防止压缩后反而比原文件大很多
  - 解压
    - 解压时生成简单的编码树，节省生成哈夫曼树的时间
    - 解压时用位运算获取每一位的信息
- 功能
  - 多个文件压缩成一个文件

这种压缩出来的文件，解压时若不用“-n”选项指定文件名，默认会被一个文件夹 tmp 包含
  - 展示压缩文件内文件结构
  - 压缩或者解压到指定文件夹
- 用户友好
  - 用命令的形式，就可以自动补全路径
  - 进度条显示
  - help显示

```
xqz - huffman compress by xq
Usage:
  xqz [OPTION...]

  -c, --compress <file/folder>  compress something
  -n, --name <name>             set output name
  -x, --extract <file>          extract .xqz file
  -t, --to <directory>         place the output into <directory>
  -s, --show <file>            show files
  -h, --help                    print help
```

## 参考

- 获取文件夹下文件目录

[https://blog.csdn.net/leo\\_888/article/details/80681184](https://blog.csdn.net/leo_888/article/details/80681184)

<https://blog.csdn.net/sagghab/article/details/81436515> - \_finddata\_t

- access()函数：用于文件读取，判断文件是否存在，并判断文件是否可写。Linux中，函数为access();标准C++中，该函数为\_access，两者的使用方法基本相同

access(const char \*pathname, int mode);//位于<unistd.h>中

int \_access(const char \*pathname, int mode);//位于<io.h>中

pathname 为文件路径或目录路径 mode 为访问权限

如果文件具有指定的访问权限，则函数返回0；如果文件不存在或者不能访问指定的权限，则返回-1。

mode的值和含义如下所示：

00——只检查文件是否存在

02——写权限

04——读权限

06——读写权限

- 进度条

<https://www.cnblogs.com/seafever/p/12345200.html>

- 命令选项解析器

<https://github.com/jarro2783/cxxopts>