

**ECE 264 Spring 2023**  
***Advanced* C Programming**

Aravind Machiry  
Purdue University

# Huffman Compression 01

## Build Tree and Compress

# Fixed-Length vs Variable-Length Code

- ASCII: fixed-length code, every character needs 8 bits
- Some characters (such as 's' and 'e') are more often than some others (such as 'q' and 'z'). Variable length can be more efficient:
  - fewer bits for frequently used characters
  - more bits for rarely used characters
  - ⇒ fewer bits per character on average
- General design principle: optimize for the frequent cases
- This is **lossless** compression. Data can be fully recovered.

# Where is data compression used?

- Everywhere
- Image, video, audio (lossy)
- File download
- When network is limited (in data rate), slow, or unstable

# Huffman Coding (Compression)

Lossless compression

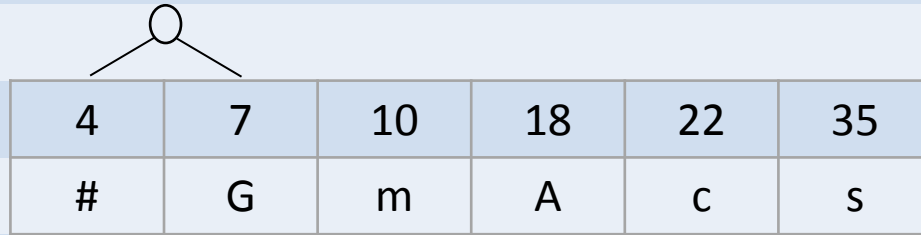
1. Count the occurrences of the characters (may include symbols and unprintable characters)
2. Sort the characters by their occurrences in the ascending order
3. Take the two least occurrences, make them left and right children of the same parent node, add the occurrences and sort in the ascending order again
4. Continue 3 until only one node is left

|            |   |    |   |    |    |    |
|------------|---|----|---|----|----|----|
| occurrence | 4 | 18 | 7 | 22 | 10 | 35 |
| letter     | # | A  | G | c  | m  | s  |

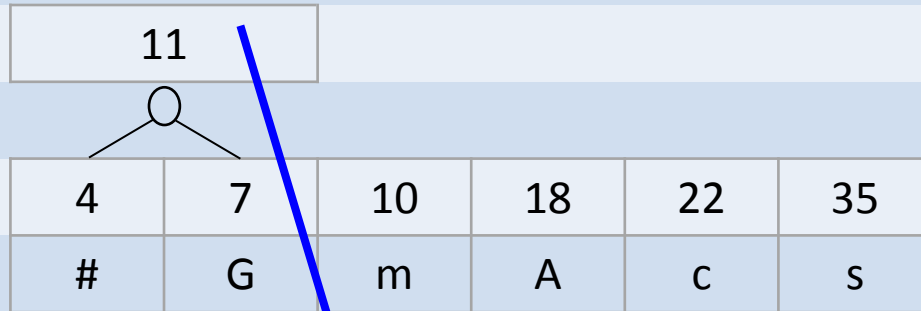
↓ Sort by the occurrences in ascending order

|   |   |    |    |    |    |
|---|---|----|----|----|----|
| 4 | 7 | 10 | 18 | 22 | 35 |
| # | G | m  | A  | c  | s  |

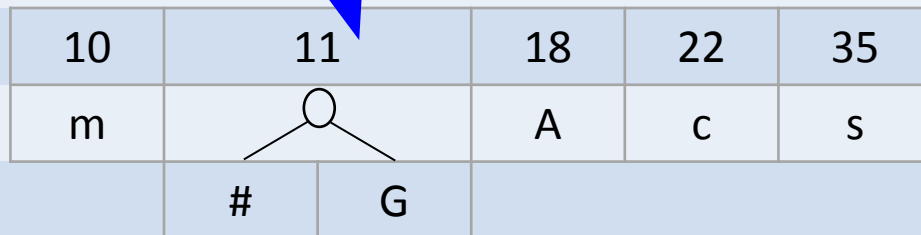
↓ Make the first two siblings of a binary tree

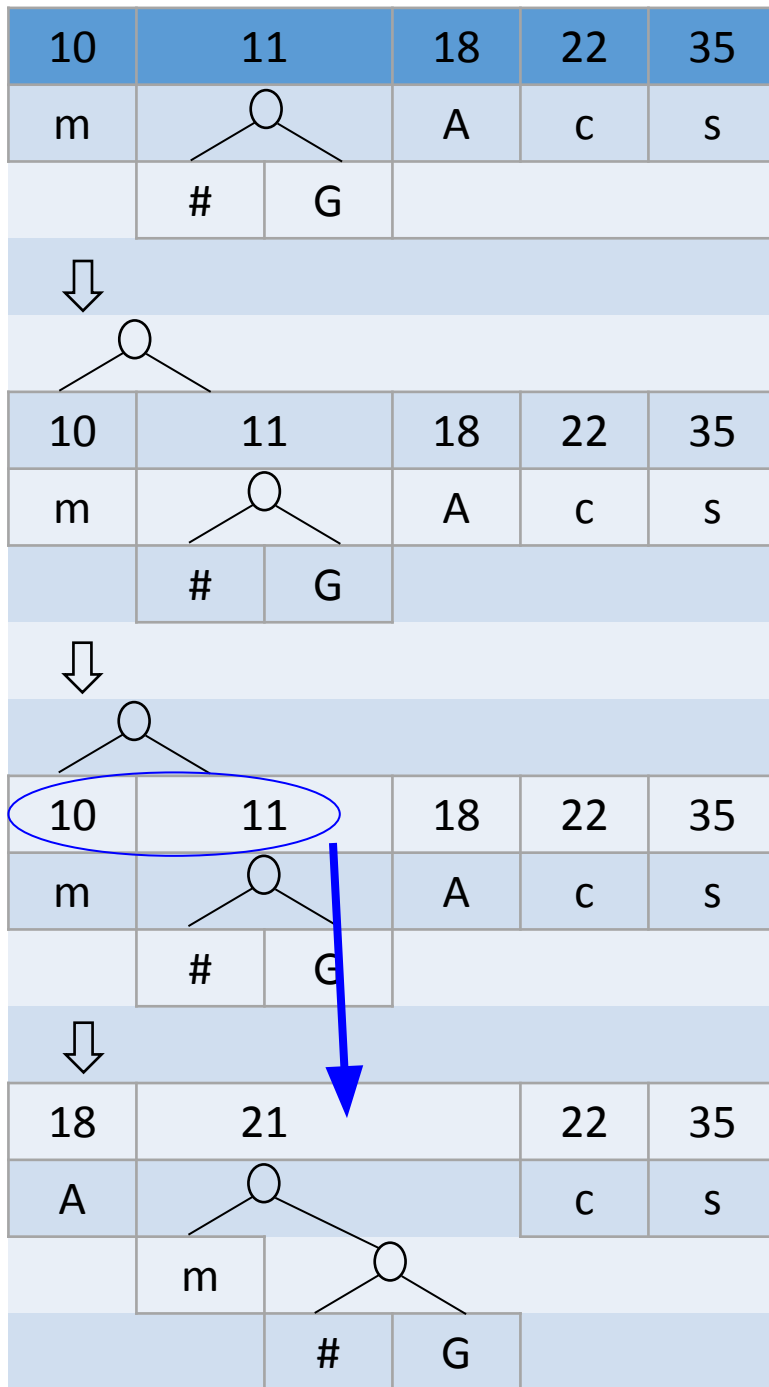


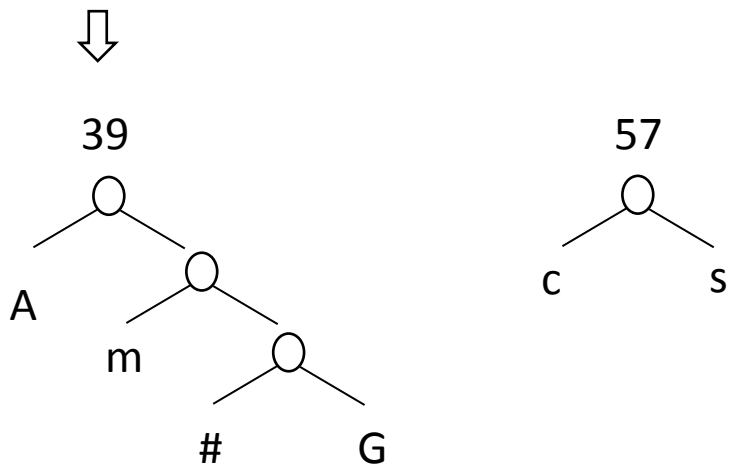
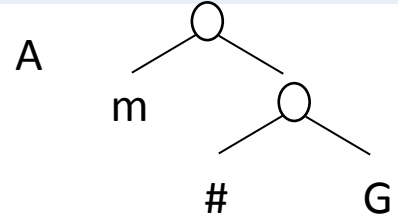
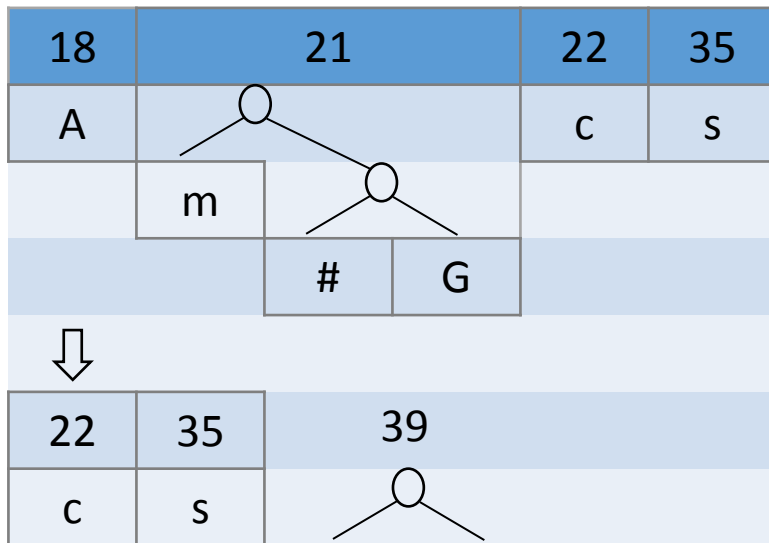
↓ The occurrence of the parent is the sum



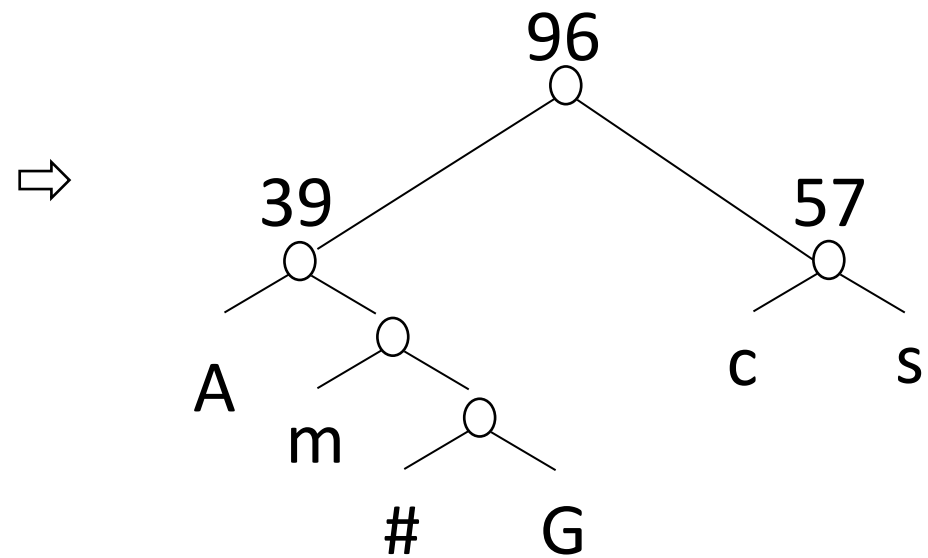
↓ Insert the parent back in ascending order



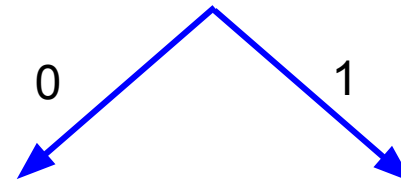
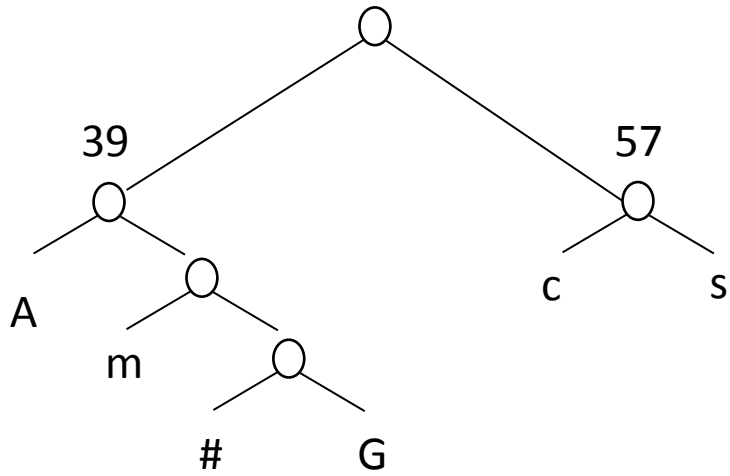




Only the leaf nodes contain characters







| character | occurrence | code |   |   |   | length |
|-----------|------------|------|---|---|---|--------|
| A         | 18         | 0    | 0 |   |   | 2      |
| m         | 10         | 0    | 1 | 0 |   | 3      |
| #         | 4          | 0    | 1 | 1 | 0 | 4      |
| G         | 7          | 0    | 1 | 1 | 1 | 4      |
| c         | 22         | 1    | 0 |   |   | 2      |
| s         | 35         | 1    | 1 |   |   | 2      |

| character | occurrence | code |   |   |   | length |
|-----------|------------|------|---|---|---|--------|
| A         | 18         | 0    | 0 |   |   | 2      |
| m         | 10         | 0    | 1 | 0 |   | 3      |
| #         | 4          | 0    | 1 | 1 | 0 | 4      |
| G         | 7          | 0    | 1 | 1 | 1 | 4      |
| c         | 22         | 1    | 0 |   |   | 2      |
| s         | 35         | 1    | 1 |   |   | 2      |

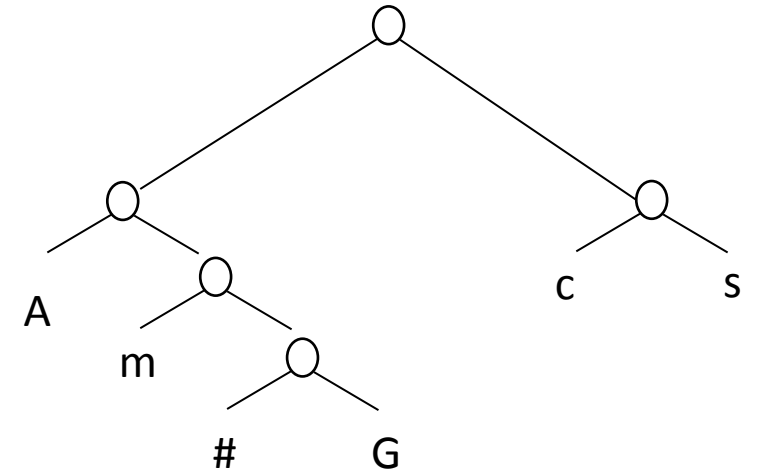
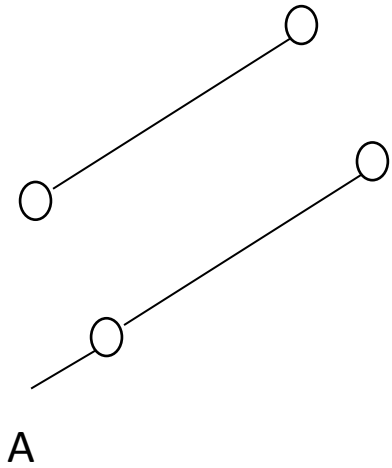
- If occurrence (X) < occurrence (Y)  
 $\Rightarrow$  code length (X)  $\geq$  code length (Y)
- code length (X) > code length (Y)  
 $\Rightarrow$  occurrence (X) < occurrence (Y) **WRONG**

| character | occurrence | code |   |   |   |  | length |
|-----------|------------|------|---|---|---|--|--------|
| A         | 18         | 0    | 0 |   |   |  | 2      |
| m         | 10         | 0    | 1 | 0 |   |  | 3      |
| #         | 4          | 0    | 1 | 1 | 0 |  | 4      |
| G         | 7          | 0    | 1 | 1 | 1 |  | 4      |
| c         | 22         | 1    | 0 |   |   |  | 2      |
| s         | 35         | 1    | 1 |   |   |  | 2      |

|        |    |    |    |    |      |     |      |    |    |    |     |
|--------|----|----|----|----|------|-----|------|----|----|----|-----|
| input  | A  | A  | c  | s  | #    | m   | G    | c  | s  | A  | ... |
| output | 00 | 00 | 10 | 11 | 0110 | 010 | 0111 | 10 | 11 | 00 | ... |

input 0000101101100100111101100

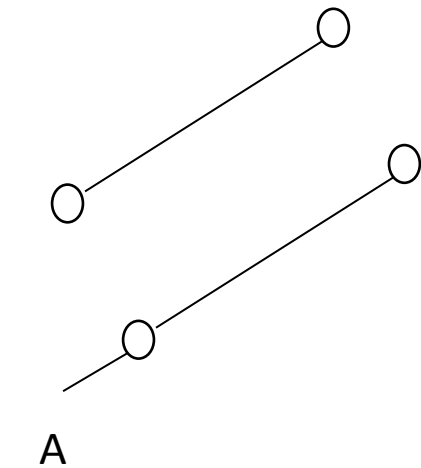
0



00 ⇒ A

|       |    |                         |
|-------|----|-------------------------|
| input | 00 | 00101101100100111101100 |
|       | A  |                         |

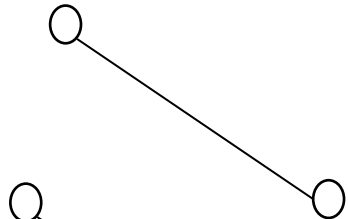
0



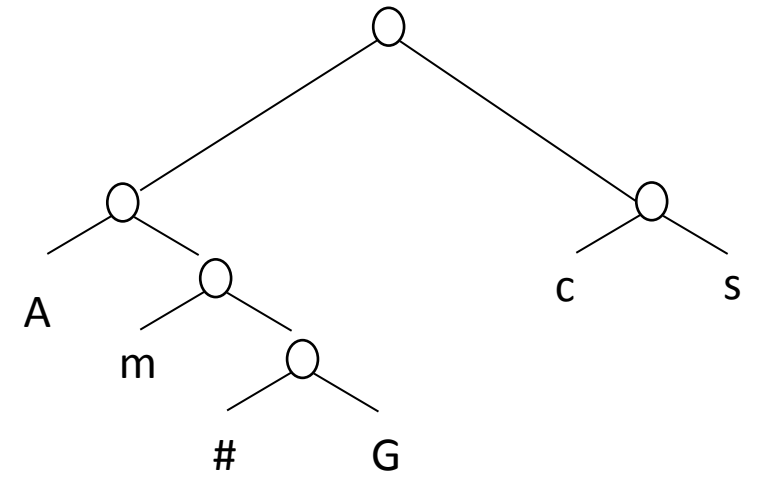
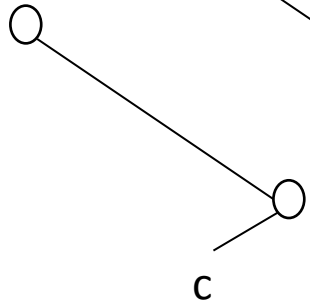
00 ⇒ A

|       |    |    |                       |
|-------|----|----|-----------------------|
| input | 00 | 00 | 101101100100111101100 |
|       | A  | A  |                       |

1

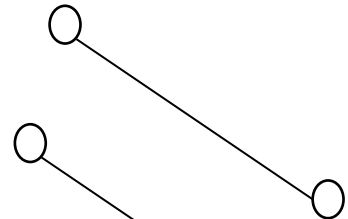


10 ⇒ c



|       |    |    |    |                     |
|-------|----|----|----|---------------------|
| input | 00 | 00 | 10 | 1101100100111101100 |
|       | A  | A  | c  |                     |

1



11 ⇒ s





|       |     |     |     |     |         |                           |
|-------|-----|-----|-----|-----|---------|---------------------------|
| input | 0 0 | 0 0 | 1 0 | 1 1 | 0 1 1 0 | 0 1 0 0 1 1 1 1 0 1 1 0 0 |
|       | A   | A   | c   | s   | #       |                           |

- 0 goes to the left
- 1 goes to the right
- If reach a leaf node, output the character
- go back to the root
  
- Characters are stored in only leaf nodes (by construction)

# Example

- Hellooo



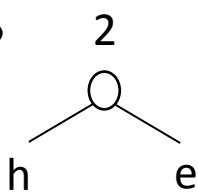
# Example

- Hellooo
- H (1) e (1) l(2) o (3)

# Example

- Hellooo

- H (1) e (1) l(2) o (3)

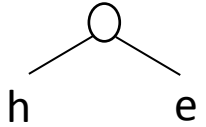
-  l(2) o (3)

# Example

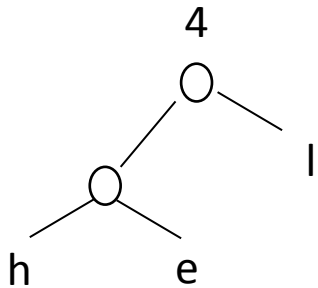
- Hellooo

- H (1) e (1) l(2) o (3)

- 2 l(2) o (3)



- 4 o (3)

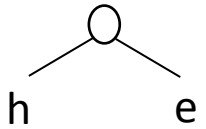


# Example

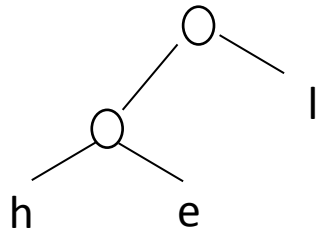
- Hellooo

- H (1) e (1) l(2) o (3)

- 2 l(2) o (3)



- o (3) 4 SHOULD BE IN ASCENDING ORDER

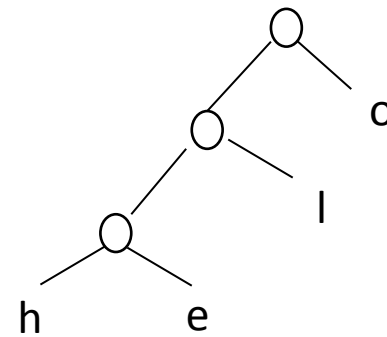
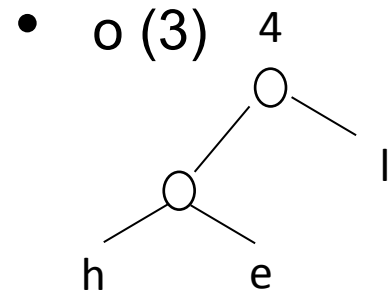
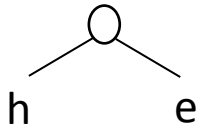


# Example

- Hellooo

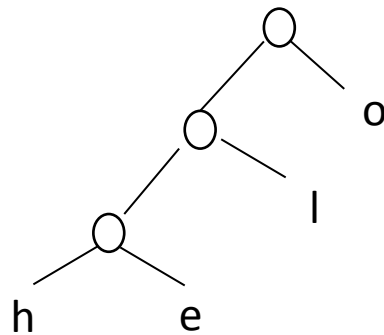
- H (1) e (1) l (2) o (3)

- 2 l (2) o (3)



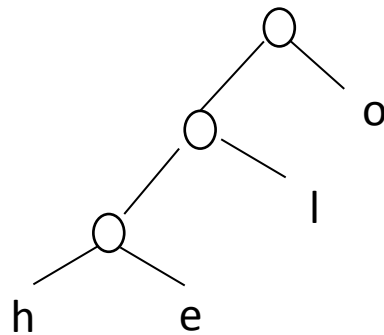
# Computing Codes

| character | occurrence | code |  |  | length |
|-----------|------------|------|--|--|--------|
| h         | 1          |      |  |  |        |
| e         | 1          |      |  |  |        |
| l         | 2          |      |  |  |        |
| o         | 3          |      |  |  |        |



# Computing Codes

| character | occurrence | code |   |   | length |
|-----------|------------|------|---|---|--------|
| h         | 1          | 0    | 0 | 0 | 3      |
| e         | 1          | 0    | 0 | 1 | 3      |
| l         | 2          | 0    | 1 |   | 2      |
| o         | 3          | 1    |   |   | 1      |

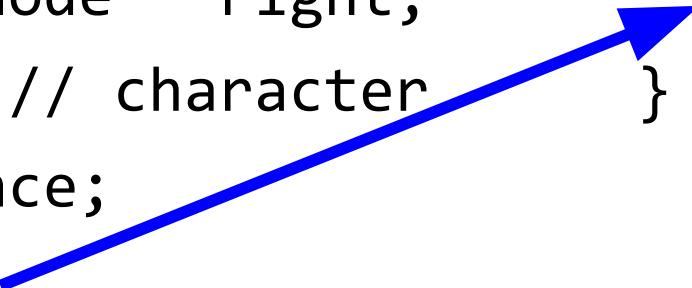


# How to build the compression tree

Ch 24 in <https://github.com/yunghsianglu/IntermediateCProgramming>

```
typedef struct treenode
{
    struct treenode * left;
    struct treenode * right;
    char value; // character
    int occurrence;
} TreeNode;

typedef struct listnode
{
    struct listnode * next;
    TreeNode * tnptr;
} ListNode;
```



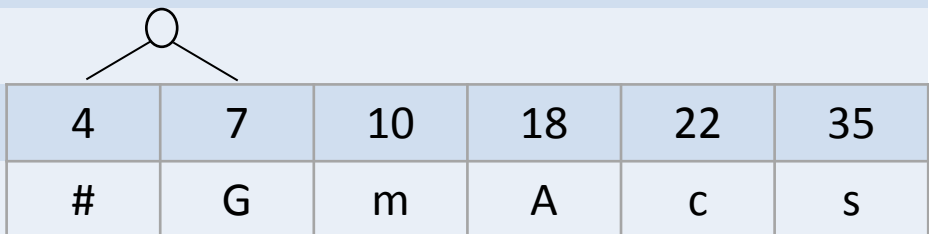


|            |   |    |   |    |    |    |
|------------|---|----|---|----|----|----|
| occurrence | 4 | 18 | 7 | 22 | 10 | 35 |
| letter     | # | A  | G | c  | m  | s  |

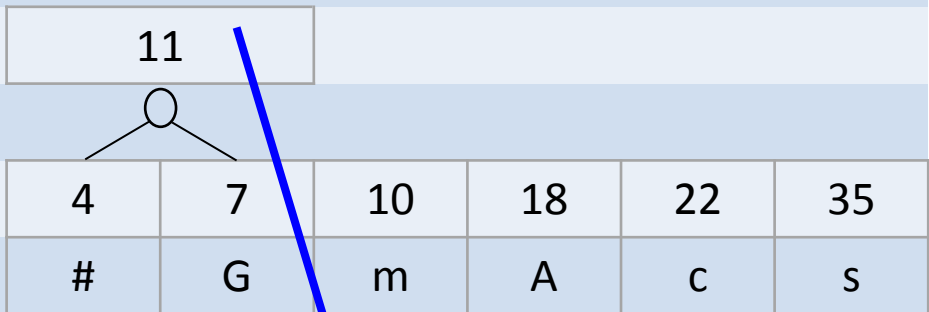
↓ Sort by the occurrences in ascending order

|   |   |    |    |    |    |
|---|---|----|----|----|----|
| 4 | 7 | 10 | 18 | 22 | 35 |
| # | G | m  | A  | c  | s  |

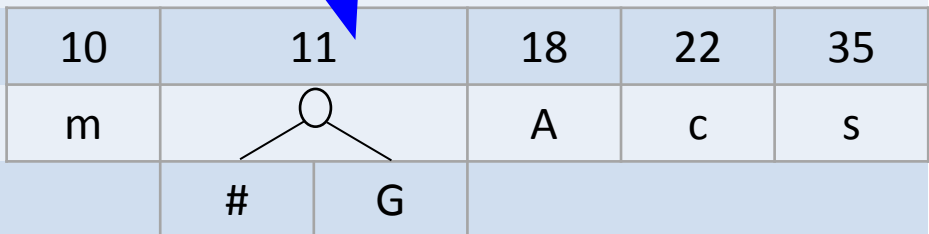
↓ Make the first two siblings of a binary tree



↓ The occurrence of the parent is the sum



↓ Insert the parent back in ascending order



□ list node

○ tree node

