

编程实验：

1. 实验目的与要求

在学习和理解二叉树的原理、构造及遍历方法的基础上，应用所学知识来解决实际问题。

本实验将通过一个实际应用问题的解决过程掌握 Huffman 树的构造、Huffman 编码的生成及基于涉及的知识点包括树的构造、遍历及 C 语言位运算和二进制文件。

2. 实验内容

Huffman 编码文件压缩

【问题描述】

编写一程序采用 Huffman 编码对一个正文文件进行压缩。具体压缩方法如下：

1. 对正文文件中字符(换行字符'\n'除外，不统计)按出现次数（即频率）进行统计。
2. 依据字符频率生成相应的 Huffman 树（未出现的字符不生成）。
3. 依据 Huffman 树生成相应字符的 Huffman 编码。
4. 依据字符 Huffman 编码压缩文件（即将源文件字符按照其 Huffman 编码输出）。

说明：

1. 只对文件中出现的字符生成 Huffman 树，注意：一定不要处理\n，即不要为其生成 Huffman 树。
2. 采用 ASCII 码值为 0 的字符作为压缩文件的结束符（即可将其出现次数设为 1 来参与编码）。
3. 在生成 Huffman 树前，初始在对字符频率权重进行（由小至大）排序时，频率相同的字符 AS 到有序权重序列中时，若出现相同权重，则将新生成的权重节点插入到原有相同权重节点之后（
4. 在生成 Huffman 树时，权重节点在前的作为左孩子节点，权重节点在后的作为右孩子节点。
5. 遍历 Huffman 树生成字符 Huffman 码时，左边为 0 右边为 1。
6. 源文件是文本文件，字符采用 ASCII 编码，每个字符占 8 个二进制位；而采用 Huffman 编码此最后输出时需要使用 C 语言中的位运算将字符的 Huffman 码依次输出到每个字节中。

【输入形式】

对当前目录下文件 input.txt 进行压缩。

【输出形式】

将压缩后结果输出到文件 output.txt 中，同时将压缩结果用十六进制形式（printf("%x",...)）输出到

3. 实验准备

1. 文件下载

从教学平台（judge.buaa.edu.cn）课程下载区下载文件 lab_tree2.rar，该文件中包括了本实验中用
1 Huffman2student.c：该文件给出本实验程序的框架，框架中部分内容未完成（见下面相关实验步骤）
程序运行后得到相应要求的运行结果；

1 input.txt：为本实验的测试数据。

2. Huffman2student.c 文件中相关数据结构说明

结构类型说明：

```
struct tnode {    //Huffman 树结构节点类型
    char c;
    int weight;
    struct tnode *left;
    struct tnode *right;
};
```

结构类型 struct tnode 用来定义 Huffman 树的节点，其中：

1) 对于树的叶节点，成员 c 和 weight 用来存放字符及其出现次数；对于非叶节点来说，c 值可不
点生成条件，若 p 为当前 Huffman 树节点指针，则有：

p->weight = p->left->weight + p->right->weight;

2) 成员 left 和 right 分别为 Huffman 树节点左右子树节点指针。

全局变量说明:

```
int Ccount[128]={0};
```

```
struct tnode *Root=NULL;
```

```
char HCode[128][MAXSIZE]={0};
```

```
int Step=0;
```

```
FILE *Src, *Obj;
```

整型数组 Ccount 存放每个字符的出现次数, 如 Ccount['a']表示字符 a 的出现次数。

变量 Root 为所生成的 Huffman 树的根节点指针。

数组 HCode 用于存储字符的 Huffman 编码, 如 HCode['a']为字符 a 的 Huffman 编码, 本实验中为

变量 Step 为实验步骤状态变量, 其取值为 1、2、3、4, 分别对应实验步骤 1、2、3、4。

变量 Src、Obj 为输入输出的文件指针, 分别用于打开输入文件 “input.txt” 和输出文件 “output.txt”。

4. 实验步骤

【步骤 1】

1) 实验要求

在程序文件 huffman2student.c 中 “//【实验步骤 1】开始” 和 “//【实验步骤 1】结束” 间编写相应代码, 统计 input.txt 中字符出现频率。

//【实验步骤 1】开始

```
void statCount()
```

```
{
```

```
}
```

//【实验步骤 1】结束

2) 实验说明

函数 statCount 用来统计输入文件 (文件指针为全局变量 Src) 中字符的出现次数 (频率), 并将统计结果存入 Ccount['a']存放字符 a 的出现次数。

注意: 在该函数中 Ccount[0]一定要置为 1, 即 Ccount[0]=1。编码值为 0 ('\0') 的字符用来作为

3) 实验结果

函数 print1()用来打印输出步骤 1 的结果, 即输出数组 Ccount 中字符出现次数多于 0 的字符及次数。完成【实验步骤 1】编码后, 本地编译并运行该程序, 并在标准输入中输入 1, 程序运行正确时在屏幕上将输出

```

NUL:1
:65
":4
':1
,:2
.:9
D:4
E:1
I:3
R:1
W:2
a:23
b:2
c:2
d:5
e:41
f:2
g:5
h:14
i:10
k:2
l:15
m:6
n:14
o:35
p:7
r:14
s:18
t:24
u:11
v:8
w:3
x:1
y:16

```

图 1 步骤 1 运行结果

在本地运行正确的情况下，将你所编写的程序文件中“//【实验步骤 1】开始”和“//【实验步骤 1】结束”间的代码【实验步骤 1】下的框中，然后点击提交按钮，若得到如下运行结果（测试数据 1 评判结果）

测试数据	评判结果
测试数据1	完全正确
测试数据2	输出错误
测试数据3	输出错误
测试数据4	输出错误

表明实验步骤 1：通过，否则：不通过。

【步骤 2】

1) 实验要求

在程序文件 huffman2student.c 中的“//【实验步骤 2】开始”和“//【实验步骤 2】结束”间编写一个根结点指针为 Root 的 Huffman 树。

//【实验步骤 2】开始

```
void createHTree()
```

```
{
```

```
}
```

//【实验步骤 2】结束

2) 实验说明

在程序文件 huffman2student.c 中函数 createHTree 将根据每个字符的出现次数（字符出现次数存放

码值为 i 的字符出现次数), 按照 Huffman 树生成规则, 生成一棵 Huffman 树。

算法提示:

1. 依据数组 $Ccount$ 中出现次数不为 0 的 (即 $Ccount[i]>0$) 项, 构造出树林 $F=\{T_0, T_1, \dots, T_m\}$, 且根结点(叶结点)的权值为相应字符的出现次数的二叉树 (每棵树结点的类型为 $struct\ tnode$, 其

```
for(i=0; i<128; i++)
```

```
if(Ccount[i]>0){
```

```
    p = (struct tnode *)malloc(sizeof(struct tnode));
```

```
    p->c = i; p->weight = Ccount[i];
```

```
    p->left = p->right = NULL;
```

```
    add p into F;
```

```
}
```

2. 对树林 F 中每棵树按其根结点的权值由小至大进行排序 (排序时, 当权值 $weight$ 相同时, 字

3. while 树个数 > 1 in F

a) 将 F 中 T_0 和 T_1 作为左、右子树合并成为一棵新的二叉树 T' , 并令 $T' \rightarrow weight = T_0 \rightarrow weight + T_1 \rightarrow weight$

b) 删除 T_0 和 T_1 from F , 同时将 T' 加入 F 。要求加入 T' 后 F 仍然有序。若 F 中有树根结点权

4. $Root = T_0$ ($Root$ 为 Huffman 树的根结点指针。循环结束时, F 中只有一个 T_0)

注: 在实现函数 `createHTree` 时, 在框中还可根据需要定义其它函数, 例如:

```
void myfun()
```

```
{
```

```
    ...
```

```
}
```

```
void createHTree()
```

```
{
```

```
    ...
```

```
    myfun();
```

```
    ...
```

```
}
```

3) 实验结果

函数 `print2()` 用来打印输出步骤 2 的结果, 即按前序遍历方式遍历步骤 2 所生成 (由全局变量 `Root` 指向) 的二叉树。遍历过程中, 遇到编码值为 0 的字符用 `NUL` 表示、空格符用 `SP` 表示、制表符用 `TAB` 表示、回车符用 `CR` 表示。程序运行正确时在屏幕上将输出如下结果:



图 2 步骤 2 运行结果

在本地运行正确的情况下, 将你在本地所编写的程序文件中 //【实验步骤 2】开始”和“//【实验步骤 2】结束”之间的代码复制到在线实验平台后所附代码【实验步骤 2】下的框中, 然后点击提交按钮, 若得到如下运行结果 (测试数据 2 评

测试数据

测试数据1

测试数据2

测试数据3

测试数据4

评判结果

完全正确

完全正确

输出错误

输出错误

表明**实验步骤 2**：通过，否则：不通过。

【步骤 3】

1) 实验要求

在程序文件 huffman2student.c 中的 “//**【实验步骤 3】**开始” 和 “//**【实验步骤 3】**结束” 间编写实现**【实验步骤 3】**中所产生的 Huffman 树为文本中出现的每个字符生成对应的 Huffman 编码。遍历右边为 1。

// **【实验步骤 3】** 开始

```
void makeHCode()
```

```
{
```

```
}
```

// **【实验步骤 3】** 结束

2) 实验说明

【步骤 3】 依据 **【步骤 2】** 所生成的根结点为 Root 的 Huffman 树生为文本中出现的每个字符生成如下：

```
char HCode[128][MAXSIZE];
```

HCode 变量用来存放每个字符的 Huffman 编码串，如 HCode['e']存放的是字母 e 的 Huffman 编码串。

算法提示：

可编写一个按前序遍历方法对根节点为 Root 的树进行遍历的递归函数，并在遍历过程中用一个字符串记录经过的路径（经过的边），经过左边时记录为 '0'，经过右边时记录为 '1'；当遍历节点为叶节点时，即执行 strcpy(HCode[p->c],路径串)。

注：在实现函数 makeHCode 时，在框中还可根据需要定义其它函数，如调用一个有类于前序遍历的函数生成 Huffman 编码：

```
void visitHTree()
```

```
{
```

```
...
```

```
}
```

```
void makeHCode()
```

```
{
```

```
...
```

```
visitHTree();
```

```
...
```

```
}
```

3) 实验结果

函数 print3()用来打印输出步骤 3 的结果，即输出步骤 3 所生成的存储在全局变量 HCode 中非空字符串。编译并运行该程序，并在标准输入中输入 3，在屏幕上将输出 ASCII 字符与其 Huffman 编码的 Huffman 编码，其中 NUL 表示 ASCII 编码为 0 的字符，SP 表示空格字符编码值为 0 的字符用

```

NUL:01001010
SP:111
":000000
':01001011
,:10011111
.:01000
D:000001
E:01011100
I:0101111
R:01011101
W:0000100
a:1000
b:0000101
c:0000110
d:010011
e:011
f:0000111
g:010110
h:10110
i:01010
k:0100100
l:11001
m:100110
n:10111
o:001
p:110100
r:11000
s:0001
t:1010
u:10010
v:110101
w:1001110
x:10011110
y:11011

```

图 3 步骤 3 运行结果

在本地运行正确的情况下，将你在本地所编写的程序文件中//【实验步骤 3】开始”和“//【实验步骤 3】后所附代码【实验步骤 3】下的框中，然后点击提交按钮，若得到如下运行结果（测试数据 3 评

测试数据	评判结果
测试数据1	完全正确
测试数据2	完全正确
测试数据3	完全正确
测试数据4	输出错误

表明实验步骤 3：通过，否则：不通过。

【步骤 4】

1) 实验要求

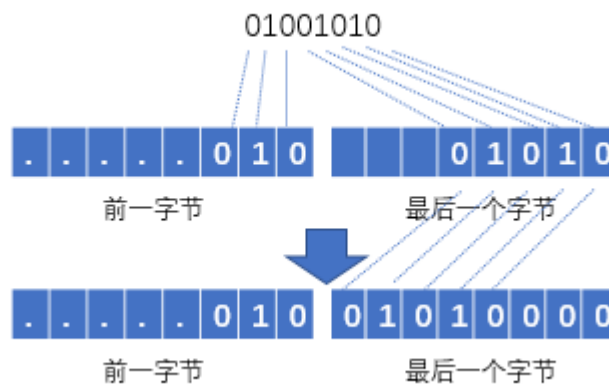
在程序文件 huffman2student.c 函数中的“//【实验步骤 4】开始”和“//【实验步骤 4】结束”间，根据【实验步骤 3】中所生成的字符 ASCII 码与 Huffman 编码对应表（存储在全局变量 HCode 中，编码串，在本实验中值为字符串”011”），将原文本文件（文件指针为 Src）内容（ASCII 字符）转换（文件指针为 Obj）中，以实现文件压缩。同时将输出结果用十六进制形式（printf(“%x”,...)）输


```
printf("%x",hc); //按十六进制输出到屏幕上
}
}
```

...

说明：

1. 当遇到源文本文件输入结束时，应将输入结束符的 Huffman 码放到 Huffman 编码串最后，即
2. 在处理完成所有 Huffman 编码串时（如上述算法结束时），处理源文本最后一个字符（文件结束符）时，可能出现如下情况：其子串“010”位于前一个字节中输出，而子串“01010”位于另（最后）左端的头，最后 3 位补 0，然后再输出最后一个字节。



注：在实现函数 `atoHZIP` 时，在框中还可根据需要定义其它函数或全局变量，如：

```
void myfun()
{
    ...
}
void atoHZIP()
{
    ...
    myfun();
    ...
}
```

1) 实验结果

函数 `print4()` 用来打印输出步骤 4 的结果，即根据输出步骤 3 所生成的存储在全局变量 `HCode` 中 Huffman 编码的 ASCII 字符转换为 Huffman 编码字符输出到文件 `output.txt` 中，同时按十六进行输出到屏幕上。运行程序，并在标准输入中输入 4，在屏幕上将输出：


```
4
5fe72b39eb2b57f665c4ccfc27aa833f14cabe541da1722bcc7fe0
c96b42cd09f31c3f568bf356781ee7b461abeccbd1e99f174fa67a
abe90ab8a175cce60af1e75a2be541dfa79c8fd992013ee6bd4247
e8fb38c577abc6fd59575beccbd9c6109f735c74774fc673ead15f
665ed1ab4027dcd71cb7a781f9375bc0fd9979d17a404b3bbd5e3e
ccb8c6c0ebfe4eafecc89ffa65cbe31beaaf6376359da0259eaf57
8fb32e31b017f89b89c6364ffd32e5fc9293aaf47d59f35678fa3c
9ba58f7b6284a
原文件大小：370B
压缩后文件大小：200B
压缩率：45.95%
```

图 4 步骤 4 运行结

说明：

从屏幕输出结果可以看出，由于采用了不定长的 Huffman 编码，且出现频率高的字符的编码长度 200 字节，文件压缩了 45.95%。

在本地运行正确的情况下，将你在本地所编写的程序文件中//【实验步骤 4】开始”和“//【实验报告后所附代码【实验步骤 4】下的框中，然后点击提交按钮，若得到如下运行结果（测试

测试数据	评判结果
测试数据1	完全正确
测试数据2	完全正确
测试数据3	完全正确
测试数据4	完全正确

表明实验步骤 4：通过，否则：不通过。
