

哈尔滨工业大学计算学部

## 实验报告

课程名称：数据结构与算法

课程类型：专业核心基础课（必修）

实验项目：树形结构及其应用

实验题目：哈夫曼编码与译码方法

实验日期：2024 年 4 月 21 日

班级：22WL022

学号：2022110829

姓名：杨明达

设计成绩	报告成绩	任课老师
		张岩

## 一、实验目的

哈夫曼编码是一种以哈夫曼树（最优二叉树，带权路径长度最小的二叉树）为基础变长编码方法。其基本思想是：将使用次数多的字符转换成长度较短的编码，而使用次数少的采用较长的编码，并且保持编码的唯一可解性。在计算机信息处理中，经常应用于数据压缩。是一种一致性编码法（又称“熵编码法”），用于数据的无损压缩。本实验要求实现一个完整的哈夫曼编码与译码系统。

## 二、实验要求及实验环境

### 实验要求：

1. 从文件中读入任意一篇英文文本文件，分别统计英文文本文件中各字符（包括标点符号和空格）的使用频率。
2. 根据已统计的字符使用频率构造哈夫曼编码树，并给出每个字符的哈夫曼编码（字符集的哈夫曼编码表）。
3. 将文本文件利用哈夫曼树进行编码，存储成压缩文件（哈夫曼编码文件）；
4. 将哈夫曼编码文件译码为文本文件，并与原文件进行比较。
5. 计算你的哈夫曼编码文件的平均编码长度和压缩率，并与实验结果比较验证。
6. 利用堆结构（优先级队列），优化哈夫曼编码算法。

实验环境：Windows 11 && Visual Studio Code

三、设计思想（本程序中的用到的所有数据类型的定义，主程序的流程图及各程序模块之间的调用关系、核心算法的主要步骤）

### 1. 逻辑设计

#### 1.1 主程序的流程图

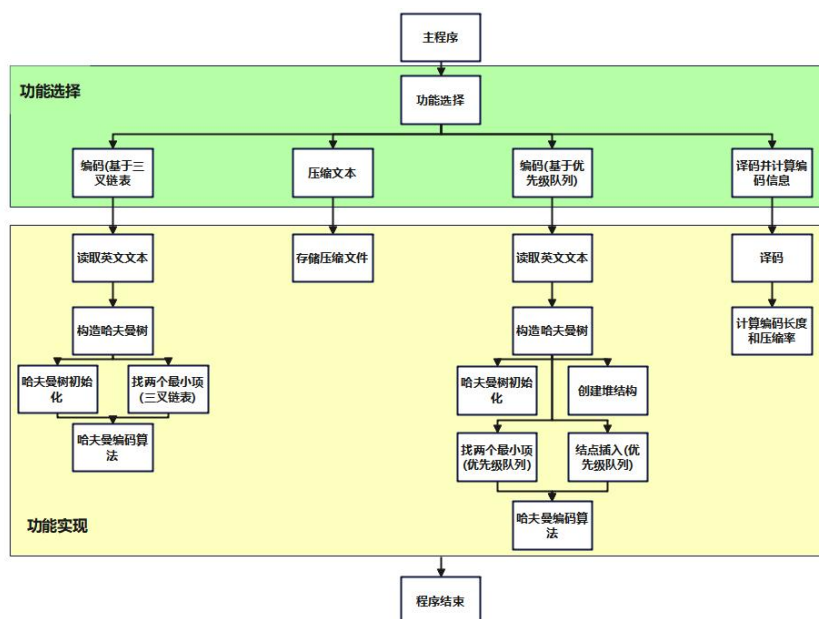


图 1 主程序流程图

## 1.2 各程序模块的调用关系

(1) 从构造哈夫曼树（基于三叉链表）函数中调用哈夫曼树初始化、找到两个最小项（基于三叉链表）模块。

(2) 从构造哈夫曼树（基于优先级队列）函数中调用哈夫曼树初始化、创建堆结构、找到两个最小项（基于优先级队列）、结点插入（优先级队列）模块。

(3) 从创建堆结构函数中调用堆的初始化、结点插入（优先级队列）模块。

(4) 从找到两个最小项（基于优先级队列）函数中调用删除最小项（基于优先级队列）模块。

(5) 从删除最小项（基于优先级队列）模块调用判断堆空模块。

(6) 从结点插入（基于优先级队列）函数调用判断堆满模块。

## 1.3 核心算法

### (1) 哈夫曼树构造（基于三叉链表）

在本次实验中，哈夫曼树的构造是实现哈夫曼编码的关键。在函数中，首先要对一个数初始化，然后进行  $n-1$  次合并，从哈夫曼树中找到两个权重最小的项，标记它的序号，之后建立双亲结点，使其的左右子树分别为这两个最小项，其权重为两个子树权重的和。

```
1. void CreatHT(HuffmanCode H, HuffmanT T) {
2.     int p1, p2, n;
3.     n = Number_Tpye;
4.     InitHT(H, T);
5.     for (int i = n; i < 2 * n - 1; i++) {
6.         SelecMin(T, i, &p1, &p2);
7.         T[p1].parent = T[p2].parent = i;
8.         T[i].lchild = p1;
9.         T[i].rchild = p2;
10.        T[i].weight = T[p1].weight + T[p2].weight;
11.    }
12. }
```

### (2) 哈夫曼树构造（基于优先级队列）

在本次实验中，哈夫曼树的构造是实现哈夫曼编码的关键。在函数中，首先要对一个数初始化，然后创建堆，再进行  $n-1$  次合并，从哈夫曼树中找到两个权重最小的项，标记它的序号，之后建立双亲结点，使其的左右子树分别为这两个最小项，其权重为两个子树权重的和。

```

1. void CreatHT_Heap(HuffmanCode H, HuffmanT T) {
2.     int p1, p2, n;
3.     n = Number_Tpye;
4.     HEAP heap;
5.     InitHT(H, T);
6.     CreatHeap(H, &heap);
7.     for (int i = n; i < 2 * n - 1; i++) {
8.         Elementype new_node;
9.         SelecMin_Heap(&heap, &p1, &p2);
10.        T[p1].parent = T[p2].parent = i;
11.        T[i].lchild = p1;
12.        T[i].rchild = p2;
13.        T[i].weight = T[p1].weight + T[p2].weight;
14.        new_node.key = i;
15.        new_node.weight = T[i].weight;
16.        Heap_Insert(&heap, new_node);
17.    }
18. }

```

### (3) 哈夫曼编码算法的实现

在本次实验中，哈夫曼编码算法是将字符转换成 01 序列的关键所在，在函数中，c 和 p 分别指示 T 中的孩子和双亲的位置，cd 临时存放编码，start 指示编码在 cd 中的位置，第五行用来定义编码结束符。进入循环，依次求叶子 T[i] 的编码：读入叶子 T[i] 对应的字符，定义编码起始位置的初值，从叶子 T[i] 开始上溯，进入循环，判断，若 T[c] 是 T[p] 的左孩子，则生成代码 0，否则生成代码 1，继续上溯直到上溯到 T[c] 是树根位置。最后复制编码位串与编码表 H。

```

1. void CharSetHuffmanEncoding(HuffmanCode H, HuffmanT T) {
2.     int c, p, i;
3.     char cd[Number_Tpye];
4.     int start;
5.     cd[Number_Tpye] = '\0';
6.     for (i = 0; i < Number_Tpye; i++) {
7.         start = Number_Tpye;
8.         c = i;
9.         while ((p = T[c].parent) >= 0) {
10.            cd[--start] = (T[p].lchild == c) ? '0' : '1';
11.            c = p;
12.        }
13.        strcpy(H[i].bits, &cd[start]);
14.    }
15. }

```

## 2. 物理设计（即存储结构设计）

本实验主要采用三叉链表存储字符，该存储结构包含四个数据，分别为权值（double 类型），左孩子链（int 类型），右孩子链（int 类型），双亲链（int 类型）。三叉链表的示意图如下图所示。



图 3 三叉链表存储

本实验存储字符结点，该存储结构包含三个数据，分别为字符（char 类型），字符编码（char 类型），字符频率（double 类型）。存储结构示意图如下图所示。



图 4 字符结点存储

本实验另一种方法主要利用优先级队列，即堆结构存储字符，该存储结构包含两个主数据，分别为节点个数（int 类型），节点数据（Elementype 类型），其中 Elementype 类型由三个两个部分组成，分别为序号（int 类型），权重（double 类型）。堆结构的示意图如下图所示。

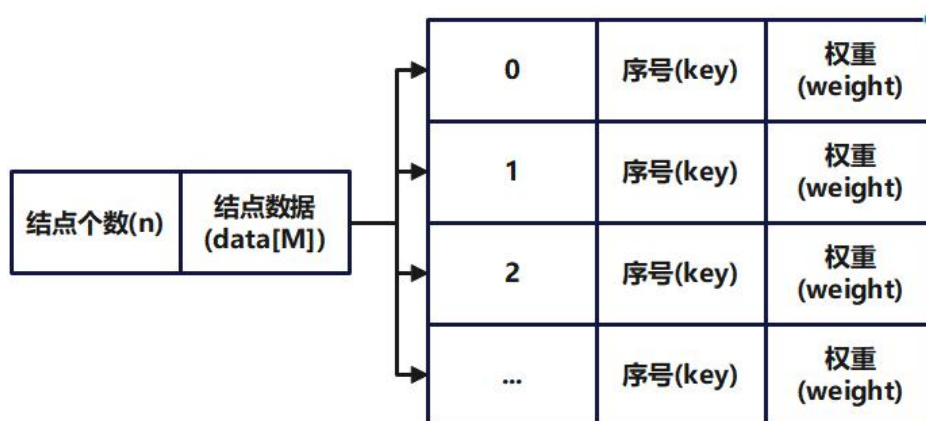


图 5 优先级队列存储

基于上述存储结构，实现哈夫曼编码的读取英文文本、构造哈夫曼树、存储压缩文本、解压等操作。

## 四、测试结果（包括测试数据、结果数据及结果的简单分析和结论，可以用截图得形式贴入此报告）

### 1. 基于三叉链表编码

```

请选择功能实现:1.编码(基于三叉链表);2.编码(基于优先级队列);3.压缩文本;4.译码并计算编码信息;0.退出
1
这个文本有5155个字符,有62种字符文本中各字符的哈夫曼编码及其频率:
字符
  编码:0100100  频率:0.006596
字符"  编码:110  频率:0.151503
字符'  编码:1010101101  频率:0.001164
字符,  编码:001001110  频率:0.001552
字符-  编码:1111111  频率:0.010669
字符.  编码:01001101001  频率:0.000582
字符0  编码:11110100  频率:0.004850
字符1  编码:101010101  频率:0.002134
字符2  编码:111111001  频率:0.002522
字符3  编码:0100110101  频率:0.000970
字符4  编码:001001001  频率:0.001358
字符5  编码:01001100000  频率:0.000388
字符6  编码:0010011110  频率:0.000776
字符7  编码:11111100000  频率:0.000582
字符8  编码:01001100001  频率:0.000388
字符9  编码:01001100010  频率:0.000388
字符A  编码:11111100001  频率:0.000582
字符B  编码:1010101100  频率:0.000970
字符C  编码:11111100010  频率:0.000582
字符D  编码:11110101  频率:0.004850
字符E  编码:01001100110  频率:0.000194
字符F  编码:11111100011  频率:0.000582
字符H  编码:001001100  频率:0.001358
字符I  编码:0010010000  频率:0.000582
字符J  编码:01001100011  频率:0.000388
字符L  编码:01001100100  频率:0.000388
字符M  编码:1010101110  频率:0.001164
字符N  编码:0010010001  频率:0.000582
字符O  编码:01001100101  频率:0.000388
字符P  编码:01001100111  频率:0.000194
字符R  编码:010011010000  频率:0.000194
字符S  编码:0100101  频率:0.006984
字符T  编码:00100101  频率:0.002716
字符U  编码:11111101  频率:0.005238
字符W  编码:010011011  频率:0.001940
字符Y  编码:010011010001  频率:0.000194
字符Z  编码:01001100110  频率:0.000388
字符a  编码:0111  频率:0.060912
字符b  编码:001000  频率:0.010863
字符c  编码:00101  频率:0.025218
字符d  编码:10010  频率:0.031814
字符e  编码:000  频率:0.090786

```

图 6 基于三叉链表编码 a

```

字符f  编码:101011  频率:0.017265
字符g  编码:100110  频率:0.015907
字符h  编码:10100  频率:0.033754
字符i  编码:1011  频率:0.074491
字符j  编码:1010101111  频率:0.001164
字符k  编码:001001101  频率:0.001358
字符l  编码:01000  频率:0.027158
字符m  编码:100111  频率:0.016877
字符n  编码:1000  频率:0.064210
字符o  编码:0101  频率:0.059360
字符p  编码:111100  频率:0.019011
字符r  编码:0011  频率:0.047915
字符s  编码:0110  频率:0.060718
字符t  编码:1110  频率:0.075267
字符u  编码:111110  频率:0.020563
字符v  编码:1010100  频率:0.008341
字符w  编码:0100111  频率:0.007565
字符x  编码:101010100  频率:0.001940
字符y  编码:1111011  频率:0.009893
字符z  编码:0010011111  频率:0.000776
请选择功能实现:1.编码(基于三叉链表);2.编码(基于优先级队列);3.压缩文本;4.译码并计算编码信息;0.退出
3
请选择功能实现:1.编码(基于三叉链表);2.编码(基于优先级队列);3.压缩文本;4.译码并计算编码信息;0.退出
4
这篇文章有 5155 个字符,共 62 种字符种类数。
哈夫曼树编码长度为:4.462658
编码文件的理论压缩率为:0.442168
编码文件的真实压缩率为:0.592214
请选择功能实现:1.编码(基于三叉链表);2.编码(基于优先级队列);3.压缩文本;4.译码并计算编码信息;0.退出
0
PS E:\C_Code\experiment_2>

```

图 7 基于三叉链表编码 b



## 2. 基于优先级队列编码

```

请选择功能实现:1.编码(基于二叉链表);2.编码(基于优先级队列);3.压缩文本;4.译码并计算编码信息;0.退出
2
这个文本有5155个字符,有62种字符文本中各字符的哈夫曼编码及其频率:
字符
  编码:0100100    频率:0.006596
字符"  编码:110    频率:0.151503
字符"  编码:100111111    频率:0.001164
字符'  编码:010011000    频率:0.001552
字符,  编码:1111111    频率:0.010669
字符-  编码:10011111100    频率:0.000582
字符.  编码:11111100    频率:0.004850
字符0  编码:100111101    频率:0.002134
字符1  编码:111101001    频率:0.002522
字符2  编码:0100110101    频率:0.000970
字符3  编码:001001001    频率:0.001358
字符4  编码:00100111010    频率:0.000388
字符5  编码:0010011111    频率:0.000776
字符6  编码:11110100011    频率:0.000582
字符7  编码:00100111000    频率:0.000388
字符8  编码:01001100111    频率:0.000388
字符9  编码:01001101001    频率:0.000582
字符A  编码:1001111100    频率:0.000970
字符B  编码:0010010000    频率:0.000582
字符C  编码:11110101    频率:0.004850
字符D  编码:010011010000    频率:0.000194
字符E  编码:0010010001    频率:0.000582
字符H  编码:001001100    频率:0.001358
字符I  编码:11110100010    频率:0.000582
字符J  编码:01001100100    频率:0.000388
字符L  编码:01001100101    频率:0.000388
字符M  编码:1001111101    频率:0.001164
字符N  编码:10011111101    频率:0.000582
字符O  编码:00100111001    频率:0.000388
字符P  编码:001001110110    频率:0.000194
字符R  编码:001001110111    频率:0.000194
字符S  编码:0100101    频率:0.006984
字符T  编码:00100101    频率:0.002716
字符U  编码:11111101    频率:0.005238
字符W  编码:010011011    频率:0.001940
字符Y  编码:010011010001    频率:0.000194
字符Z  编码:01001100110    频率:0.000388
字符a  编码:0111    频率:0.060912
字符b  编码:001000    频率:0.010863
字符c  编码:00101    频率:0.025218
字符d  编码:10010    频率:0.031814
字符e  编码:000    频率:0.090786

```

图 8 基于优先级队列编码 a

```

字符f  编码:101011    频率:0.017265
字符g  编码:100110    频率:0.015907
字符h  编码:10100    频率:0.033754
字符i  编码:1011    频率:0.074491
字符j  编码:1111010000    频率:0.001164
字符k  编码:001001101    频率:0.001358
字符l  编码:01000    频率:0.027158
字符m  编码:101010    频率:0.016877
字符n  编码:1000    频率:0.064210
字符o  编码:0101    频率:0.059360
字符p  编码:111100    频率:0.019011
字符r  编码:0011    频率:0.047915
字符s  编码:0110    频率:0.060718
字符t  编码:1110    频率:0.075267
字符u  编码:111110    频率:0.020563
字符v  编码:1001110    频率:0.008341
字符w  编码:0100111    频率:0.007565
字符x  编码:100111100    频率:0.001940
字符y  编码:1111011    频率:0.009893
字符z  编码:0010011110    频率:0.000776
请选择功能实现:1.编码(基于二叉链表);2.编码(基于优先级队列);3.压缩文本;4.译码并计算编码信息;0.退出
3
请选择功能实现:1.编码(基于二叉链表);2.编码(基于优先级队列);3.压缩文本;4.译码并计算编码信息;0.退出
4
这篇文章有 5155 个字符,共 62 种字符种类数。
哈夫曼树编码长度为:4.462658
编码文件的理论压缩率为:0.442168
编码文件的真实压缩率为:0.592214
请选择功能实现:1.编码(基于二叉链表);2.编码(基于优先级队列);3.压缩文本;4.译码并计算编码信息;0.退出
0

```

图 9 基于优先级队列编码 b

### 3. 输入文本

Trade tension between China and the United States, simmering over Washington's repeated tariff increases and policy restrictions against Chinese exporters and enterprises, has escalated again with Washington's launch of a Section 301 investigation targeting China's maritime, logistics and shipbuilding sectors.

However, the old playbook of unilateralism and protectionism pursued by the White House will not only fail to bring about the reshoring of manufacturing desired by the US, but will also result in the US facing the challenge of a more expensive supply chain, experts said.

China expressed on Wednesday night its strong dissatisfaction with and firm opposition to the US Trade Representative's initiation of the Section 301 investigation, emphasizing that it will keep abreast of the progress of the probe and take all necessary measures to defend its own rights and interests, the Ministry of Commerce said in a statement.

Section 301 of the Trade Act of 1974, as amended, is used to respond to so-called unjustifiable, unreasonable or discriminatory foreign government practices that burden or restrict US commerce.

The investigation petition was jointly filed on March 12 by five US labor unions. "The number of commercial shipyards in the US has plunged by more than 70 percent, tens of thousands of jobs have been lost, and the US now produces only a fraction of 1 percent of the world's commercial vessels, falling to 19th place," they said in the petition.

Multiple US research reports have shown, however, that the US shipbuilding industry lost its competitive edge years ago due to excessive protectionism. While the US provides discriminatory subsidies amounting to hundreds of billions of dollars to its own industries, it accuses China of adopting so-called nonmarket practices, the Chinese ministry said.

Noting that the US petition is rife with unfounded allegations, distorting normal trade and investment activities as threats to US national security and business interests, the ministry said that "blaming China for the US' own industrial issues lacks factual basis and contradicts economic common sense".

The main anticipated objective of the petitioner unions is to urge the US government to enhance policy support and subsidies for its domestic maritime and shipbuilding industries, said Sun Lei, a trade lawyer at Beijing Dacheng Law Offices, adding that the investigation is expected to last approximately one year.

In an election year, supporting labor unions has become an effortless decision for US President Joe Biden, who, in his bid for reelection, hopes to garner support from the working class, said He Weiwen, executive director of the China Association of International Trade.

The investigation, using an old narrative to position China as a threat, aims to redirect attention away from domestic issues that have long been unaddressed, He said.

The Chinese shipbuilding sector began its path toward internationalization in the 1980s, constructing a large number of vessels for European and Asian shipowners. But there has never been direct competition between the US and Chinese shipbuilding industries, said Zhao Yifei, a professor of Antai College of Economics and Management, which is part of Shanghai Jiao Tong University.

China's shipbuilding industry secured the top position for the 14th consecutive year in 2023 in terms of completed vessels, which accounted for 50.2 percent of the global total, new orders, at 66.6 percent of the world total, and order backlog, at 55 percent, according to data from the China Association of the National Shipbuilding Industry.

Currently, orders from US shipowners account for less than 5 percent of the total new orders received by Chinese shipyards each year. Even if the US government were to impose additional Section 301 tariffs, that would have a minimal impact on Chinese shipyards, Zhao said.

Moreover, the investigation comes as the ability of the US to reindustrialize and regain its position as a global manufacturing powerhouse has become a focal point of the White House in the lead-up to the November election.

The White House also announced on Wednesday that it plans to triple the existing Section 301 tariffs on Chinese steel and aluminum products and collaborate with neighboring countries to prevent indirect exports of Chinese steel and aluminum to the US.

China has urged the US to confront its internal challenges, while also demanding the immediate removal of imposed tariffs, according to a separate statement by China's Commerce Ministry on Thursday.

Since 2018, the US has selectively imposed tariffs on steel and aluminum products from its global trading partners under the guise of "national security". The World Trade Organization has ruled that the US measures violate global trade rules, the ministry said.

However, the US has abused the Section 301 tariff review process, openly demanding arbitrary adjustments to tariffs on Chinese products, the ministry said, adding that US pressure on other countries to restrict Chinese products will further disrupt the security and stability of global industrial and supply chains.

图 9 输入文本

### 4. 输出文本

Trade tension between China and the United States, simmering over Washington's repeated tariff increases and policy restrictions against Chinese exporters and enterprises, has escalated again with Washington's launch of a Section 301 investigation targeting China's maritime, logistics and shipbuilding sectors.

However, the old playbook of unilateralism and protectionism pursued by the White House will not only fail to bring about the reshoring of manufacturing desired by the US, but will also result in the US facing the challenge of a more expensive supply chain, experts said.

China expressed on Wednesday night its strong dissatisfaction with and firm opposition to the US Trade Representative's initiation of the Section 301 investigation, emphasizing that it will keep abreast of the progress of the probe and take all necessary measures to defend its own rights and interests, the Ministry of Commerce said in a statement.

Section 301 of the Trade Act of 1974, as amended, is used to respond to so-called unjustifiable, unreasonable or discriminatory foreign government practices that burden or restrict US commerce.

The investigation petition was jointly filed on March 12 by five US labor unions. "The number of commercial shipyards in the US has plunged by more than 70 percent, tens of thousands of jobs have been lost, and the US now produces only a fraction of 1 percent of the world's commercial vessels, falling to 19th place," they said in the petition.

Multiple US research reports have shown, however, that the US shipbuilding industry lost its competitive edge years ago due to excessive protectionism. While the US provides discriminatory subsidies amounting to hundreds of billions of dollars to its own industries, it accuses China of adopting so-called nonmarket practices, the Chinese ministry said.

Noting that the US petition is rife with unfounded allegations, distorting normal trade and investment activities as threats to US national security and business interests, the ministry said that "blaming China for the US' own industrial issues lacks factual basis and contradicts economic common sense".

The main anticipated objective of the petitioner unions is to urge the US government to enhance policy support and subsidies for its domestic maritime and shipbuilding industries, said Sun Lei, a trade lawyer at Beijing Dacheng Law Offices, adding that the investigation is expected to last approximately one year.

In an election year, supporting labor unions has become an effortless decision for US President Joe Biden, who, in his bid for reelection, hopes to garner support from the working class, said He Weiwen, executive director of the China Association of International Trade.

The investigation, using an old narrative to position China as a threat, aims to redirect attention away from domestic issues that have long been unaddressed, He said.

The Chinese shipbuilding sector began its path toward internationalization in the 1980s, constructing a large number of vessels for European and Asian shipowners. But there has never been direct competition between the US and Chinese shipbuilding industries, said Zhao Yifei, a professor of Antai College of Economics and Management, which is part of Shanghai Jiao Tong University.

China's shipbuilding industry secured the top position for the 14th consecutive year in 2023 in terms of completed vessels, which accounted for 50.2 percent of the global total, new orders, at 66.6 percent of the world total, and order backlog, at 55 percent, according to data from the China Association of the National Shipbuilding Industry.

Currently, orders from US shipowners account for less than 5 percent of the total new orders received by Chinese shipyards each year. Even if the US government were to impose additional Section 301 tariffs, that would have a minimal impact on Chinese shipyards, Zhao said.

Moreover, the investigation comes as the ability of the US to reindustrialize and regain its position as a global manufacturing powerhouse has become a focal point of the White House in the lead-up to the November election.

The White House also announced on Wednesday that it plans to triple the existing Section 301 tariffs on Chinese steel and aluminum products and collaborate with neighboring countries to prevent indirect exports of Chinese steel and aluminum to the US.

China has urged the US to confront its internal challenges, while also demanding the immediate removal of imposed tariffs, according to a separate statement by China's Commerce Ministry on Thursday.

Since 2018, the US has selectively imposed tariffs on steel and aluminum products from its global trading partners under the guise of "national security". The World Trade Organization has ruled that the US measures violate global trade rules, the ministry said.

However, the US has abused the Section 301 tariff review process, openly demanding arbitrary adjustments to tariffs on Chinese products, the ministry said, adding that US pressure on other countries to restrict Chinese products will further disrupt the security and stability of global industrial and supply chains.

图 10 输出文本

## 五、经验体会与不足

**经验体会：**通过实现一个完整的哈夫曼编码与译码系统，熟练掌握了三叉链表链式存储结构和优先级队列（即堆结构）的实现及相关操作，巩固了对存储结构的认识，熟练 c 语言文件操作函数，熟练理解哈夫曼编码的原理到算法实现再到编码文件的压缩与解压的全过程。

**不足：**程序中的一些算法编写的不够高效简洁，还存在时间和空间的不必要消耗，并且程序的输出端口编写的也存在不规范之处。



## 六、附录：源代码（带注释）

```

1.  #include <ctype.h>
2.  #include <stdio.h>
3.  #include <stdlib.h>
4.  #include <string.h>
5.  #include <sys/stat.h>
6.
7.  #define N 128
8.  #define M 2 * N - 1
9.  // 三叉链表存储结构
10. typedef struct {
11.     double weight;
12.     int lchild;
13.     int rchild;
14.     int parent;
15. } HTNODE;
16. typedef HTNODE HuffmanT[M];
17. // 编码表的存储结构
18. typedef struct {
19.     char ch;
20.     char bits[M];
21.     double frequency;
22. } CodeNode;
23. typedef CodeNode HuffmanCode[N];
24. // 最小堆存储结构
25. typedef struct {
26.     int key;
27.     double weight;
28. } Elementype;
29. typedef struct {
30.     Elementype data[M];
31.     int n;
32. } HEAP;
33.
34. long long Number_Character; // 字符总数
35. int Number_Tpye;           // 字符种类数
36. // 堆的初始化
37. void Heap_Init(HEAP *heap) { heap->n = 0; }
38. // 判断堆空
39. int Heap_Empty(HEAP *heap) { return (!heap->n); }
40. // 判断堆满
41. int Heap_Full(HEAP *heap) { return (heap->n == N - 1); }

```

```

42. // 结点插入(基于优先级队列)
43. void Heap_Insert(HEAP *heap, Elementype item) {
44.     int i;
45.     if (!Heap_Full(heap)) {
46.         heap->n++;
47.         i = heap->n;
48.         while ((i != 1) && (item.weight < heap->data[i /
49. 2].weight)) {
50.             heap->data[i] = heap->data[i / 2];
51.             i /= 2;
52.         }
53.         heap->data[i].key = item.key;
54.         heap->data[i].weight = item.weight;
55.     }
56. // 删除最堆的最小项
57. Elementype Heap_DeleteMin(HEAP *heap) {
58.     int parent = 1;
59.     int child = 2;
60.     Elementype item, tmp;
61.     if (!Heap_Empty(heap)) {
62.         item = heap->data[1];
63.         tmp = heap->data[heap->n--];
64.         while (child <= heap->n) {
65.             if ((child < heap->n) &&
66.                 (heap->data[child].weight > heap->data[ch
67. ild + 1].weight)) {
68.                 child++;
69.             }
70.             if (tmp.weight <= heap->data[child].weight) {
71.                 break;
72.             }
73.             heap->data[parent] = heap->data[child];
74.             parent = child;
75.             child *= 2;
76.         }
77.         heap->data[parent] = tmp;
78.         return item;
79.     }
80. // 读取英文文本, 统计字符总数和字符种类数
81. void ReadText(HuffmanCode H) {
82.     int flag;

```

```

83.     int number_total, type_total;
84.     char Single_Character;
85.     int Character_Count[N] = {0};
86.     flag = 0;
87.     number_total = 0;
88.     type_total = 0;
89.     FILE *file = fopen("./Input/InputText.txt", "r");
90.     if (file == NULL) {
91.         printf("打开文件失败\n");
92.         exit(1);
93.     }
94.     while ((Single_Character = fgetc(file)) != EOF) {
95.         if (isascii(Single_Character)) {
96.             Character_Count[(int)Single_Character]++;
97.             number_total++;
98.         }
99.     }
100.    fclose(file);
101.    for (int i = 0; i < N; i++) {
102.        if (Character_Count[i]) {
103.            H[flag].ch = i;
104.            H[flag].frequency = (double)Character_Count[i]
/ number_total;
105.            flag++;
106.            type_total++;
107.        }
108.    }
109.    Number_Character = number_total;
110.    Number_Tpye = type_total;
111. }
112. // 哈夫曼树初始化
113. void InitHT(HuffmanCode H, HuffmanT T) {
114.     for (int i = 0; i < Number_Tpye; i++) {
115.         T[i].weight = H[i].frequency;
116.     }
117.     for (int i = 0; i < M; i++) {
118.         T[i].parent = -1;
119.         T[i].lchild = -1;
120.         T[i].rchild = -1;
121.     }
122. }
123. // 找到两个最小项(基于三叉链表)
124. void SelecMin(HuffmanT T, int n, int *p1, int *p2) {
125.     int i, j, temp;

```

```

126.     for (i = 0; i < n; i++) {
127.         if (T[i].parent == -1) {
128.             *p1 = i;
129.             break;
130.         }
131.     }
132.     for (j = i + 1; j < n; j++) {
133.         if (T[j].parent == -1) {
134.             *p2 = j;
135.             break;
136.         }
137.     }
138.     for (i = 0; i < n; i++) {
139.         if ((T[*p1].weight > T[i].weight) && (T[i].parent
140.             == -1) &&
141.             (*p2 != i)) {
142.             *p1 = i;
143.         }
144.         for (j = 0; j < n; j++) {
145.             if ((T[*p2].weight > T[j].weight) && (T[j].parent
146.                 == -1) &&
147.                 (*p1 != j)) {
148.                 *p2 = j;
149.             }
150.             if (T[*p1].weight > T[*p2].weight) {
151.                 temp = *p1;
152.                 *p1 = *p2;
153.                 *p2 = temp;
154.             }
155.         }
156.         // 构造哈夫曼树(基于三叉链表)
157.         void CreatHT(HuffmanCode H, HuffmanT T) {
158.             int p1, p2, n;
159.             n = Number_Tpye;
160.             InitHT(H, T);
161.             for (int i = n; i < 2 * n - 1; i++) {
162.                 SelecMin(T, i, &p1, &p2);
163.                 T[p1].parent = T[p2].parent = i;
164.                 T[i].lchild = p1;
165.                 T[i].rchild = p2;
166.                 T[i].weight = T[p1].weight + T[p2].weight;
167.             }

```

```

168. }
169. // 哈夫曼编码算法的实现
170. void CharSetHuffmanEncoding(HuffmanCode H, HuffmanT T) {
171.     int c, p, i;
172.     char cd[Number_Tpye];
173.     int start;
174.     cd[Number_Tpye] = '\0';
175.     for (i = 0; i < Number_Tpye; i++) {
176.         start = Number_Tpye;
177.         c = i;
178.         while ((p = T[c].parent) >= 0) {
179.             cd[--start] = (T[p].lchild == c) ? '0' : '1';
180.             c = p;
181.         }
182.         strcpy(H[i].bits, &cd[start]);
183.     }
184. }
185. // 利用哈夫曼树进行编码，存储成压缩文件
186. void Compressed(HuffmanCode H, HuffmanT T) {
187.     char Single_Character;
188.     FILE *file;
189.     file = fopen("./Input/InputText.txt", "r");
190.     if (!file) {
191.         printf("打开文件失败\n");
192.         exit(1);
193.     }
194.     FILE *bitfile;
195.     bitfile = fopen("./Output//Output.Huffman", "wb");
196.     if (!bitfile) {
197.         printf("打开文件失败\n");
198.         exit(1);
199.     }
200.     // 压入字符种类数
201.     fwrite(&Number_Tpye, sizeof(int), 1, bitfile);
202.     // 压入各字符
203.     for (int i = 0; i < Number_Tpye; i++) {
204.         fwrite(&H[i].ch, sizeof(char), 1, bitfile);
205.     }
206.     // 压入哈夫曼树节点
207.     for (int i = Number_Tpye; i < 2 * Number_Tpye - 1; i+
+) {

```



```

208.         fwrite(&T[i].lchild, sizeof(unsigned char), 1, bi
        tfile);
209.         fwrite(&T[i].rchild, sizeof(unsigned char), 1, bi
        tfile);
210.     }
211.     // 压入字符数
212.     fwrite(&Number_Character, sizeof(Number_Character), 1,
        bitfile);
213.     int flag;
214.     int bitflag;
215.     int i;
216.     unsigned char bit = 0;
217.     bitflag = 0;
218.     while ((Single_Character = fgetc(file)) != EOF) {
219.         for (i = 0; i < Number_Tpye; i++) {
220.             if (H[i].ch == Single_Character) {
221.                 break;
222.             }
223.         }
224.         for (flag = 0; H[i].bits[flag]; flag++) {
225.             if (H[i].bits[flag] == '0') {
226.                 bit = bit << 1;
227.             } else {
228.                 bit = (bit << 1) + 1;
229.             }
230.             bitflag++;
231.             if (bitflag == 8) {
232.                 fwrite(&bit, sizeof(unsigned char), 1, bi
                tfile);
233.                 bitflag = 0;
234.                 bit = 0;
235.             }
236.         }
237.     }
238.     if (bitflag) {
239.         bit = bit << (8 - bitflag);
240.         fwrite(&bit, sizeof(unsigned char), 1, bitfile);
241.         fwrite(&bitflag, sizeof(unsigned char), 1, bitfil
        e);
242.     }
243.     fclose(file);
244.     fclose(bitfile);
245. }

```

```

246. // 将哈夫曼编码文件译码为文本文件
247. void Decode(HuffmanCode H, HuffmanT T) {
248.     int type_total;
249.     long long number_total;
250.     FILE *bfile = fopen("./Output/Output.Huffman", "rb");

251.     FILE *file = fopen("./Output/OutputText.txt", "w");
252.     if (!bfile) {
253.         printf("打开文件失败\n");
254.         exit(1);
255.     }
256.     // 读取字符种类数
257.     fread(&type_total, sizeof(type_total), 1, bfile);
258.     // 读取各字符
259.     for (int i = 0; i < type_total; i++) {
260.         fread(&H[i].ch, sizeof(H[i].ch), 1, bfile);
261.     }
262.     // 读取哈夫曼树节点
263.     for (int i = type_total; i < 2 * type_total - 1; i++)
264.     {
265.         unsigned char p1, p2;
266.         fread(&p1, sizeof(p1), 1, bfile);
267.         fread(&p2, sizeof(p2), 1, bfile);
268.         T[i].lchild = (int)p1;
269.         T[i].rchild = (int)p2;
270.     }
271.     // 读取字符数
272.     fread(&number_total, sizeof(number_total), 1, bfile);

273.     printf("这篇文章有 %lld 个字符，共 %d 种字符种类数。
274.         \n", number_total,
275.             type_total);
276.     unsigned char b1;
277.     fread(&b1, sizeof(b1), 1, bfile);
278.     int i = 0;
279.     int j = 2 * type_total - 2;
280.     int total = 0;
281.     while (total < number_total) {
282.         if (((b1 >> (7 - i)) & 1))
283.             j = T[j].rchild;
284.         else
285.             j = T[j].lchild;
286.         i++;
287.         if (j < type_total) {

```

```

286.         fprintf(file, "%c", H[j].ch);
287.         total++;
288.         j = 2 * type_total - 2;
289.     }
290.     if (i == 8) {
291.         i = 0;
292.         fread(&b1, sizeof(b1), 1, bfile);
293.     }
294. }
295. fclose(bfile);
296. fclose(file);
297. }
298. // 计算哈夫曼树的编码长度和编码文件的压缩率
299. void Calculate(HuffmanCode H, HuffmanT T) {
300.     double lenth = 0;
301.     for (int i = 0; i < Number_Tpye; i++) {
302.         lenth = lenth + strlen(H[i].bits) * H[i].frequenc
303.     }
304.     printf("哈夫曼树编码长度为:%lf\n", lenth);
305.     double ratio_theoretical;
306.     double ratio_real;
307.     ratio_theoretical = 1 - lenth / 8;
308.     printf(" 编 码 文 件 的 理 论 压 缩 率
309.     为:%lf\n", ratio_theoretical);
310.     struct _stat text;
311.     struct _stat bittext;
312.     double Size_Text;
313.     double Size_BitText;
314.     _stat("./Input/InputText.txt", &text);
315.     _stat("./Output/Output.Huffman", &bittext);
316.     Size_Text = text.st_size;
317.     Size_BitText = bittext.st_size;
318.     ratio_real = Size_BitText / Size_Text;
319.     printf("编码文件的真实压缩率为:%lf\n", ratio_real);
320. }
321. // 创建堆结构
322. void CreatHeap(HuffmanCode H, HEAP *heap) {
323.     Heap_Init(heap);
324.     for (int i = 0; i < Number_Tpye; i++) {
325.         Elementype new_node;
326.         new_node.key = i;
327.         new_node.weight = H[i].frequency;
328.         Heap_Insert(heap, new_node);

```

```

328.     }
329. }
330. // 找到两个最小项(基于优先级队列)
331. void SelecMin_Heap(HEAP *heap, int *p1, int *p2) {
332.     *p1 = Heap_DeleteMin(heap).key;
333.     *p2 = Heap_DeleteMin(heap).key;
334. }
335. // 构造哈夫曼树(基于优先级队列)
336. void CreatHT_Heap(HuffmanCode H, HuffmanT T) {
337.     int p1, p2, n;
338.     n = Number_Tpye;
339.     HEAP heap;
340.     InitHT(H, T);
341.     CreatHeap(H, &heap);
342.     for (int i = n; i < 2 * n - 1; i++) {
343.         Elementype new_node;
344.         SelecMin_Heap(&heap, &p1, &p2);
345.         T[p1].parent = T[p2].parent = i;
346.         T[i].lchild = p1;
347.         T[i].rchild = p2;
348.         T[i].weight = T[p1].weight + T[p2].weight;
349.         new_node.key = i;
350.         new_node.weight = T[i].weight;
351.         Heap_Insert(&heap, new_node);
352.     }
353. }
354. // 主函数
355. int main() {
356.     HuffmanCode H;
357.     HuffmanT T;
358.     HuffmanCode H_decode;
359.     HuffmanT T_decode;
360.     while (1) {
361.         int Function_Options;
362.         printf(
363.             "请选择功能实现:1.编码(基于三叉链表);2.编码(基于
              优先级队列);3."
364.             "压缩文本;"
365.             "4.译码并计算编码信息;0.退出\n");
366.         scanf("%d", &Function_Options);
367.         if (Function_Options == 1) {
368.             ReadText(H);
369.             CreatHT(H, T);
370.             CharSetHuffmanEncoding(H, T);

```

```

371.             printf("这个文本有%d 个字符,有%d 种字符
    ", Number_Character,
372.               Number_Tpye);
373.             printf("文本中各字符的哈夫曼编码及其频率:\n");
374.             for (int i = 0; i < Number_Tpye; i++) {
375.                 printf(" 字符%c      编 码 :%s      频
    率:%lf\n", H[i].ch, H[i].bits,
376.                   H[i].frequency);
377.             }
378.         } else if (Function_Options == 2) {
379.             ReadText(H);
380.             CreatHT_Heap(H, T);
381.             CharSetHuffmanEncoding(H, T);
382.             printf("这个文本有%d 个字符,有%d 种字符
    ", Number_Character,
383.               Number_Tpye);
384.             printf("文本中各字符的哈夫曼编码及其频率:\n");
385.             for (int i = 0; i < Number_Tpye; i++) {
386.                 printf(" 字符%c      编 码 :%s      频
    率:%lf\n", H[i].ch, H[i].bits,
387.                   H[i].frequency);
388.             }
389.         } else if (Function_Options == 3) {
390.             Compressed(H, T);
391.         } else if (Function_Options == 4) {
392.             Decode(H_decode, T_decode);
393.             Calculate(H, T);
394.         } else if (Function_Options == 0) {
395.             return 0;
396.         } else {
397.             printf("输入值非法,请重新输入\n");
398.         }
399.     }
400.     return 0;

```