

L15

4102

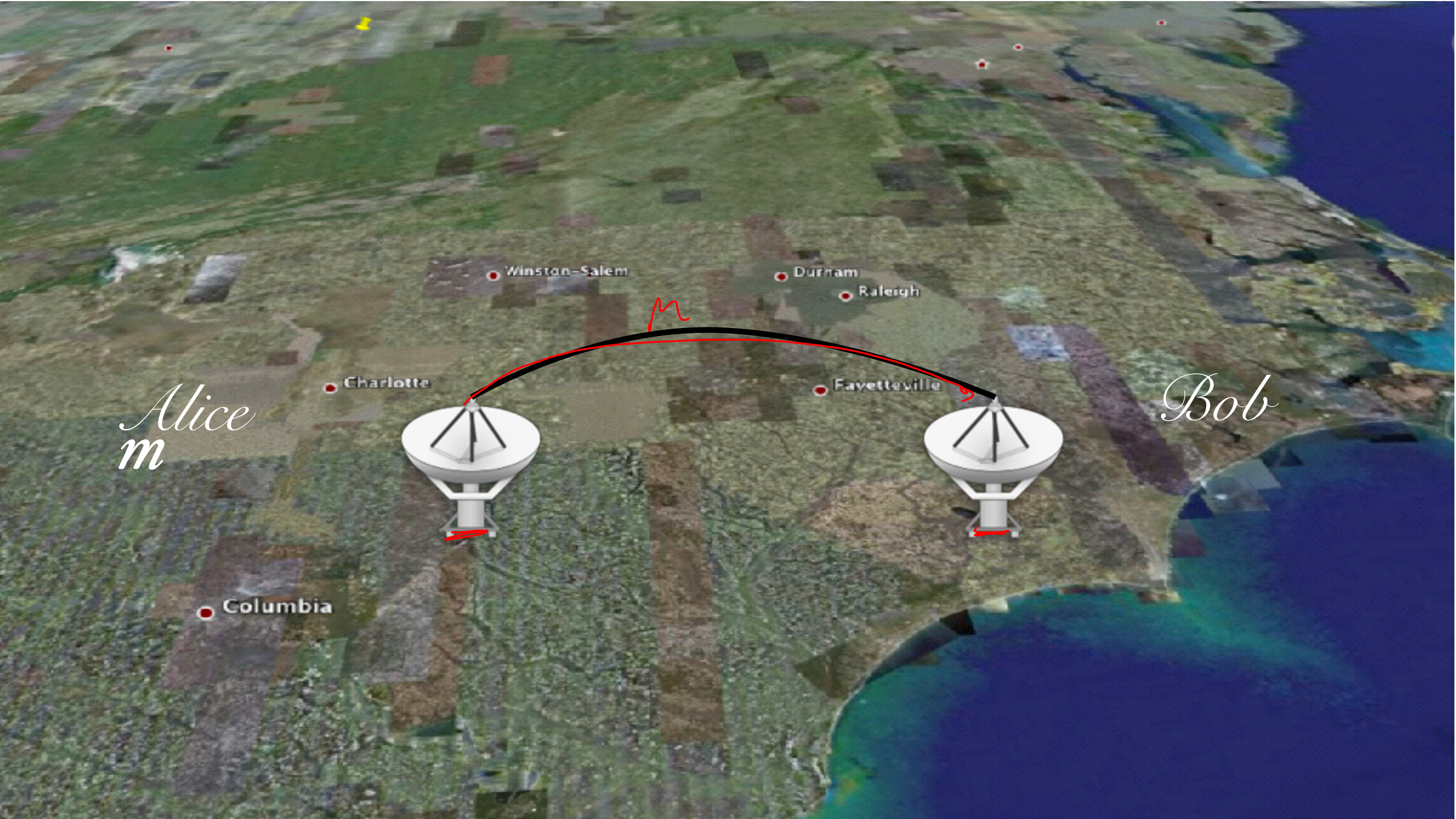
3.17.2016

abhi shelat

Huffman



SAM MORSE



Alice
m

Bob

Columbia

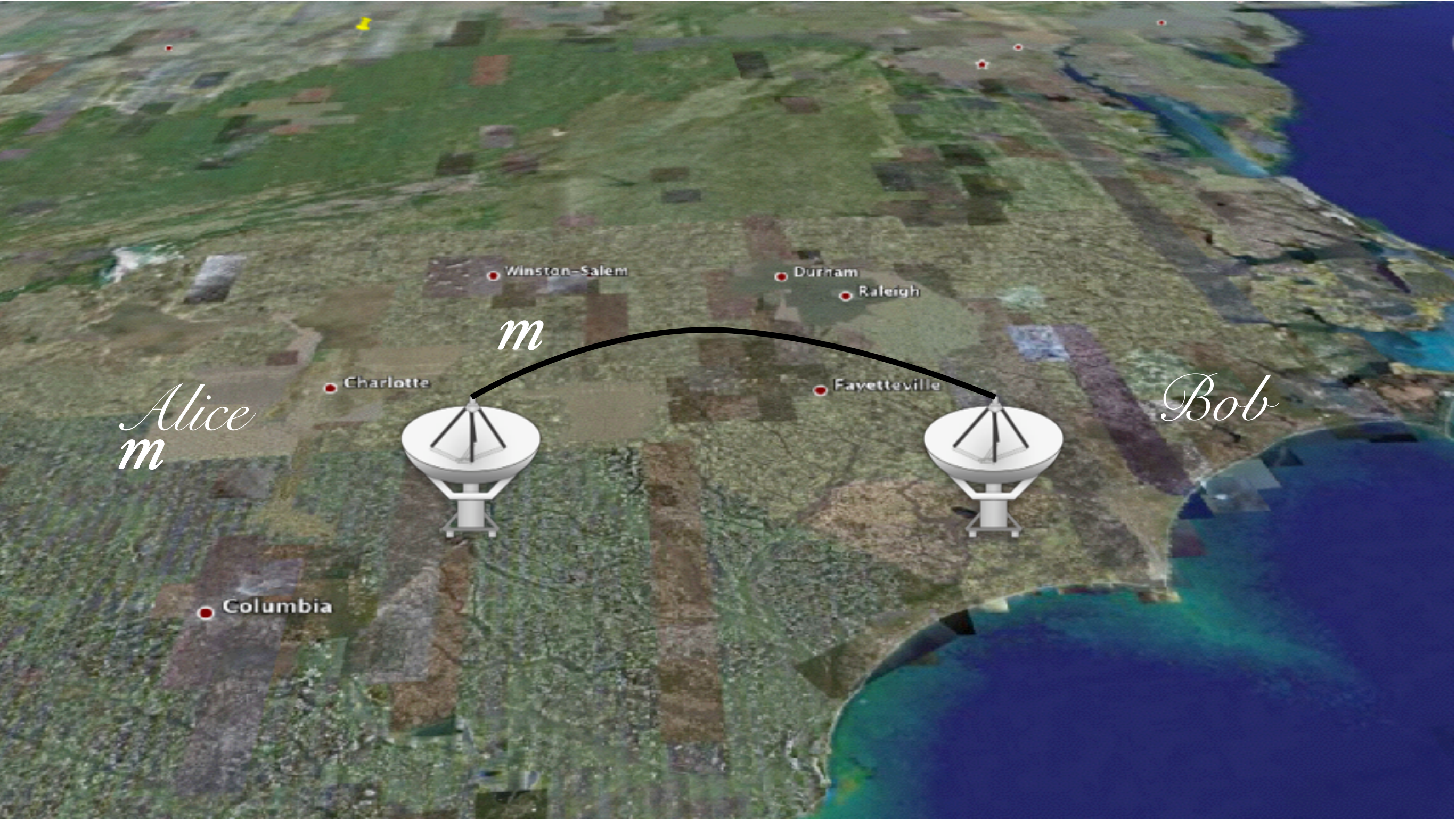
Charlotte

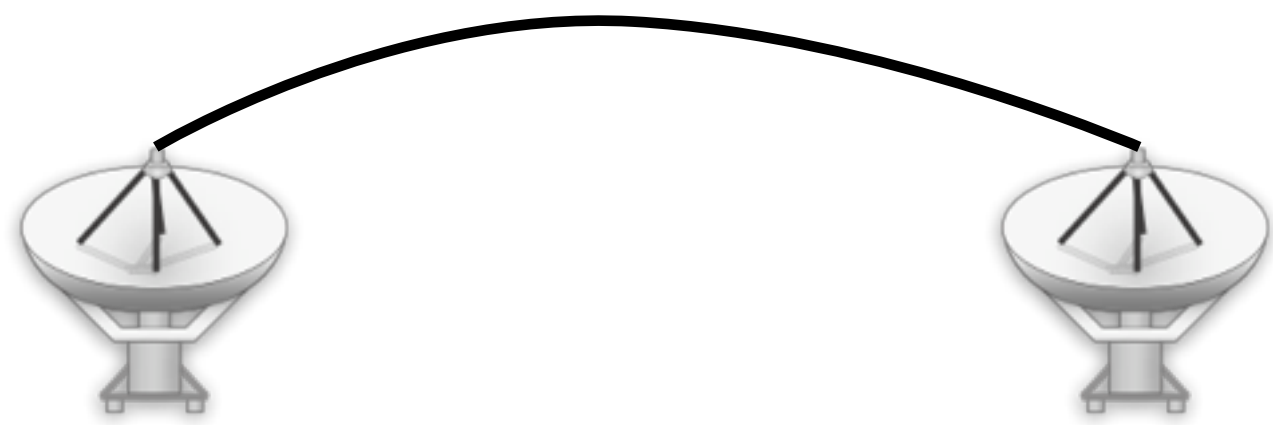
Winston-Salem

Durham

Raleigh

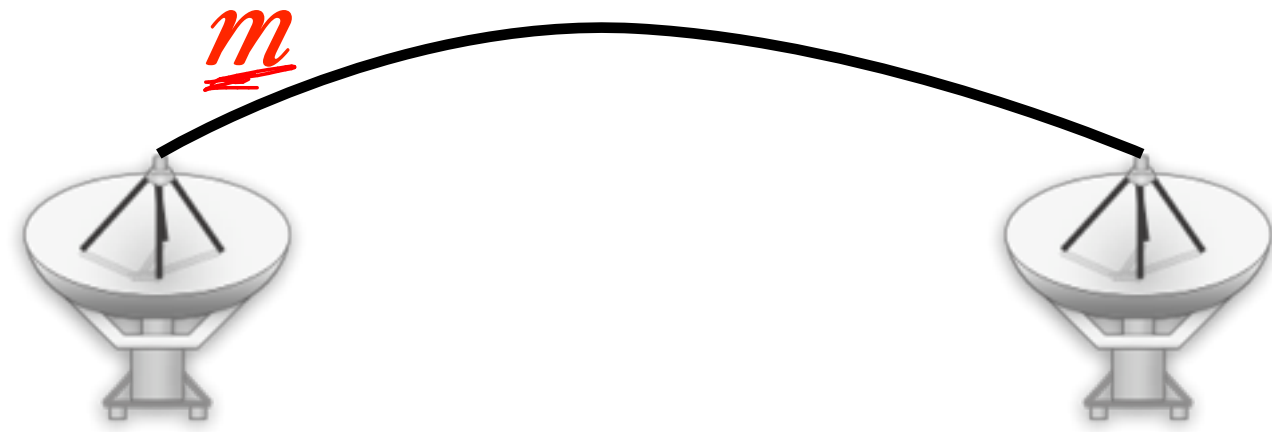
Fayetteville





MOSCOW — President Vladimir V. Putin's typically theatrical order to withdraw the bulk of Russian forces from Syria, a process that the Defense Ministry said it began on Tuesday, seemingly caught Washington, Damascus and everybody in between off guard — just the way the Russian leader likes it.

By all accounts, Mr. Putin delights at creating surprises, reinforcing Russia's newfound image as a sovereign, global heavyweight and keeping him at the center of



MOSCOW — President Vladimir V. Putin's typically theatrical order to withdraw the bulk of Russian forces from Syria, a process that the Defense Ministry said it began on Tuesday, seemingly caught Washington, Damascus and everybody in between off guard — just the way the Russian leader likes it.

alphabet with
60 characters
in it.

28.

By all accounts, Mr. Putin delights at creating surprises, reinforcing Russia's newfound image as a sovereign, global heavyweight and keeping him at the center of

$c \in C$	f_c	T
e	235	000
i	200	001
o	170	010
u	87	011
p	78	100
g	47	:
b	40	:
f	24	111

881

frequency

C

c ∈ C f_c

T

l_c

e:	<u>235</u>	<u>000</u>	<u>3</u>
i:	<u>200</u>	001	<u>3</u>
o:	170	010	3
u:	87	011	3
p:	78	100	3
g:	47	101	3
b:	40	110	3
f:	<u>24</u>	<u>111</u>	3

881

.

3

=

2643

=

def: cost of an encoding

encoding

frequency for each $c \in C$

$$\underline{B(T, \{f_c\})} = \sum_{c \in C} \underline{f_c} \cdot \underline{l_c}$$

table

$c \in C$	f_c	T	l_c
<u>e</u> : 235		<u>000</u>	3
i: 200		<u>001</u>	3
o: 170		<u>010</u>	3
u: 87		011	3
p: 78		100	3
g: 47		101	3
b: 40		110	3
f: 24		111	3

$$235 \cdot 3 + 200 \cdot 3 + 170 \cdot 3 + \dots +$$

$$24 \cdot 3 = 881 \cdot 3 = 2643$$

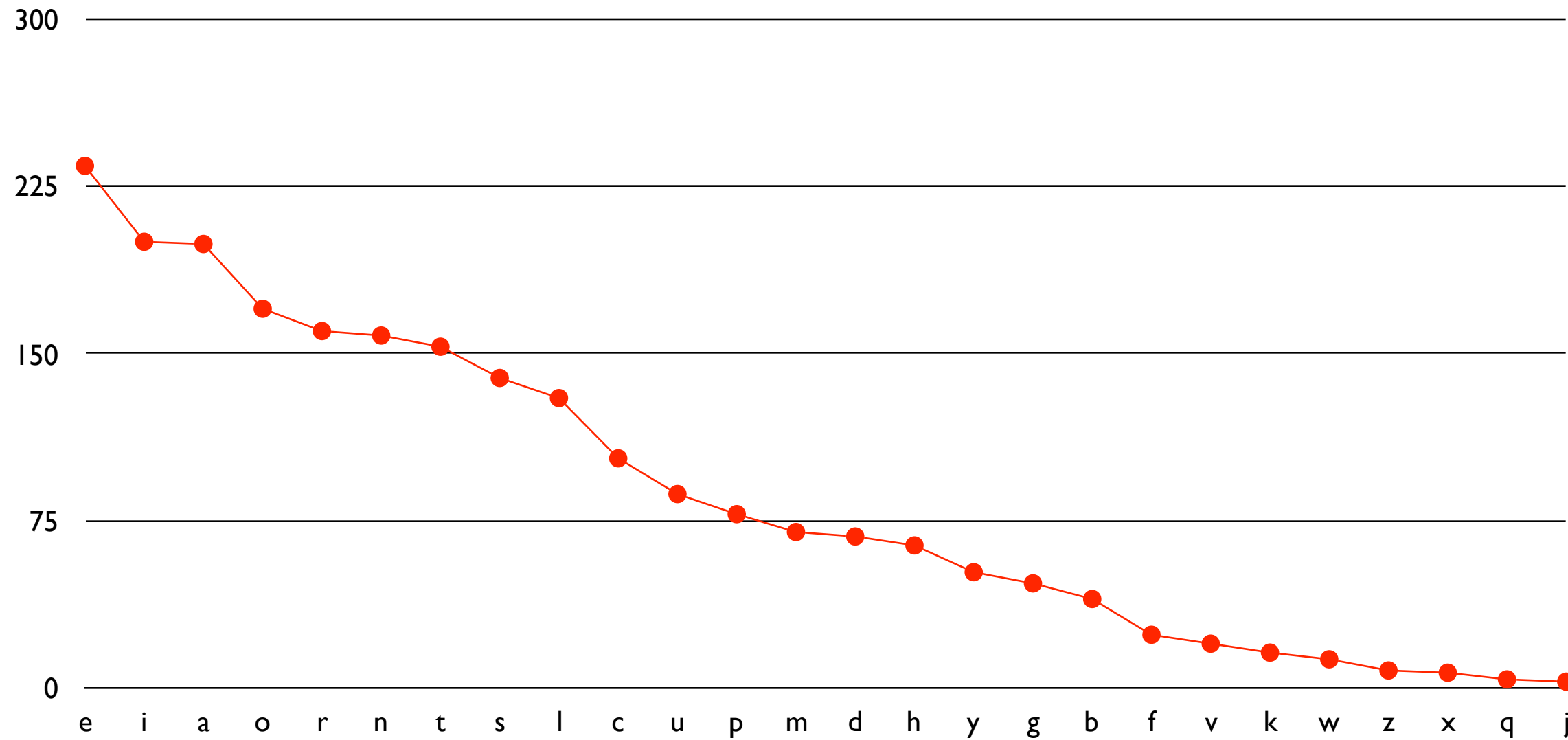
881

$$\cdot 3 = 881 \cdot 3$$

character frequency

e: 234803
i: 200613
a: 198938
o: 170392
r: 160491
n: 158281
t: 152570
s: 139238
l: 130172
c: 103307
u: 87211
p: 78077
m: 70504
d: 68007
h: 64165
y: 51527
g: 47011
b: 40351
f: 24110
v: 20103
k: 16012
w: 13825
z: 8439
x: 6926
q: 3729
j: 3075

$\{f_c\}$



e

q

morse code

International Morse Code

- 1 dash = 3 dots.
- The space between parts of the same letter = 1 dot.
- The space between letters = 3 dots.
- The space between words = 7 dots.

A	• —	V	• • • —
B	— • • •	W	— — • •
C	— • — •	X	— • • —
D	— • • •	Y	— • — —
E	•	Z	— — • •
F	• • — •	.	• — • — • —
G	— — •	,	— — • • — —
H	• • • •	?	• • — — • •
I	• •	/	— • • — •
J	• — — —	@	• — — • • •
K	— • —	1	• — — — —
L	• — • •	2	• • — — —
M	— —	3	• • • — —
N	— •	4	• • • • —
O	— — —	5	• • • • •
P	• — — •	6	— • • • •
Q	— — • —	7	— — • • •
R	• — • •	8	— — — • •
S	• • •	9	— — — — •
T	—	0	— — — — —
U	• • —		

image http://en.wikipedia.org/wiki/Morse_code

code

decide to send several messages

— this is a problem

A A

E T E T

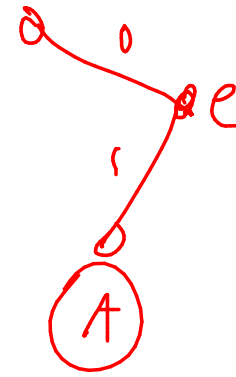
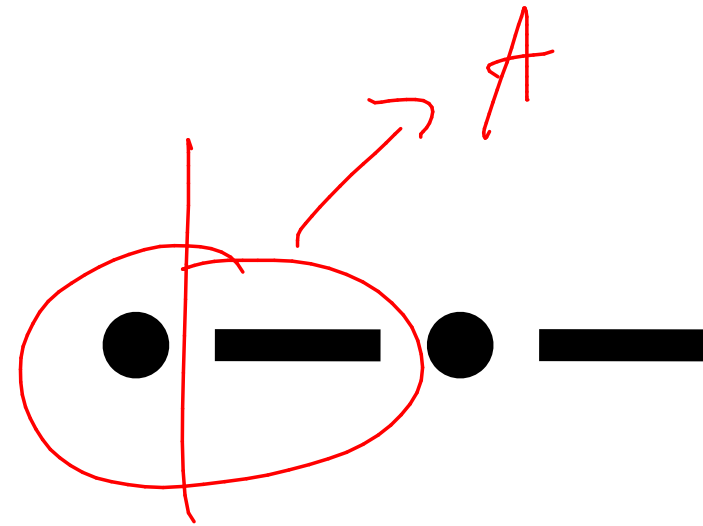
R I

Morse code

International Morse Code

- 1 dash = 3 dots.
- The space between parts of the same letter = 1 dot.
- The space between letters = 3 dots.
- The space between words = 7 dots.

A	• —	V	• • • —
B	— • • •	W	— — • •
C	— • — •	X	— • • —
D	— • • •	Y	— • — —
E	•	Z	— — • •
F	• • — •	.	• — • — • —
G	— — •	,	— — • • — —
H	• • • •	?	• • — • • •
I	• •	/	— • • — •
J	• — — —	@	• — — • • •
K	— • —	1	• — — — —
L	• — • •	2	• • — — —
M	— —	3	• • • — —
N	— •	4	• • • • —
O	— — —	5	• • • • •
P	• — — •	6	— • • • •
Q	— — • —	7	— — • • •
R	• — • •	8	— — — • •
S	• • •	9	— — — — •
T	—	0	— — — — —
U	• • —		



encoding e - @ IS A

PREFIX of the encoding

for A.

def: prefix-free code

Code^{for} C such that for any two $x, y \in C$

if $x \neq y$ then code(x) is NOT A prefix of
code(y).

def: prefix-free code

$\forall x, y \in C, x \neq y \implies \text{CODE}(x)$ not a prefix of $\text{CODE}(y)$

def: prefix code

$\forall x, y \in C, x \neq y \implies \text{CODE}(x) \text{ not a prefix of } \text{CODE}(y)$

e:	235	0
i:	200	10
o:	170	110
u:	87	1110
p:	78	11110
g:	47	111110
b:	40	1111110
f:	24	11111110

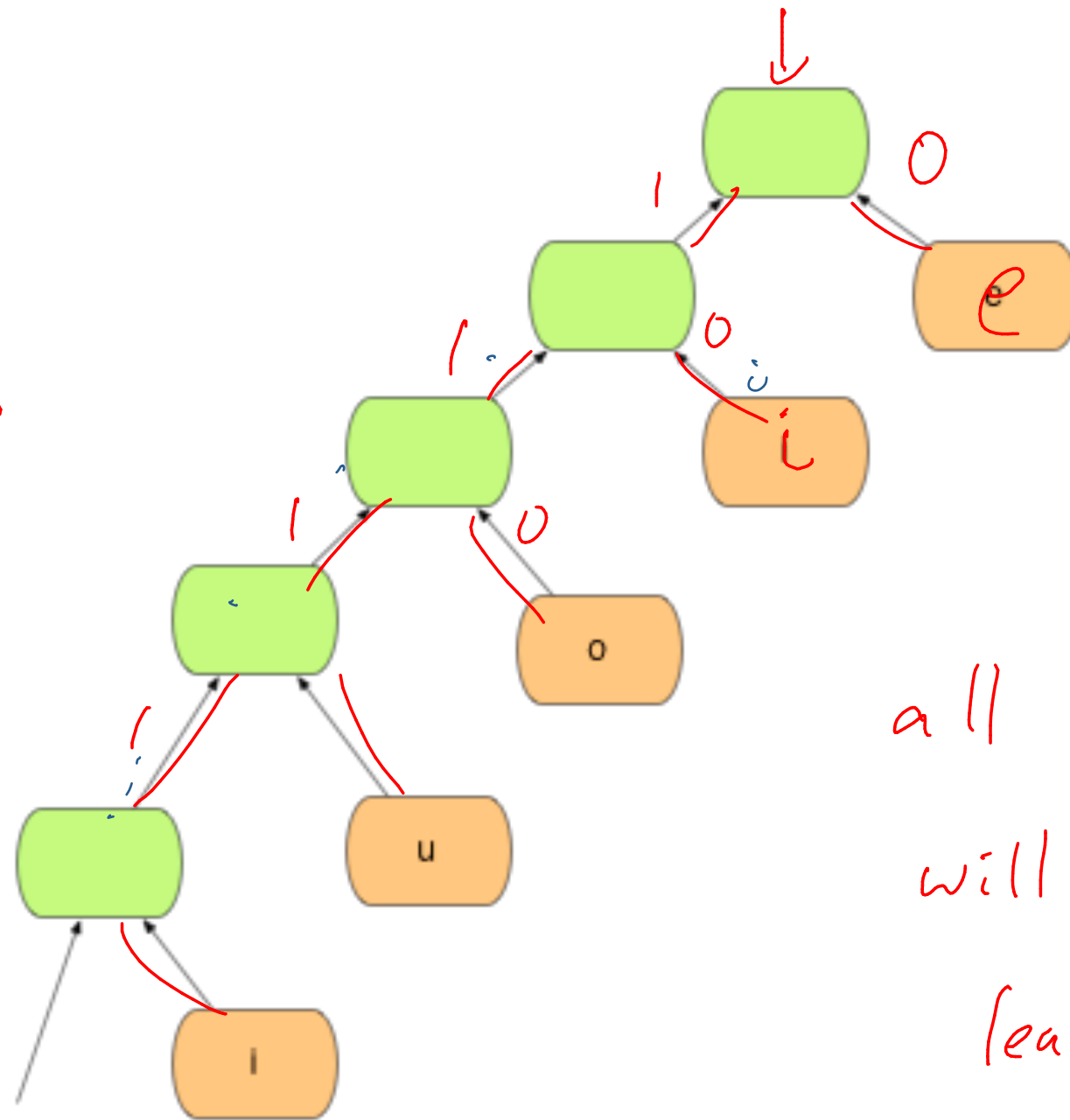
decoding a prefix code

e:	235	0
i:	200	10
o:	170	110
u:	87	1110
p:	78	11110
g:	47	111110
b:	40	<u>1111110</u>
f:	24	11111110

111111010111110
b i j

code to binary tree

e: 235	<u>0</u>
i: 200	<u>10</u>
o: 170	110
u: 87	1110
p: 78	11110
g: 47	111110
b: 40	1111110
f: 24	11111110

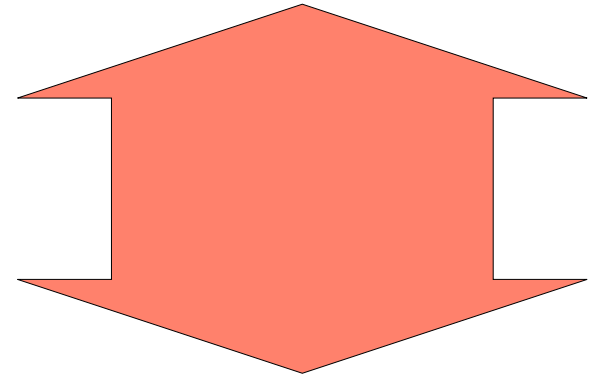


all symbols in C
will be unique
leaves in the tree,

111111010111110

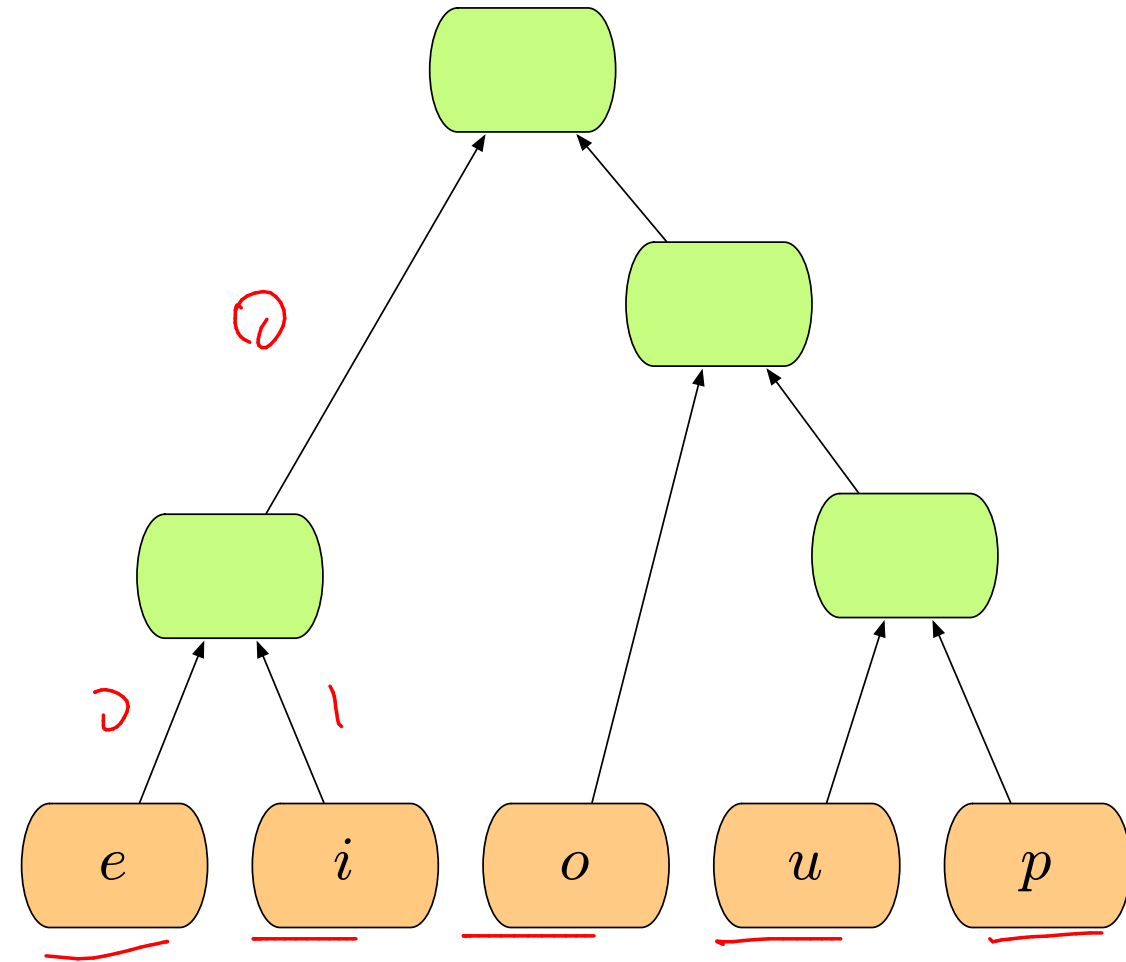
prefix code

free



binary tree

use tree to encode



$c \in C$	f_c	T	l_c
e	235	00	2
i	200	01	2
o	170	10	2
u	87	110	3
p	78	111	3

Goal

input
given the

frequencies

$\{f_c\}$

for an alphabet C ,

find the optimal prefix-free code T

that

minimizes

$B(T, \{f_c\})$

output

Goal

given the character frequencies $\{f_c\}_{c \in C}$

produce a prefix code T with smallest cost

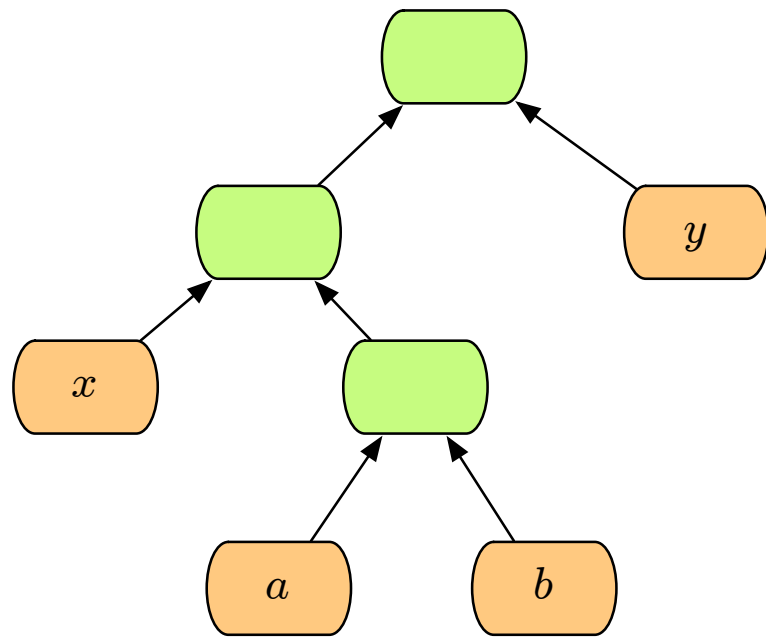
$$\min_T B(T, \{f_c\})$$

Property

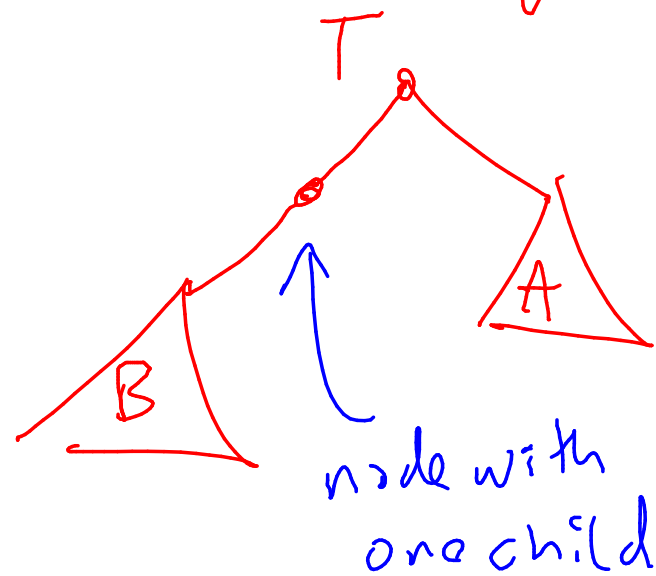
* assume all frequencies > 0 .

Lemma: Optimal tree is full.

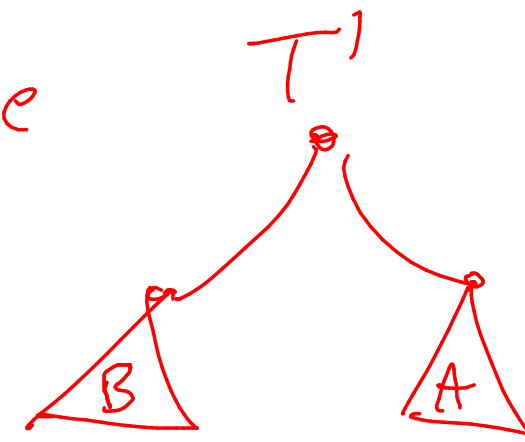
Every node in the optimal tree has either 0 or 2 children.



Proof: Suppose optimal T has a node with $\neq 2$ child



Consider tree



claim is that

$$B(T) > B(T')$$

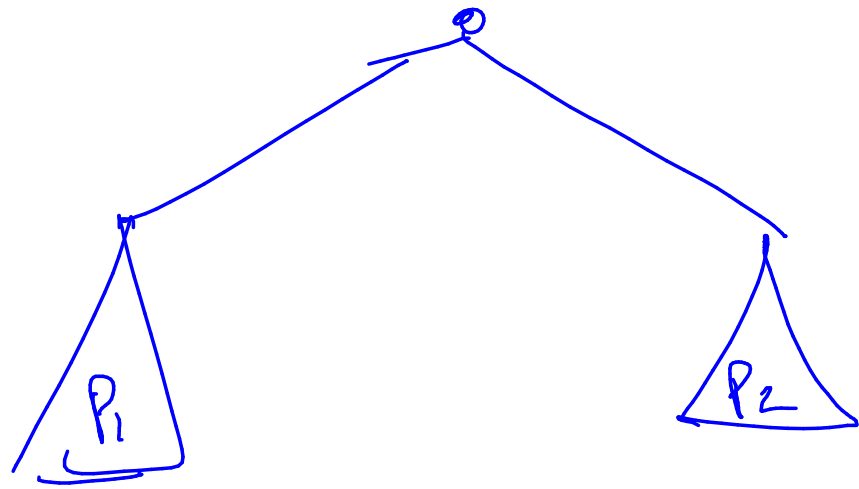
because all symbols in subtree B have

encodings that are 1 bit longer than in T' ,

Divide & conquer?

Partition the frequencies into 2 sets

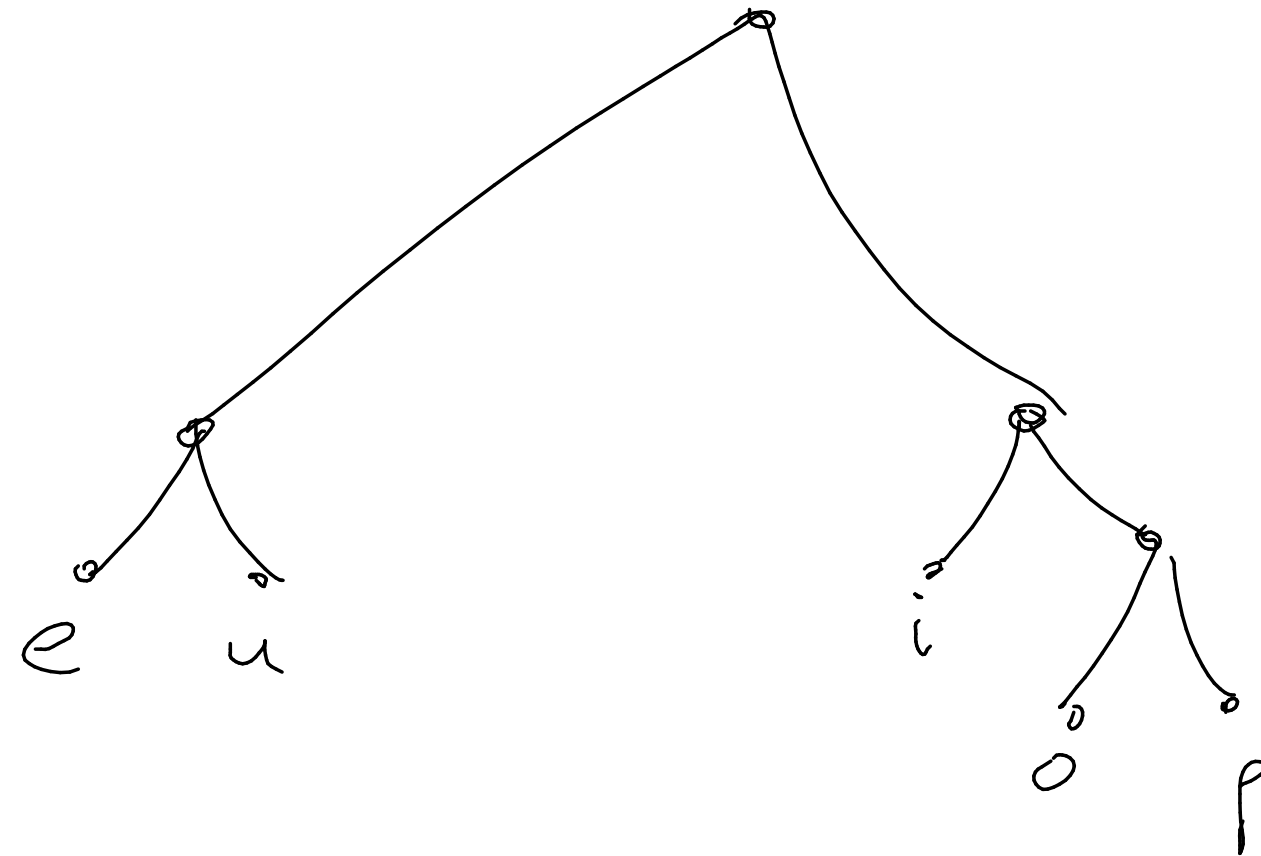
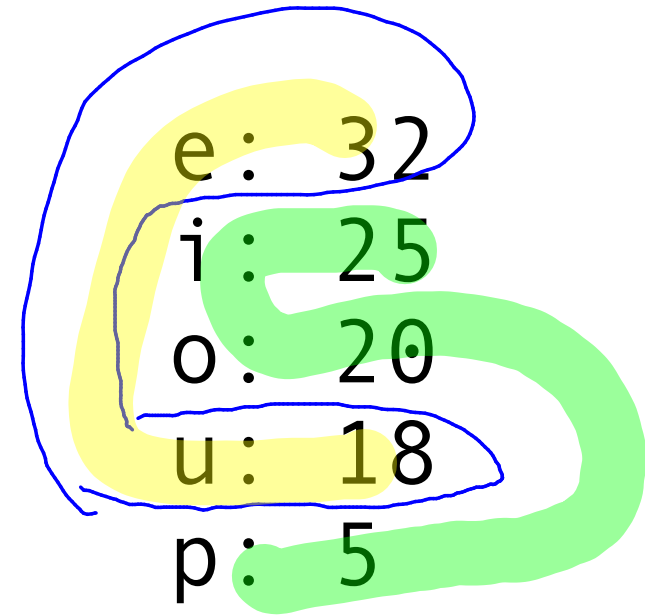
P_1 P_2 , Build trees for these & combine.



counter-example

$e, u = 50$

$i, o, p = 50$



$e = 00$

$u = 01$

$i = 10$

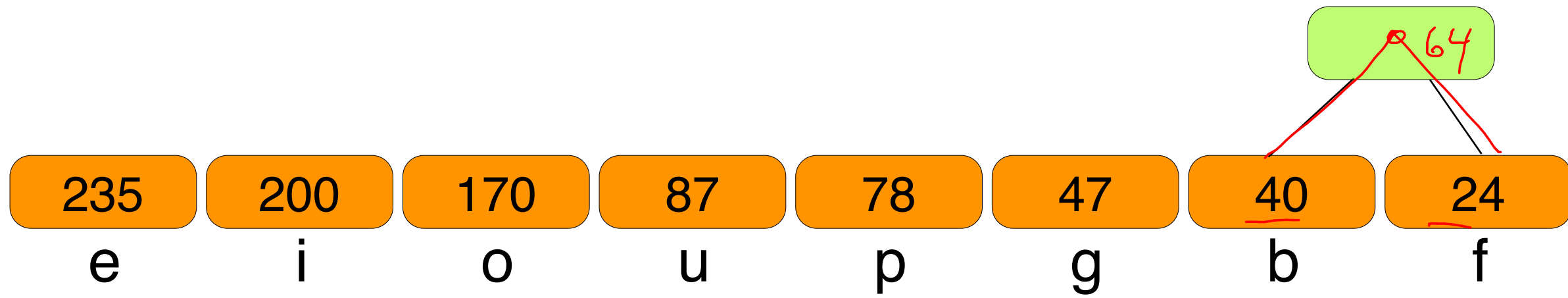
$o = 110$

$p = 111$

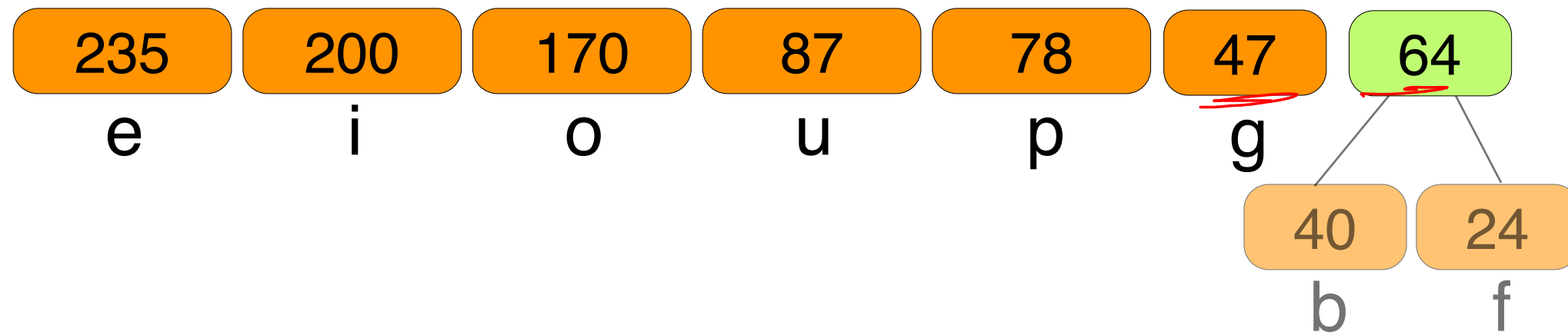
Swapping the codes for o and u results

in a better code b/c

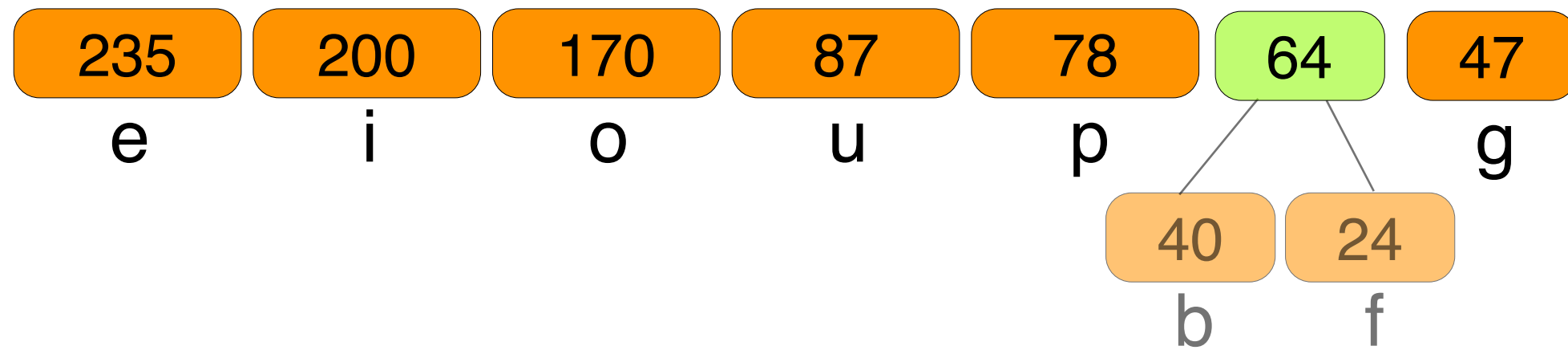
$f_o > f_u$

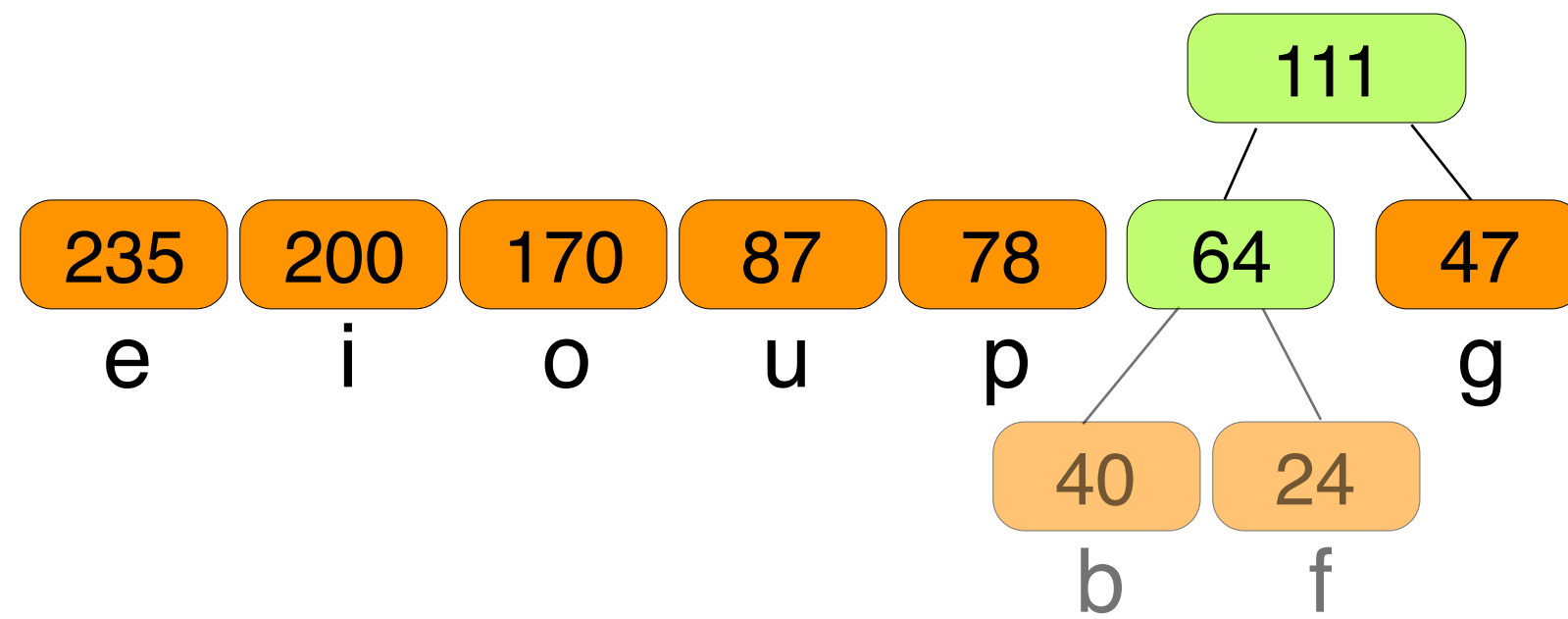


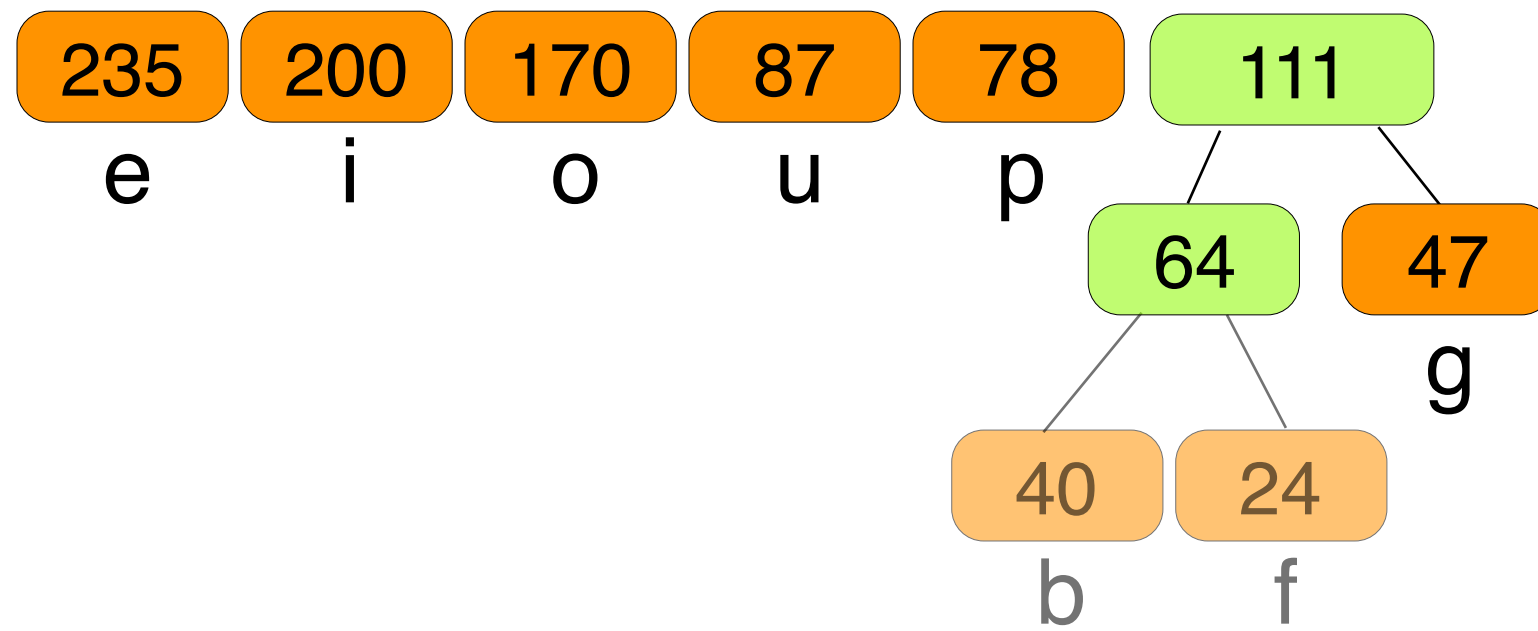
n characters
in C

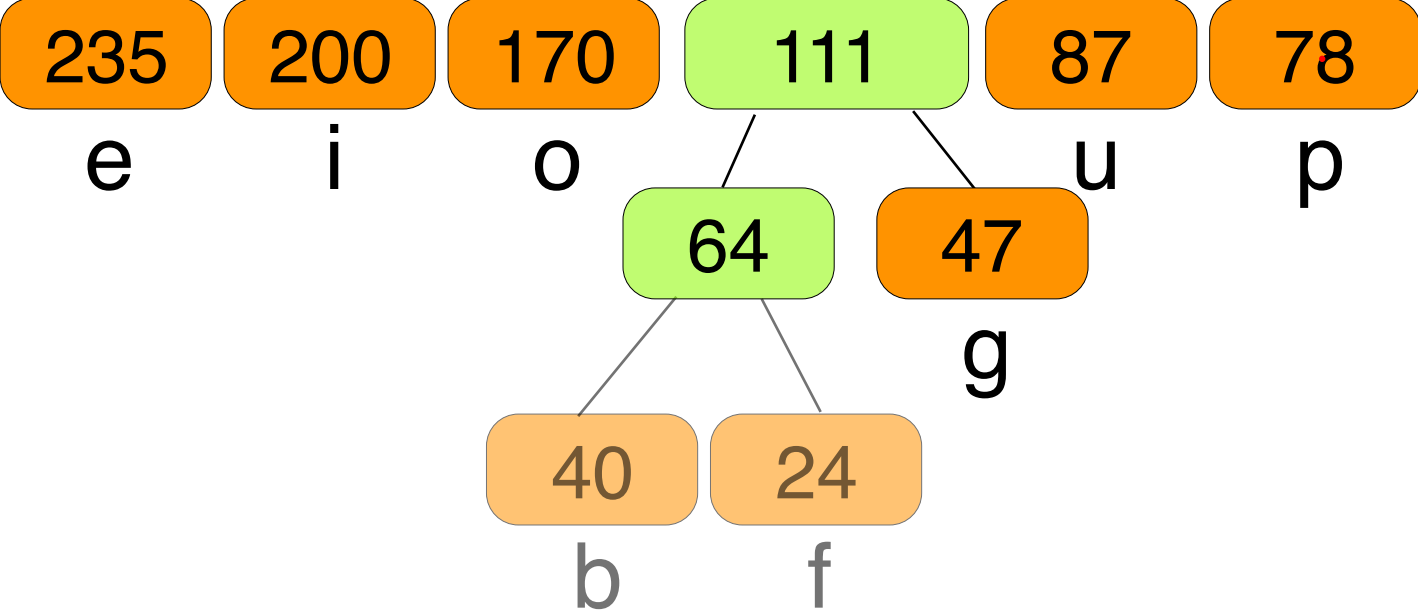


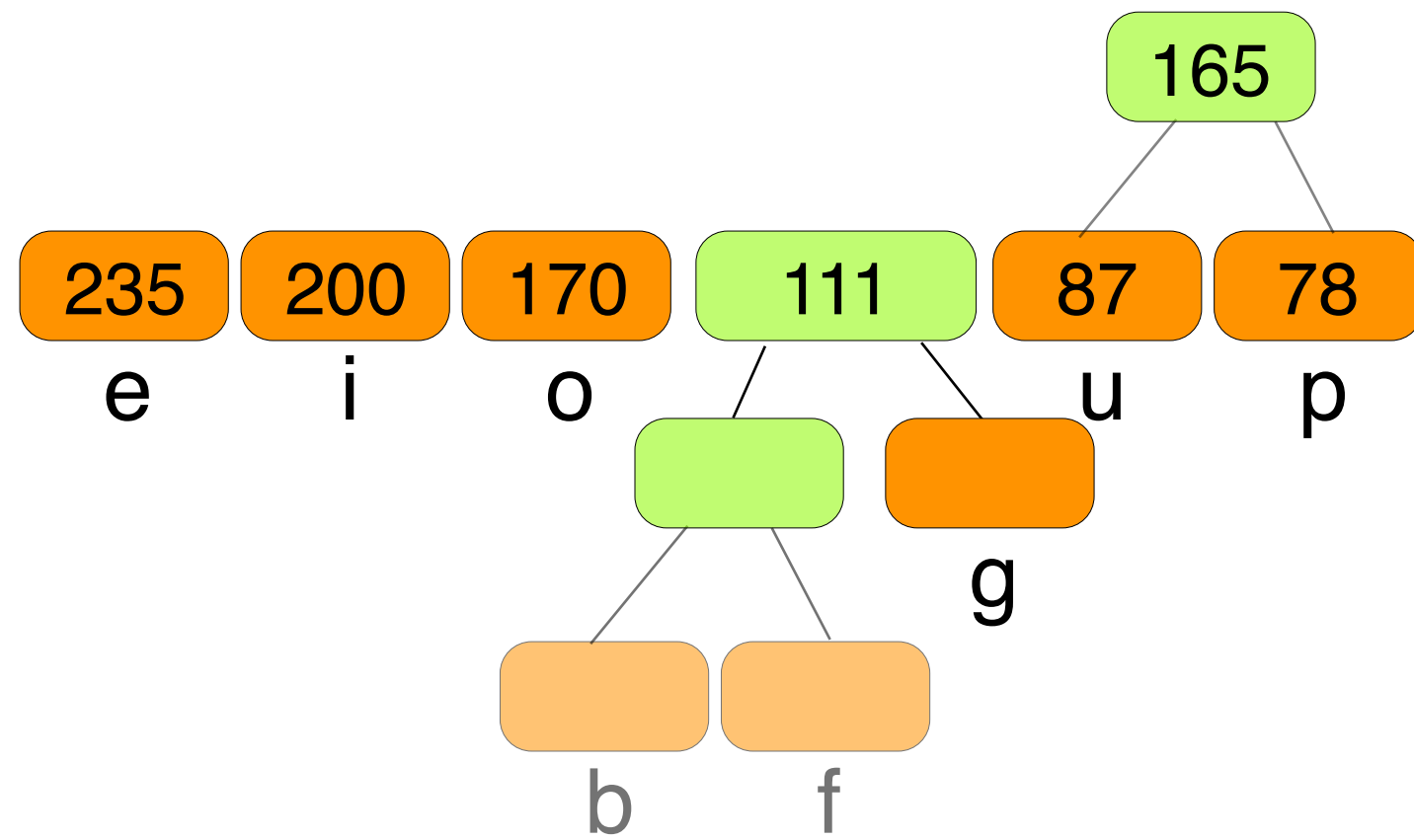
} n-1 characters
in c

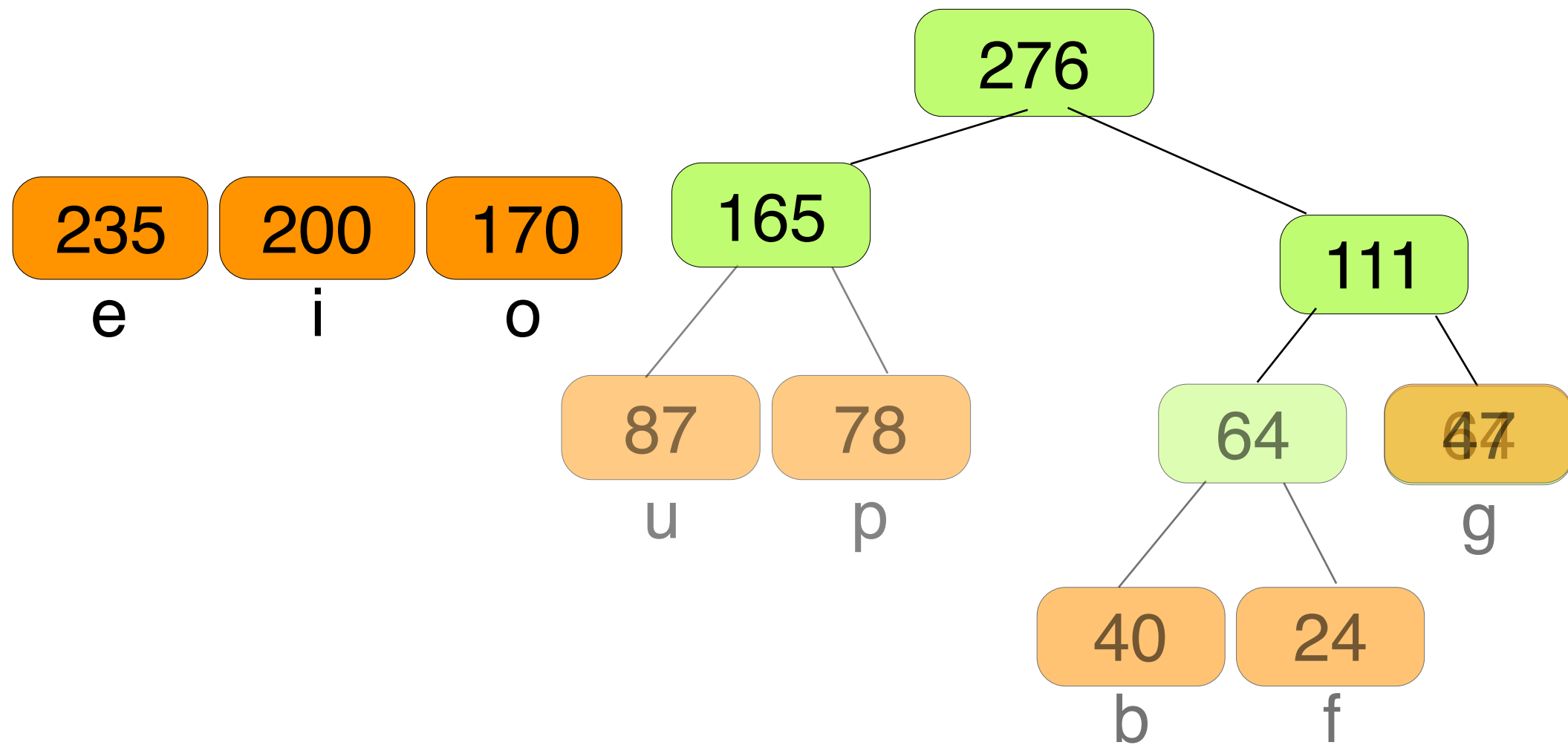


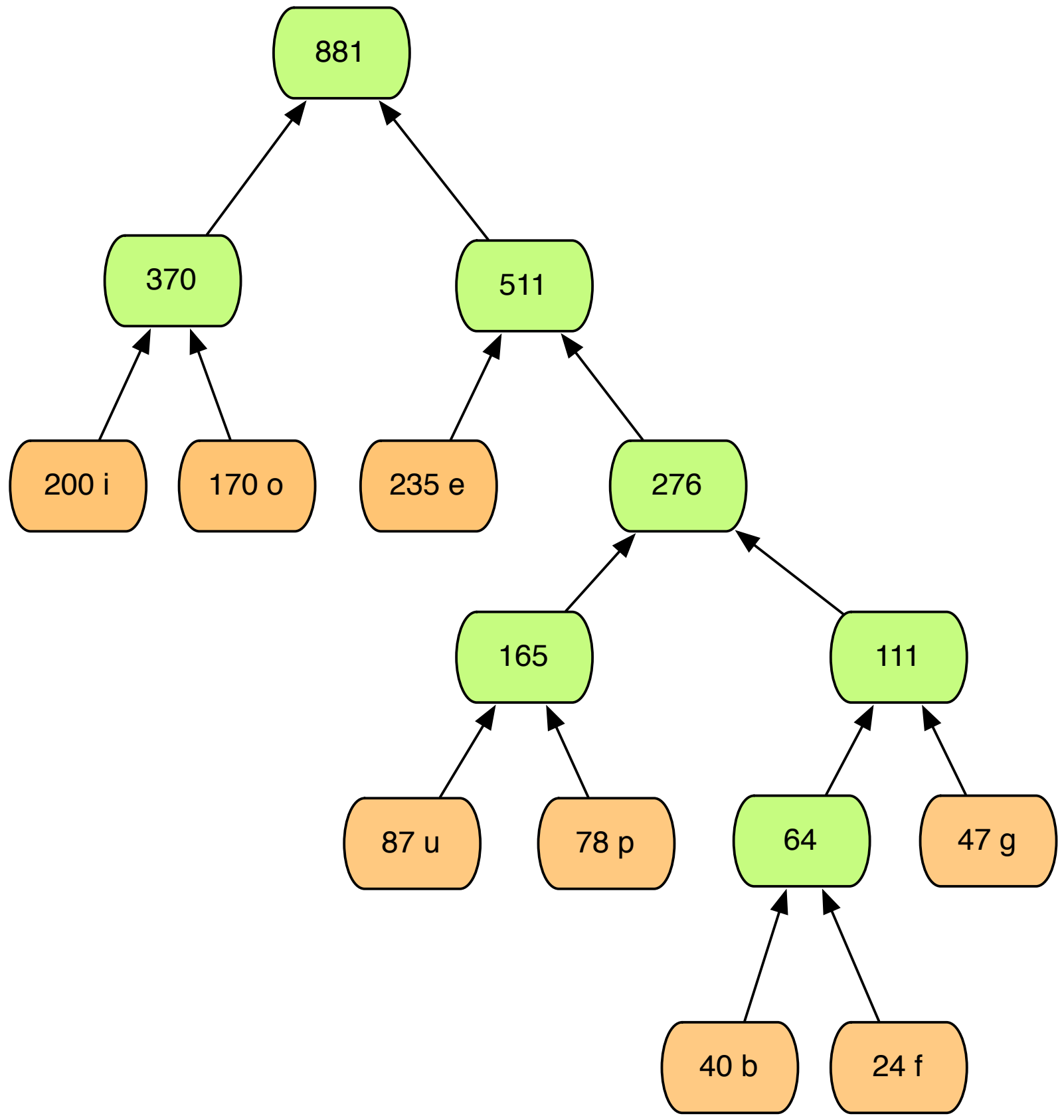


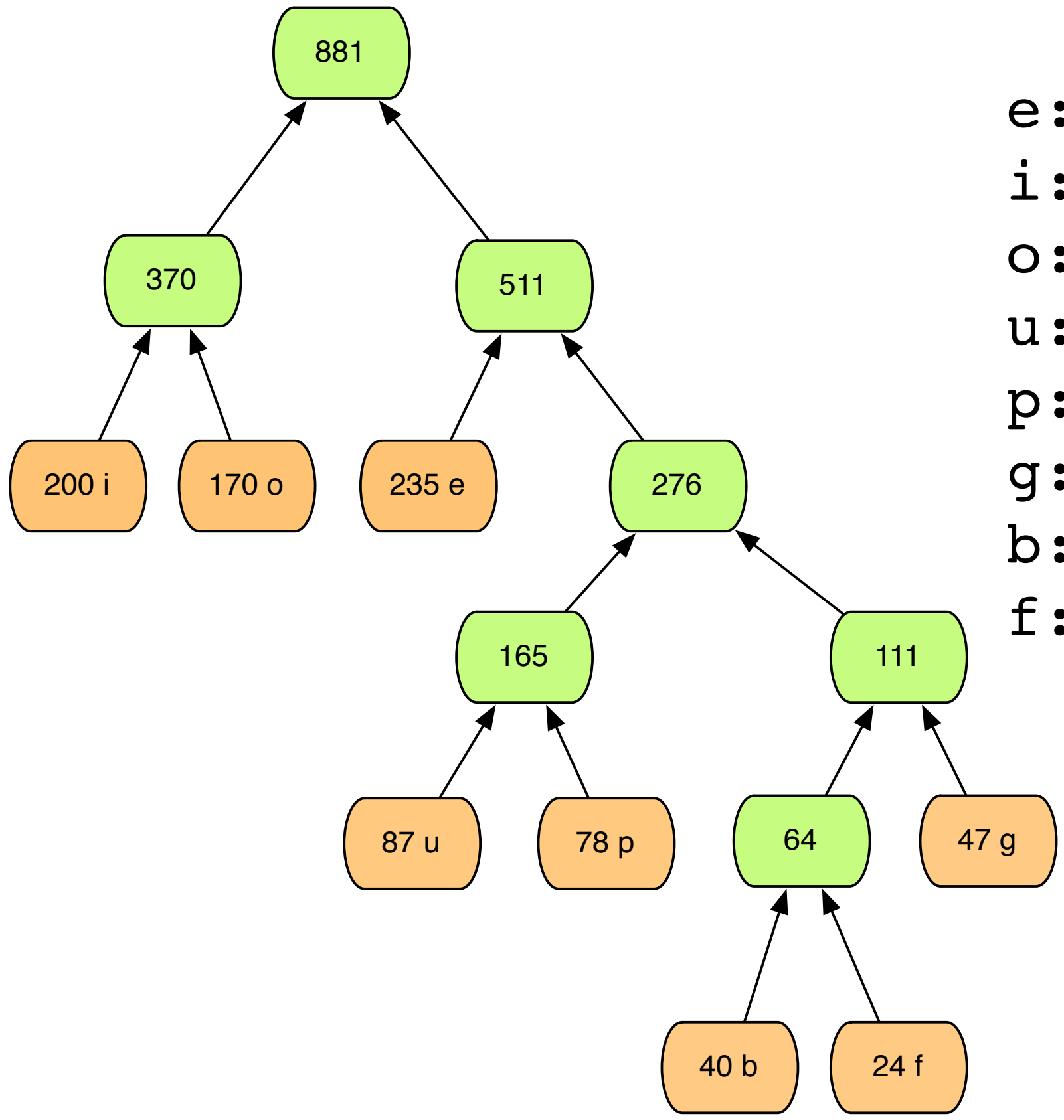




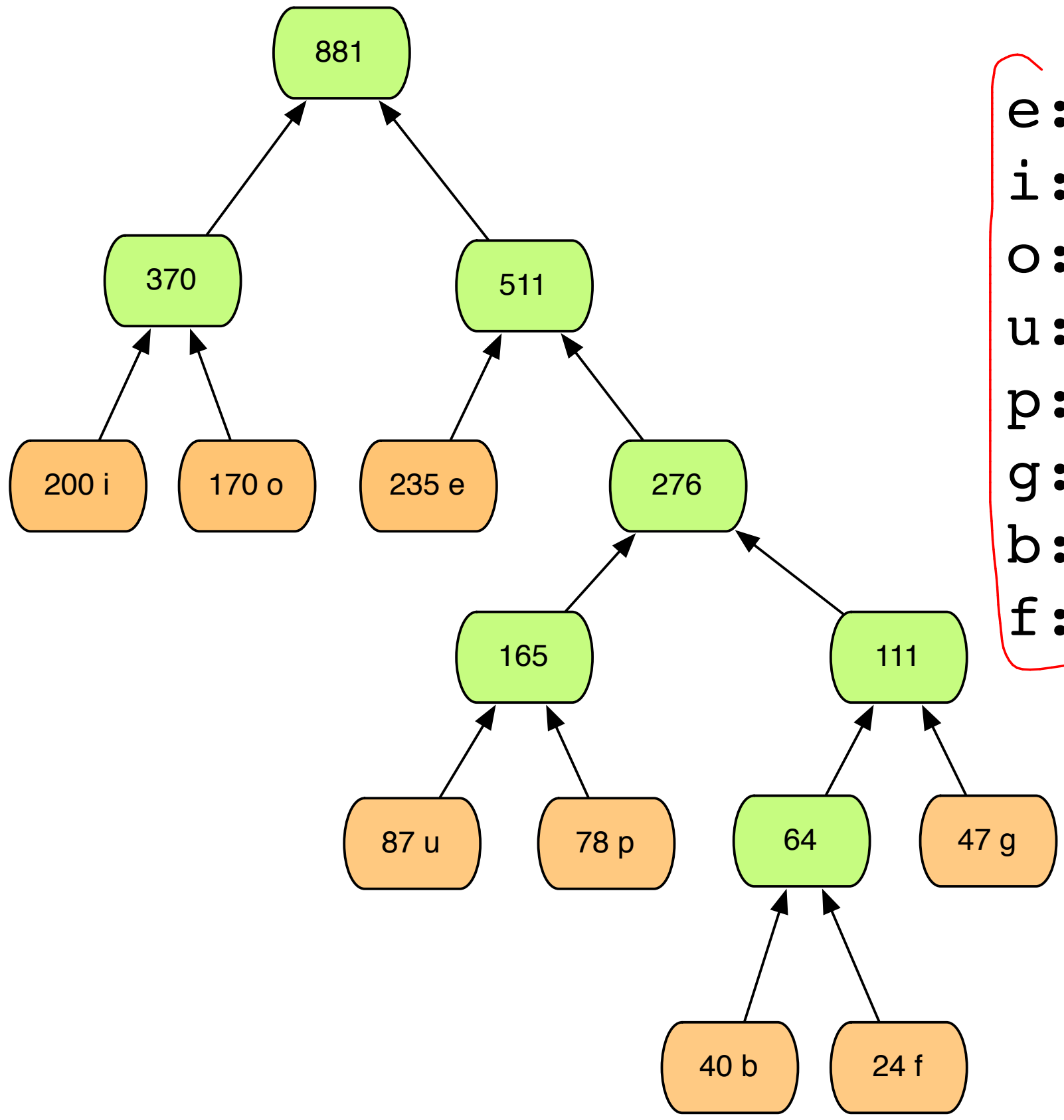








e: 235 01
 i: 200 11
 o: 170 10
 u: 87 0011
 p: 78 0010
 g: 47 0000
 b: 40 00011
 f: 24 00010



e:	235	01	470
i:	200	11	400
o:	170	10	340
u:	87	0011	312
p:	78	0010	188
g:	47	0000	200
b:	40	00011	120
f:	24	00010	2378

compare with slide I,
 2378 vs 2643.

objective

Given $\{f_c\}$, we can compute a prefix-free code using Huffman's algorithm.

Prove the resulting tree is optimal.

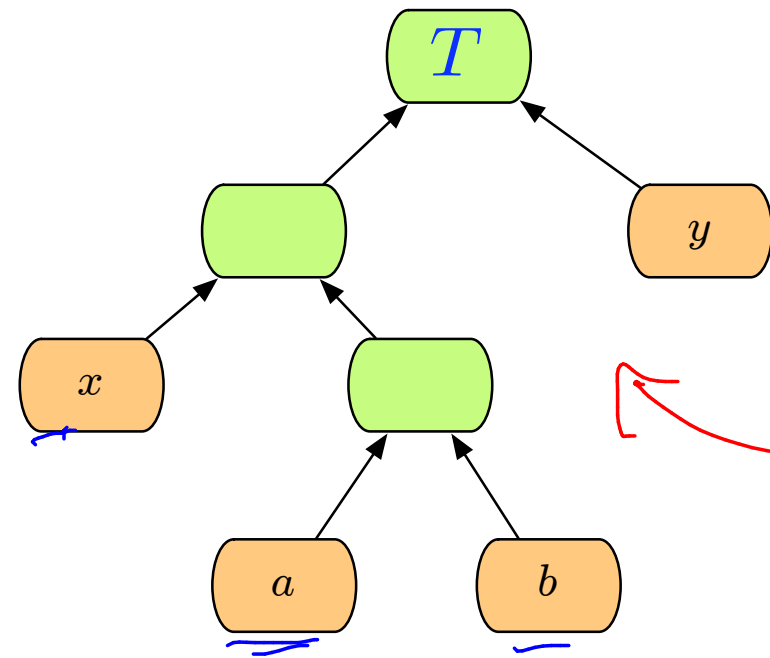
Exchange argument

Lemma: Let f_x and f_y be the symbols with the smallest frequencies.
There exists an optimal tree in which x and y are siblings.

Exchange argument

LEMMA:

Let $x, y \in C$ be characters with smallest frequencies f_x, f_y . There exists an optimal prefix code T'' for C in which x, y are siblings. That is, the codes for x, y have the same length and only differ in the last bit.



Proof: Let T be an opt. p.f. code for C .

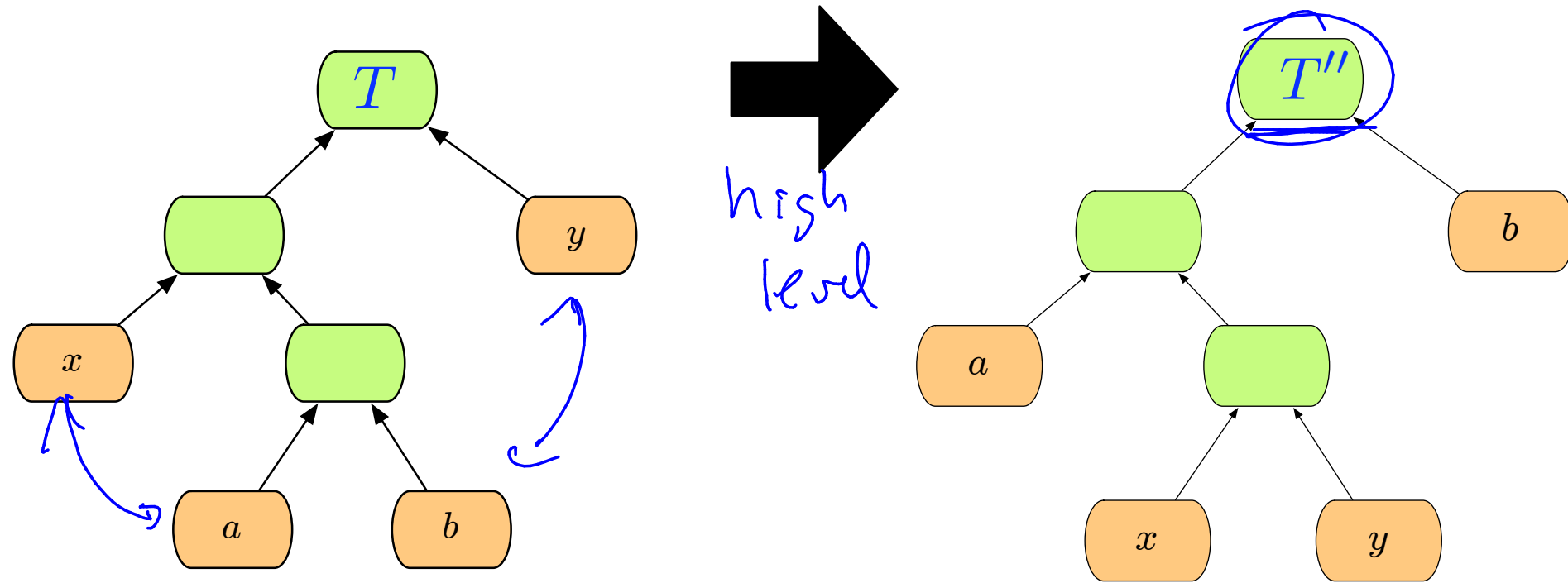
If x, y are siblings, then the lemma holds.

Since they are not, let a, b be the 2 characters of lowest depth that are siblings. *

Why do a, b exist ?? ANS: b/c the tree is full.

LEMMA:

Let $x, y \in C$ be characters with smallest frequencies f_x, f_y . There exists an optimal prefix code T'' for C in which x, y are siblings. That is, the codes for x, y have the same length and only differ in the last bit.



wlog f_x, f_y have the smallest frequencies.
 $f_x \leq f_a$ $f_y \leq f_b$



$$B(T) = C + f_x \cdot l_x + f_a \cdot l_a$$

$$B(T') = C + \underline{f_x \cdot l_a} + \underline{f_a \cdot l_x}$$

$$\underline{B(T) - B(T')} = f_x \cdot l_x + f_a \cdot l_a - f_x \cdot l_a - f_a \cdot l_x$$

$$= f_x(\underline{l_x - l_a}) + f_a(l_a - l_x) = \underbrace{(l_a - l_x)}_{\geq 0} \underbrace{(f_a - f_x)}_{\geq 0}$$

$$\geq 0$$

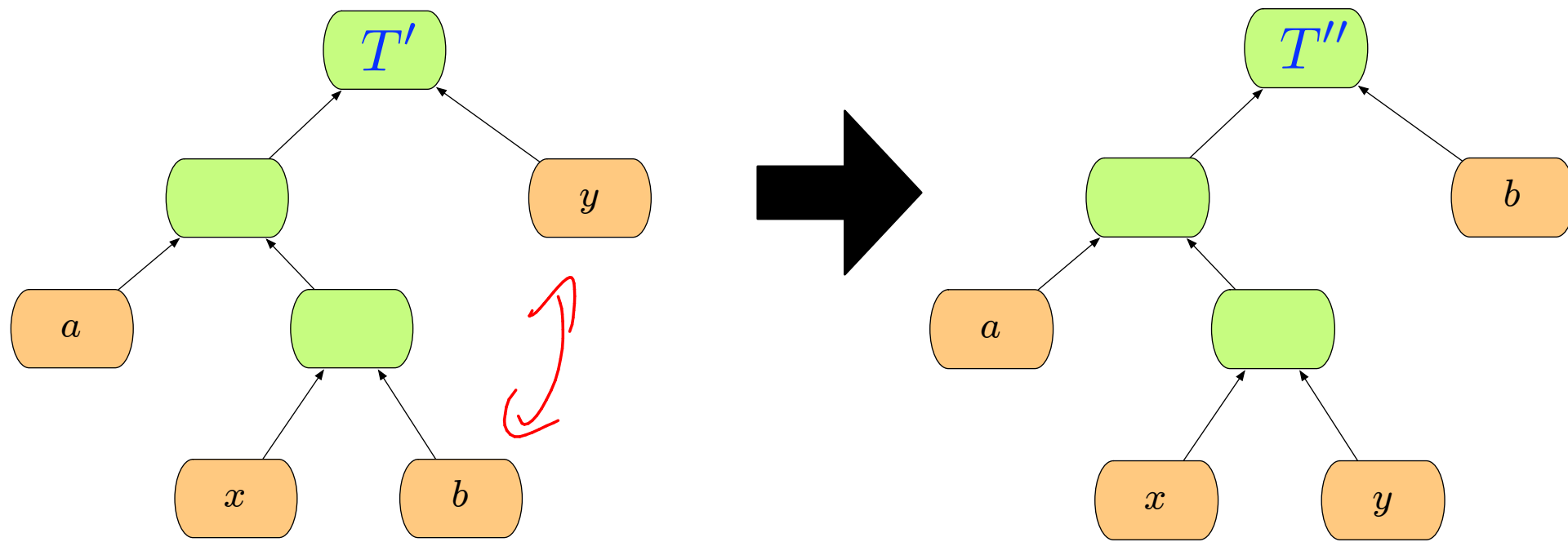
T is optimal \Rightarrow difference must be 0 \Rightarrow T' is optimal.



$$B(T) = \sum_c f_c l_c + f_x l_x + f_a l_a \quad B(T') = \sum_c f_c l'_c + f_x l'_x + f_a l'_a$$

$$B(T) - B(T') \geq 0$$

exchange argument

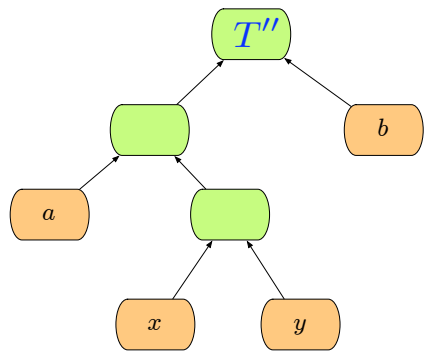
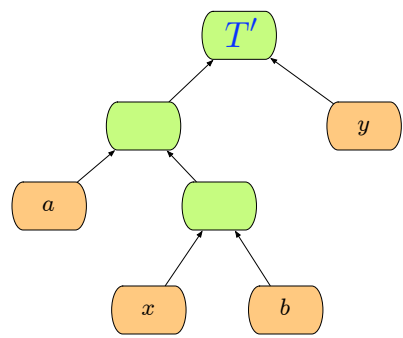
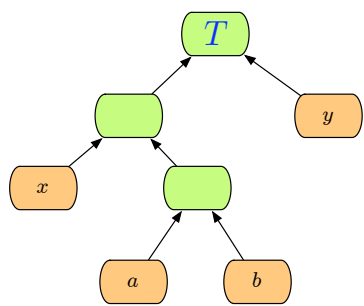


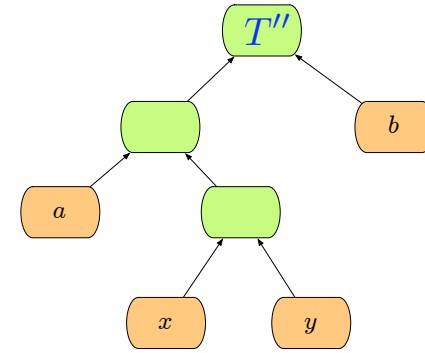
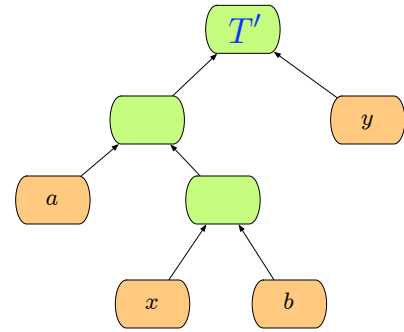
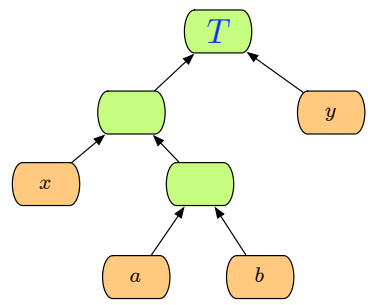
$$B(T') - B(T'') \geq 0$$

do the same argument with y and b .

$\Rightarrow T''$ is optimal, & it has x, y as siblings.







$$B(T) - B(T') \geq 0$$

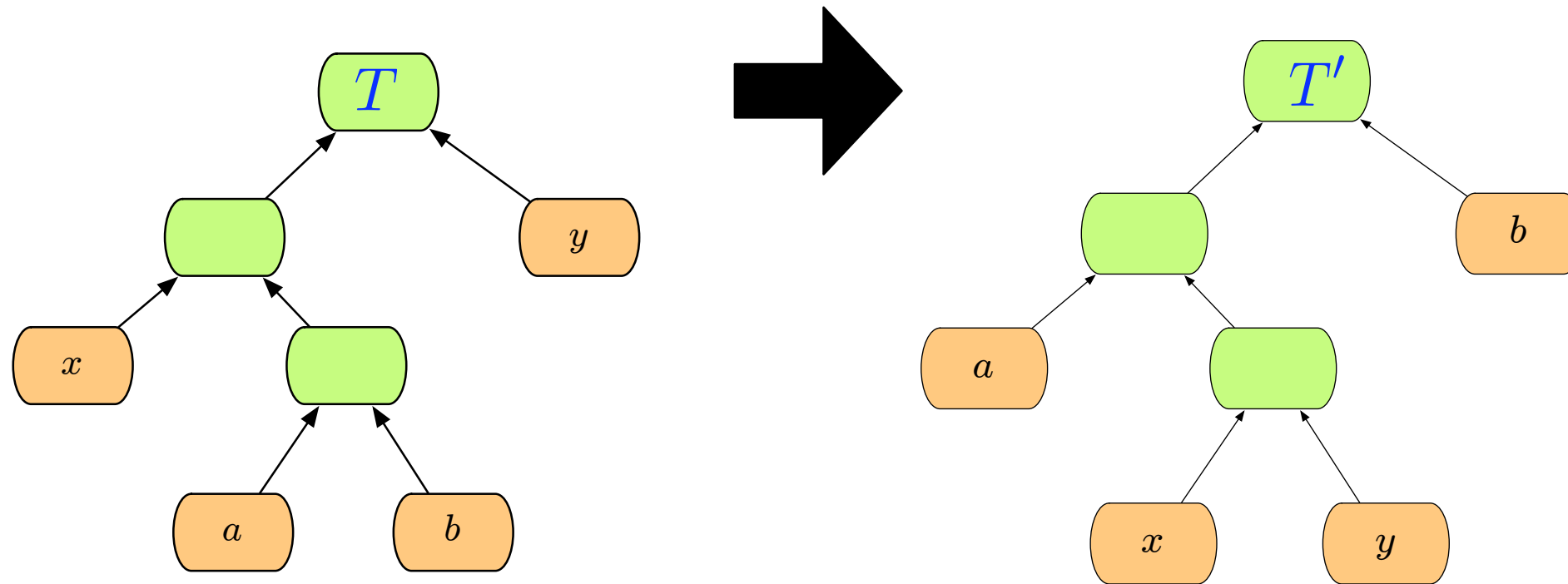
$$B(T') - B(T'') \geq 0$$

T'' is also optimal

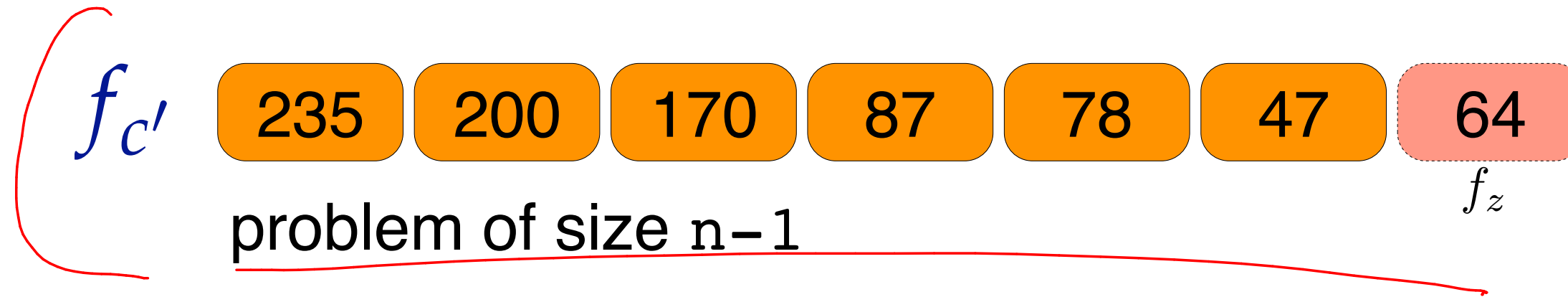
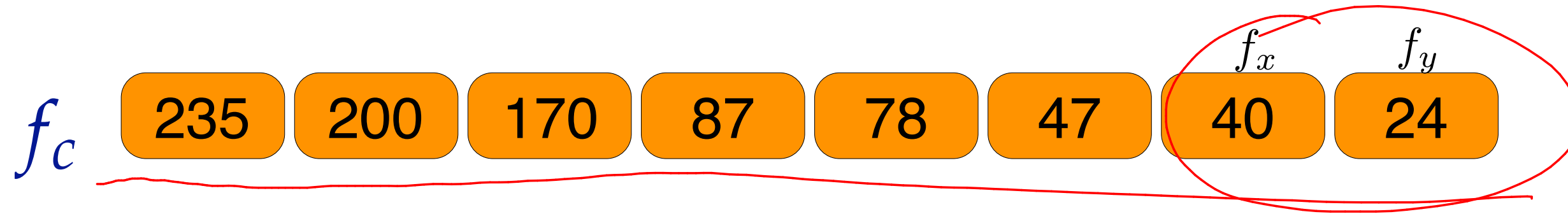
exchange argument

LEMMA:

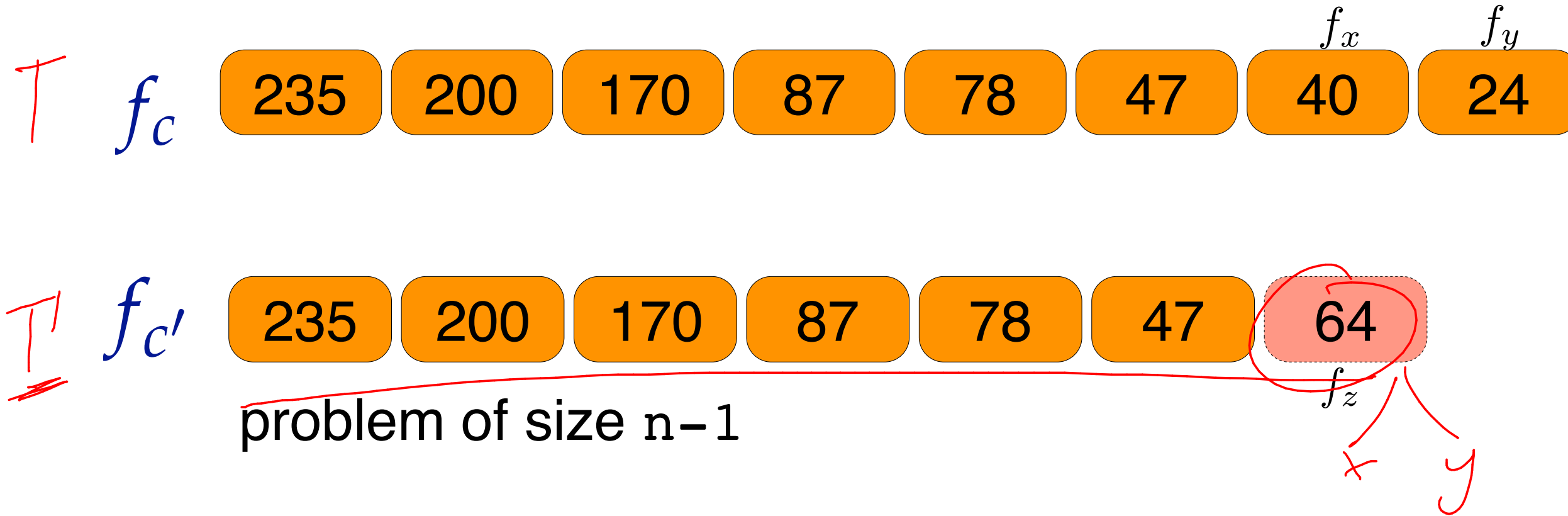
Let $x, y \in C$ be characters with smallest frequencies f_x, f_y . There exists an optimal prefix code T'' for C in which x, y are siblings. That is, the codes for x, y have the same length and only differ in the last bit.



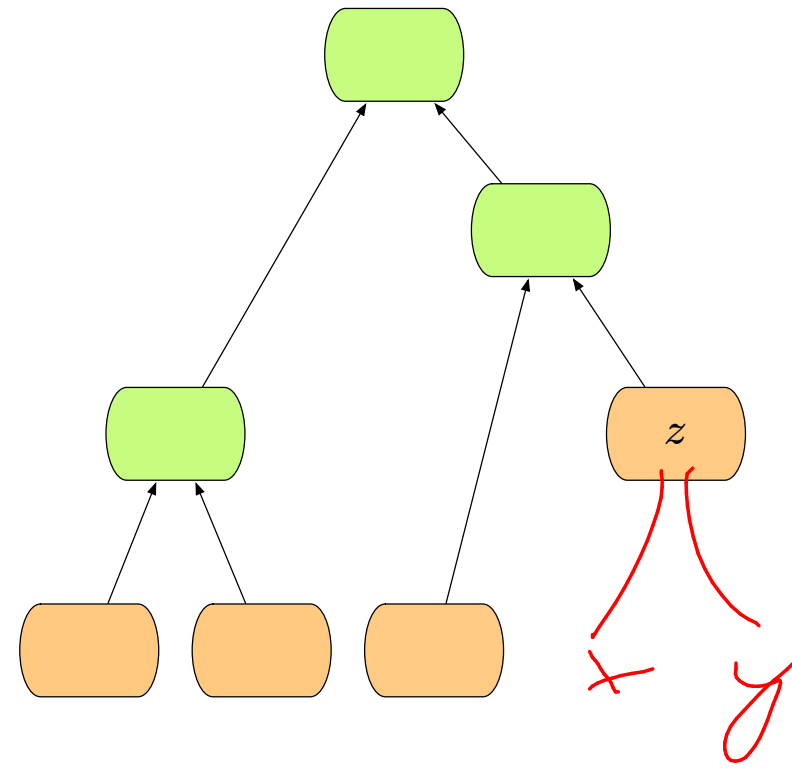
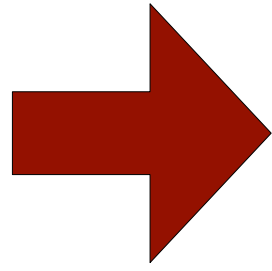
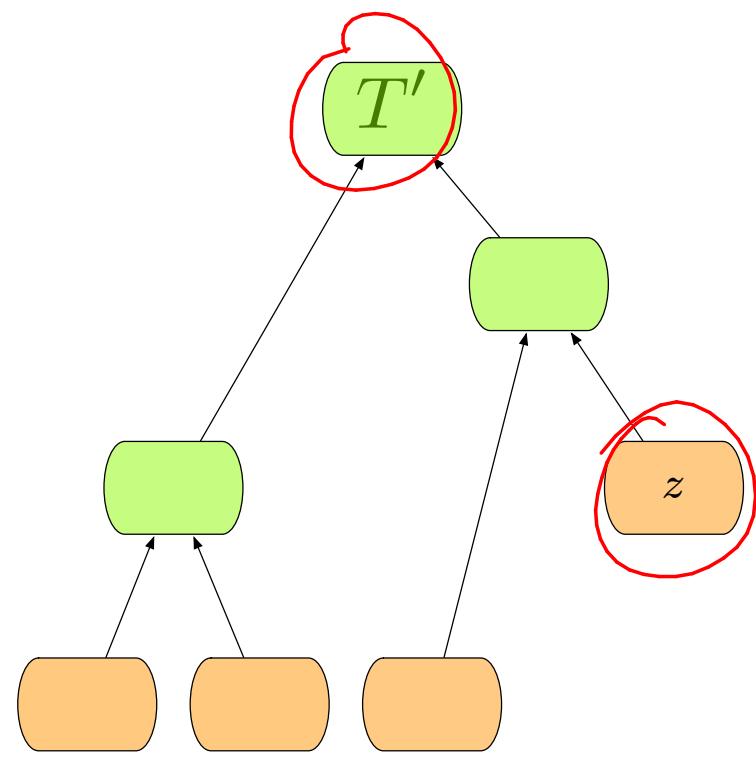
optimal sub-structure



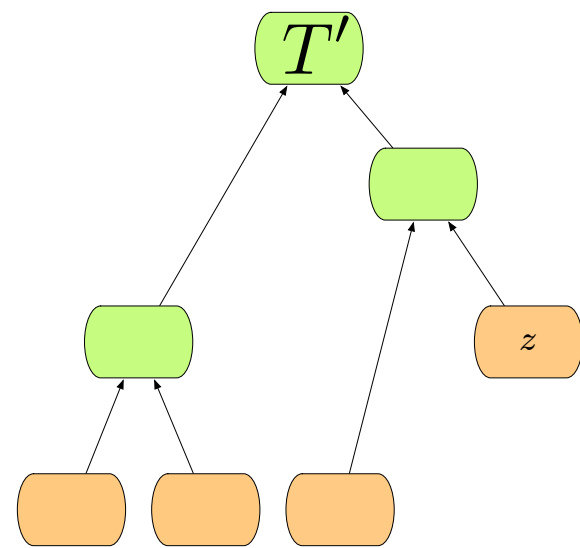
optimal sub-structure



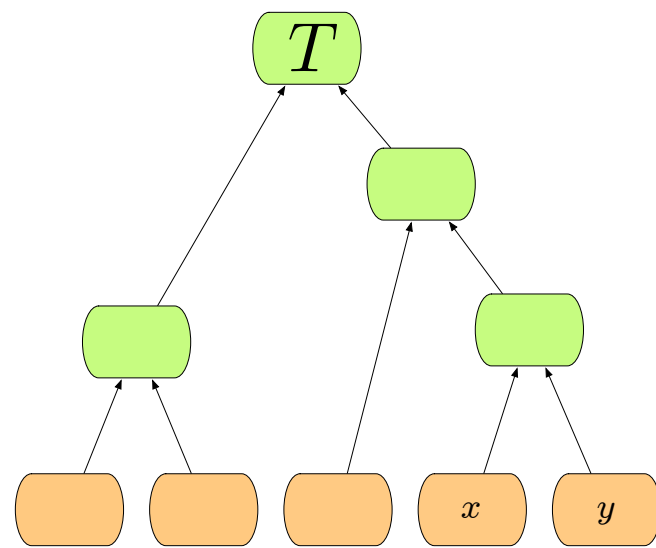
Lemma: The optimal solution for T consists of computing an optimal solution for T' and replacing the left z with a node having children x, y .



is optimal



$B(T')$

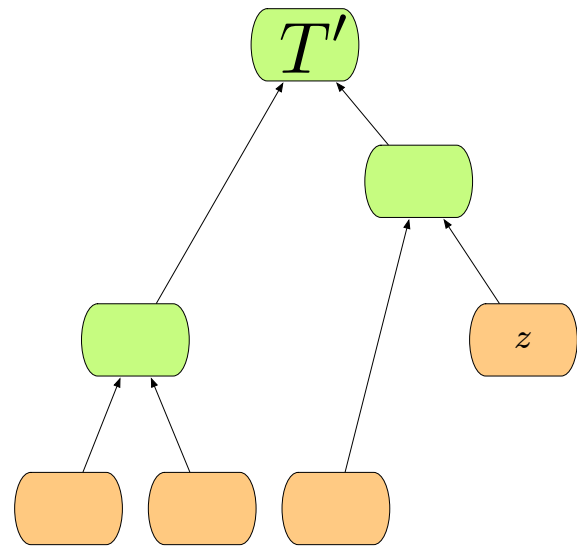


$B(T)$

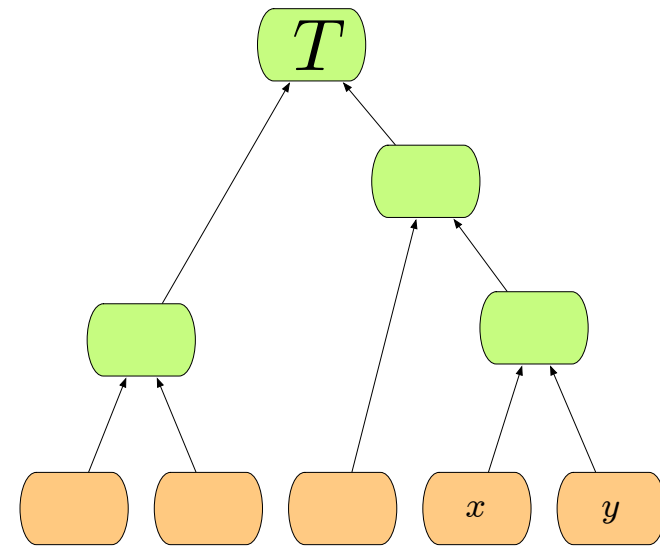
$$l_x = l_z(x)$$

$$l_y = l_z(x)$$

$$\underline{B(T')} = B(T) - f_x - f_y$$



$B(T')$



$B(T)$

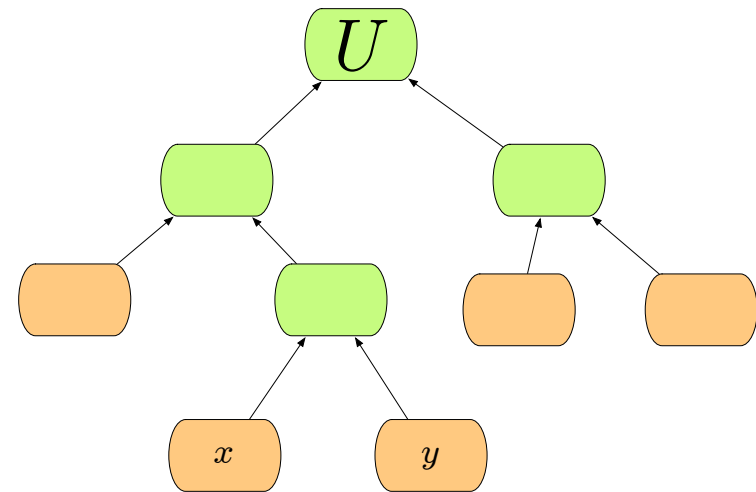
$$B(T') = B(T) - f_x - f_y$$

Suppose T is not *Optimal*.

I.e. \exists another tree U s.t. $B(U) < B(T)$

Suppose T is not optimal,

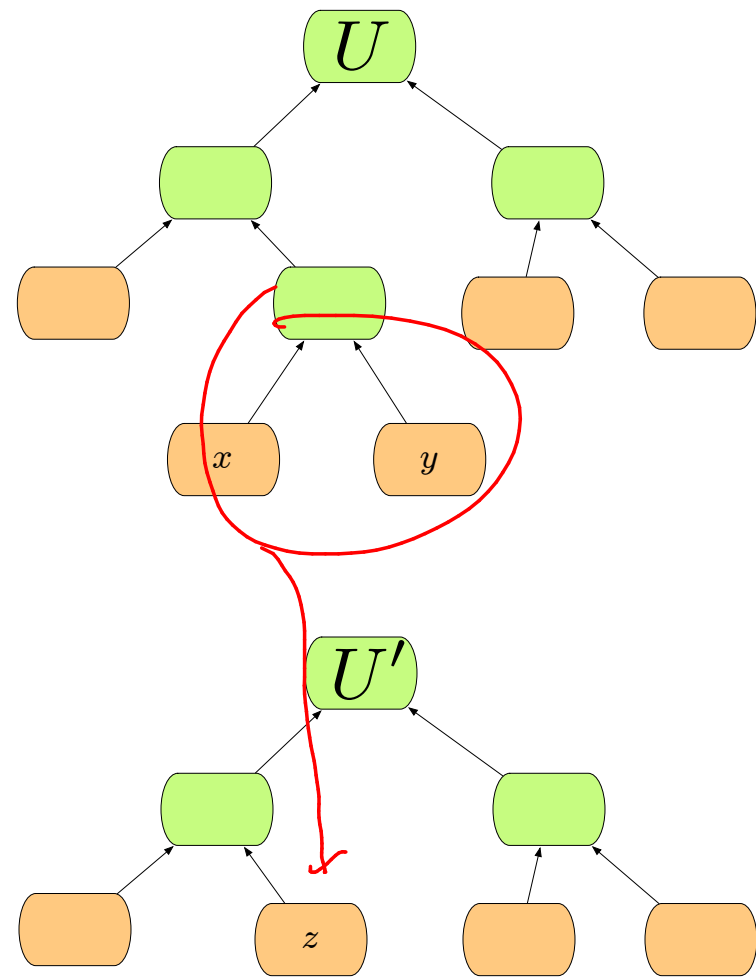
i.e. $\exists u$ s.t.



$$B(U) < B(T)$$

wlog, we can assume that x, y are siblings in u by Lemma I.

Suppose T is not



$$\underline{B(U) < B(T)}$$

Define a tree U' s.t. x & y are combined.

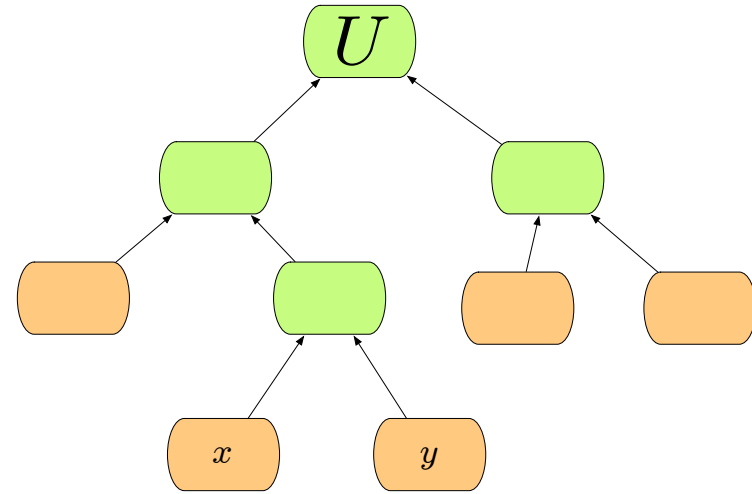
$$B(U') = B(U) - f_x - f_y$$

$$< \underline{B(T) - f_x - f_y}$$

$$= B(T')$$

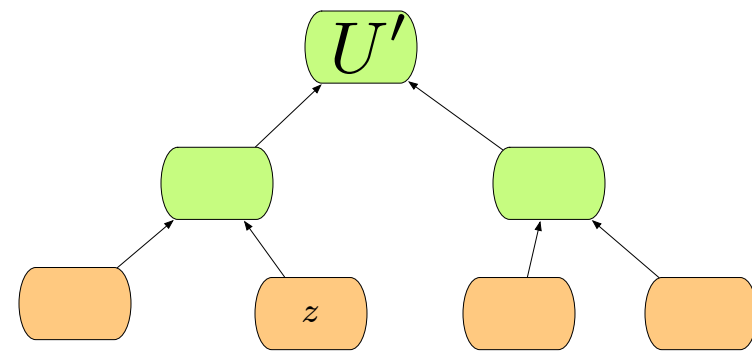
contradicts the assumption that T' was optimal $\Rightarrow T$ must have also been optimal \square

Suppose T is not



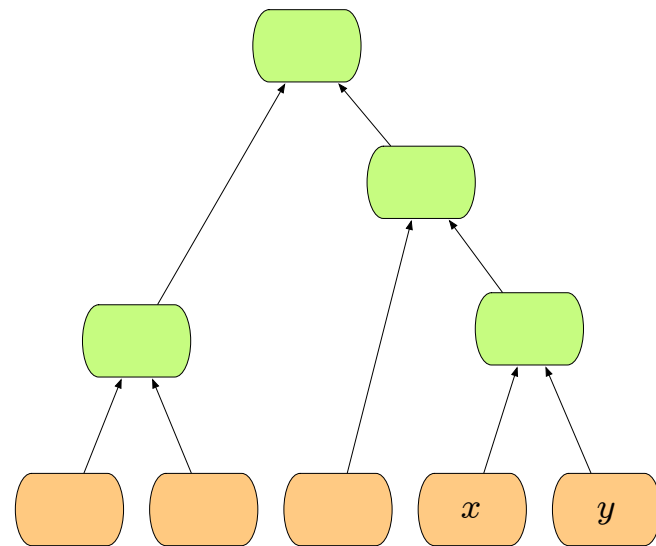
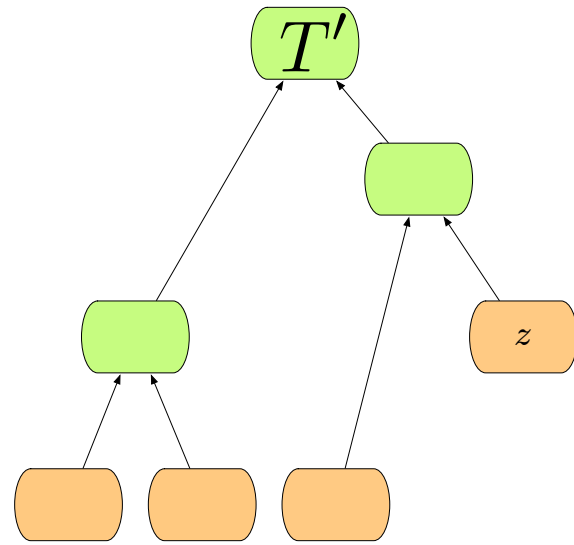
$$B(U) < B(T)$$

$$B(U') = B(U) - f_x - f_y$$
$$< B(t) - f_x - f_y$$



But this implies that $B(T')$ was not optimal

therefore



summary of argument

MST

connecting houses



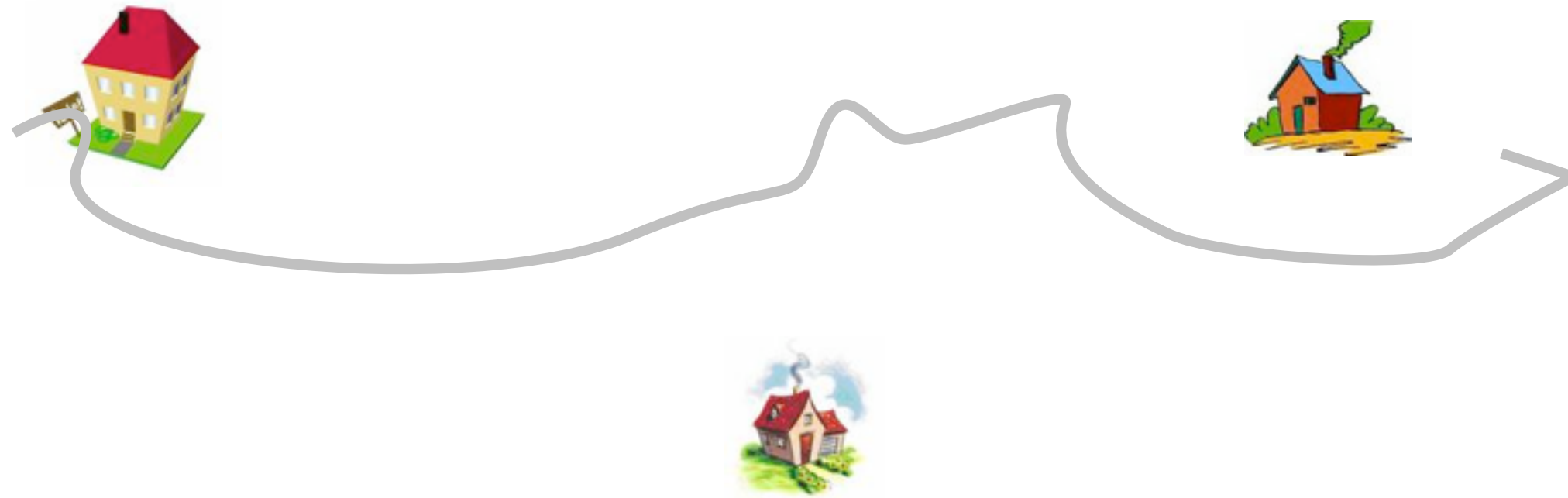
connecting houses



connecting houses



connecting houses



connecting houses



connecting houses

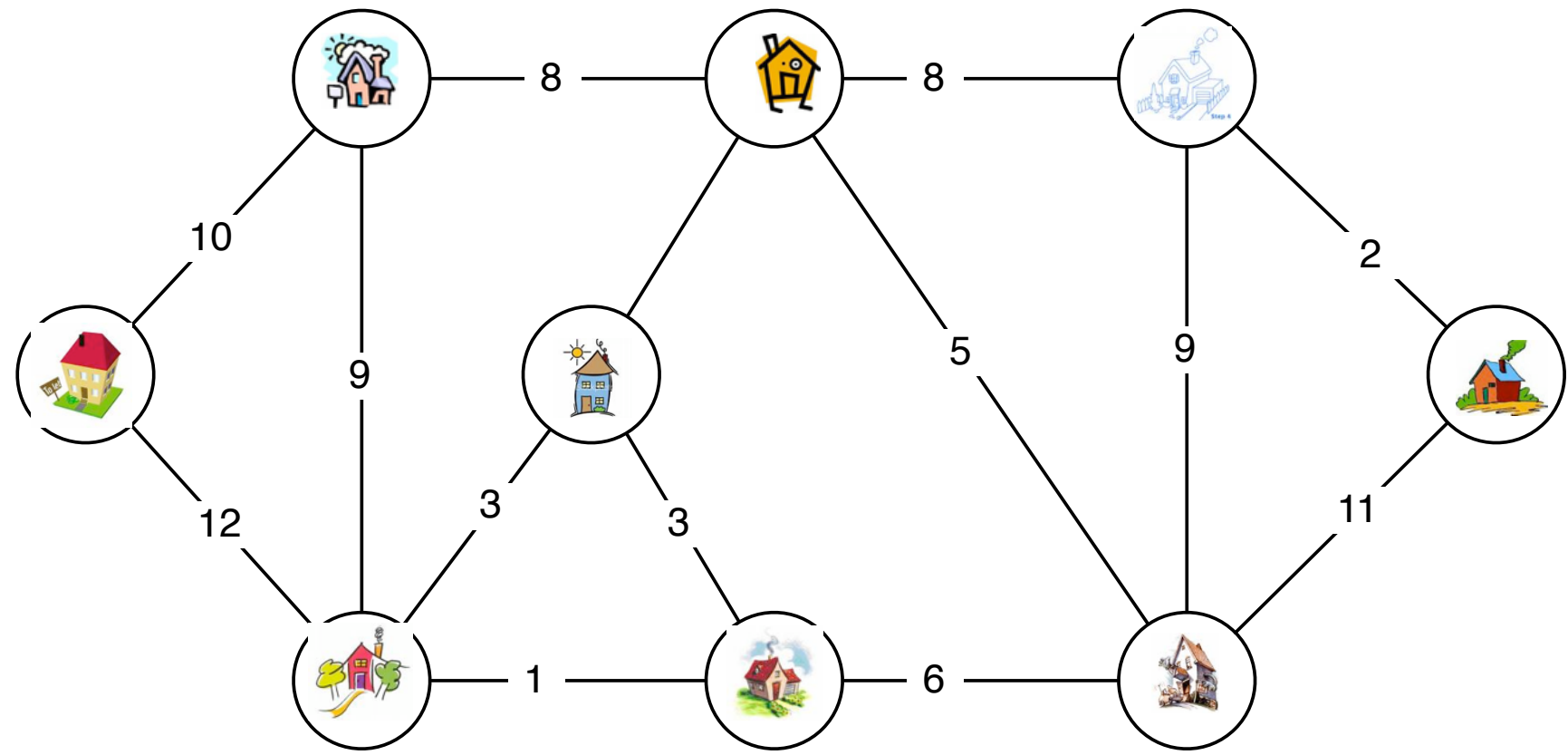


connecting houses



connecting houses

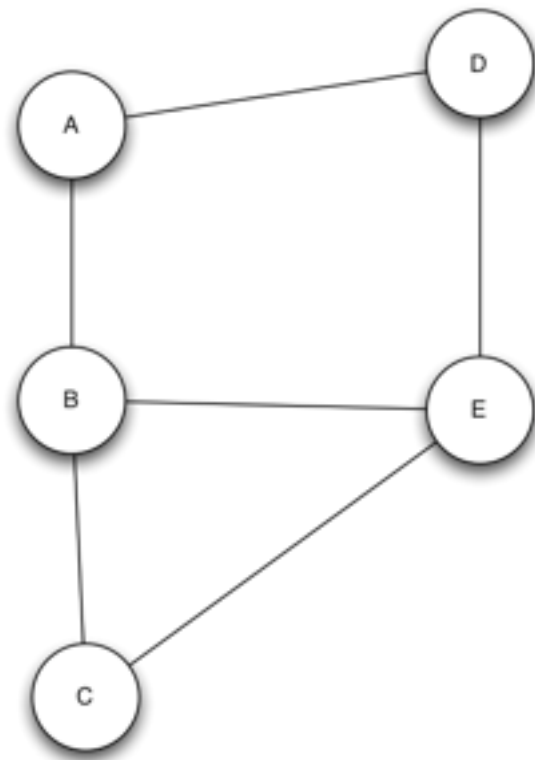






clrs [ch 22]

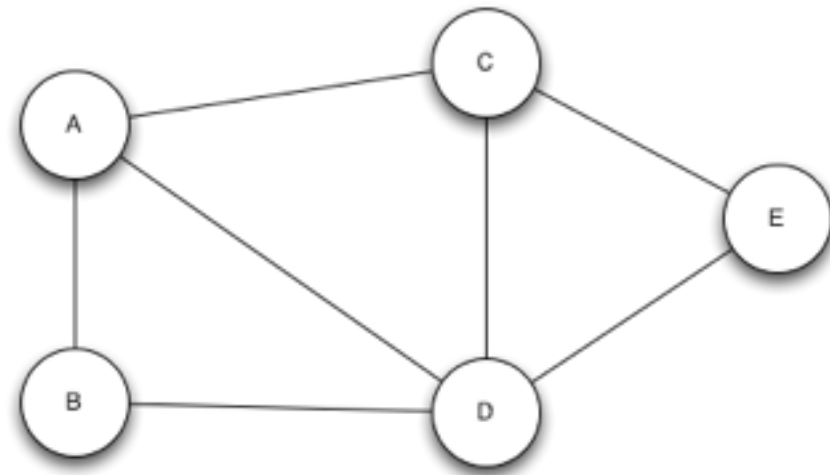
$$G = (V, E)$$



representation

$$G = (V, E)$$

adjacency list



space:

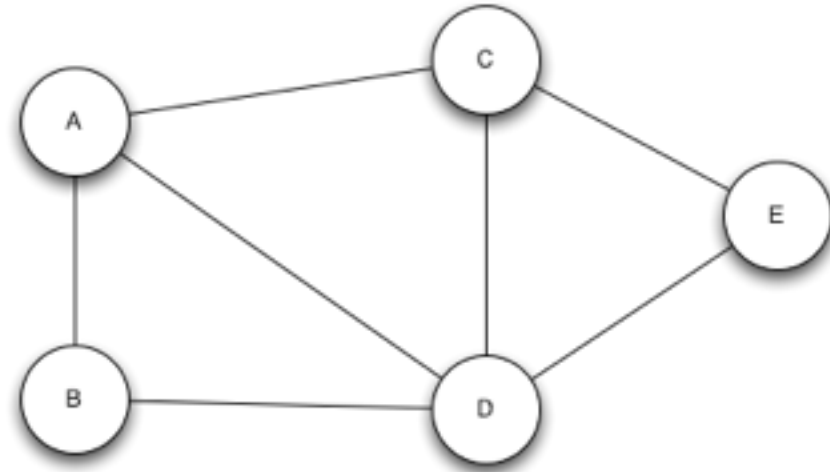
time list neighbors:

time check ~~an~~ edge:

representation

$$G = (V, E)$$

adjacency matrix



space:

time list neighbors:

time check an edge:

definition: path

a sequence of nodes v_1, v_2, \dots, v_k
with the property that $(v_i, v_{i+1}) \in E$

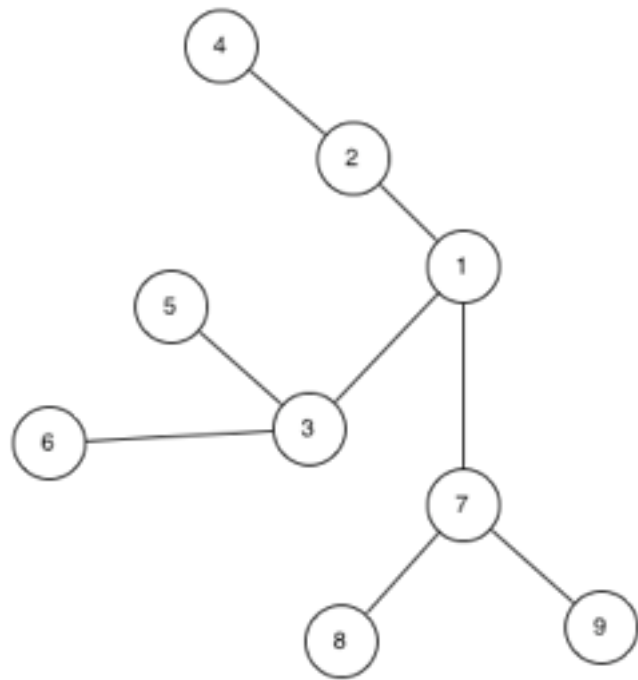
simple path:

cycle:

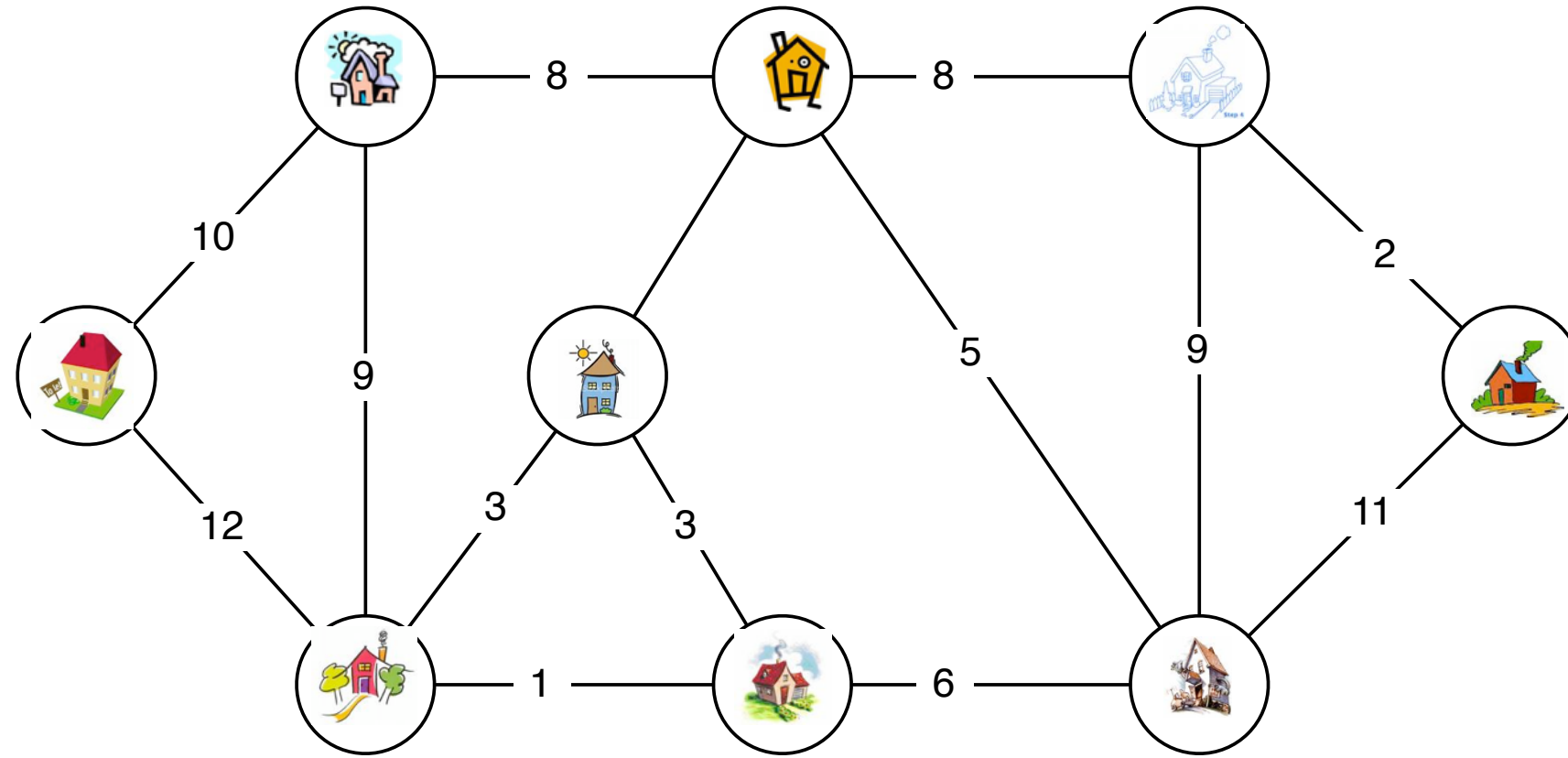
definition:tree

connected graph:

a tree is



what we want:



minimum spanning tree

looking for a set of edges $T \subseteq E$ that

(a) connects all vertices

(b) has the least cost

$$\min \sum_{(u,v) \in T} w(u,v)$$

facts

looking for a set of edges $T \subseteq E$ that

(a) connects all vertices

(b) has the least cost

$$\min \sum_{(u,v) \in T} w(u,v)$$

how many edges does solution have ?

does solution have a cycle?

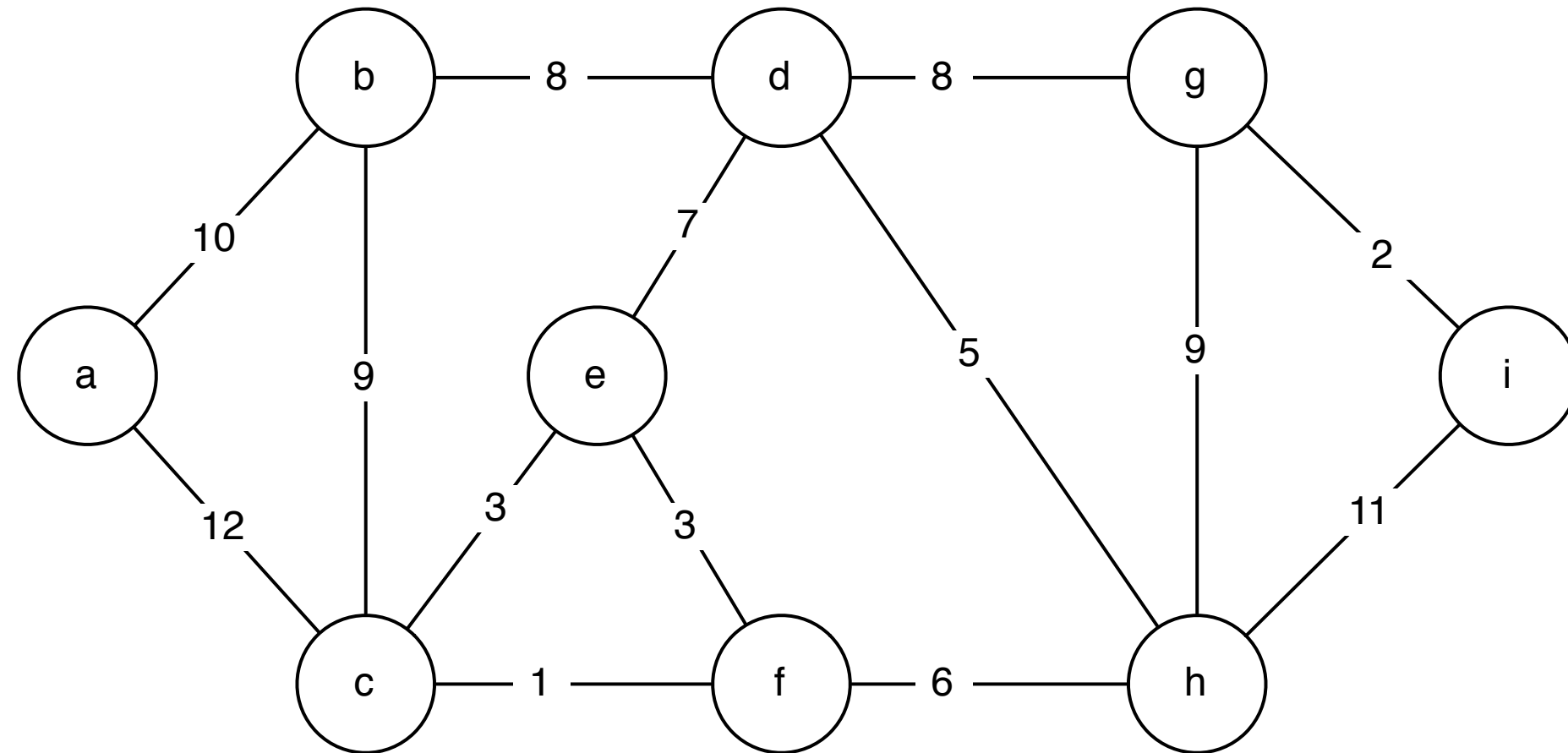
strategy

start with an empty set of edges A

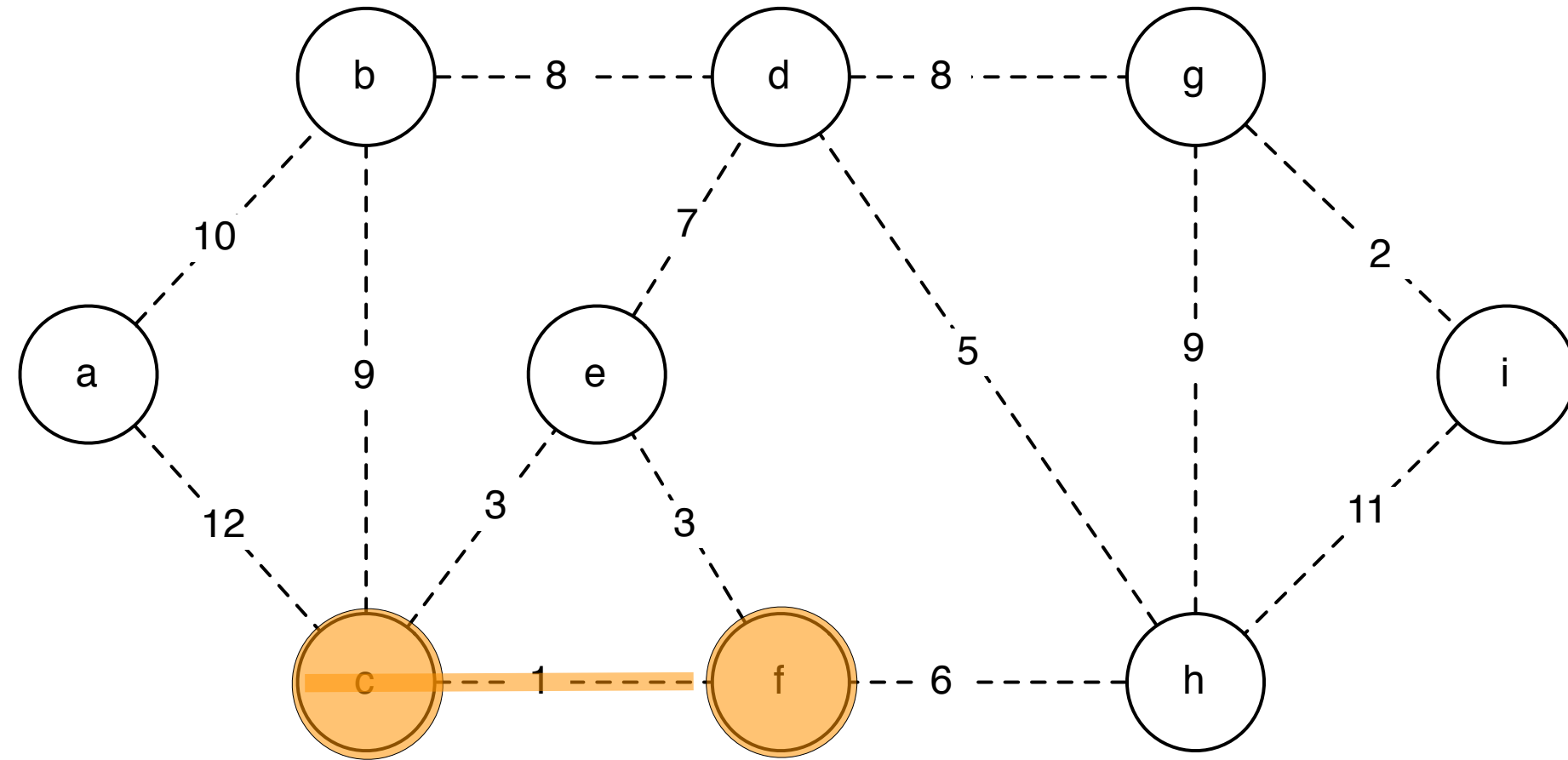
repeat for $v-1$ times:

 add lightest edge that does not create a cycle

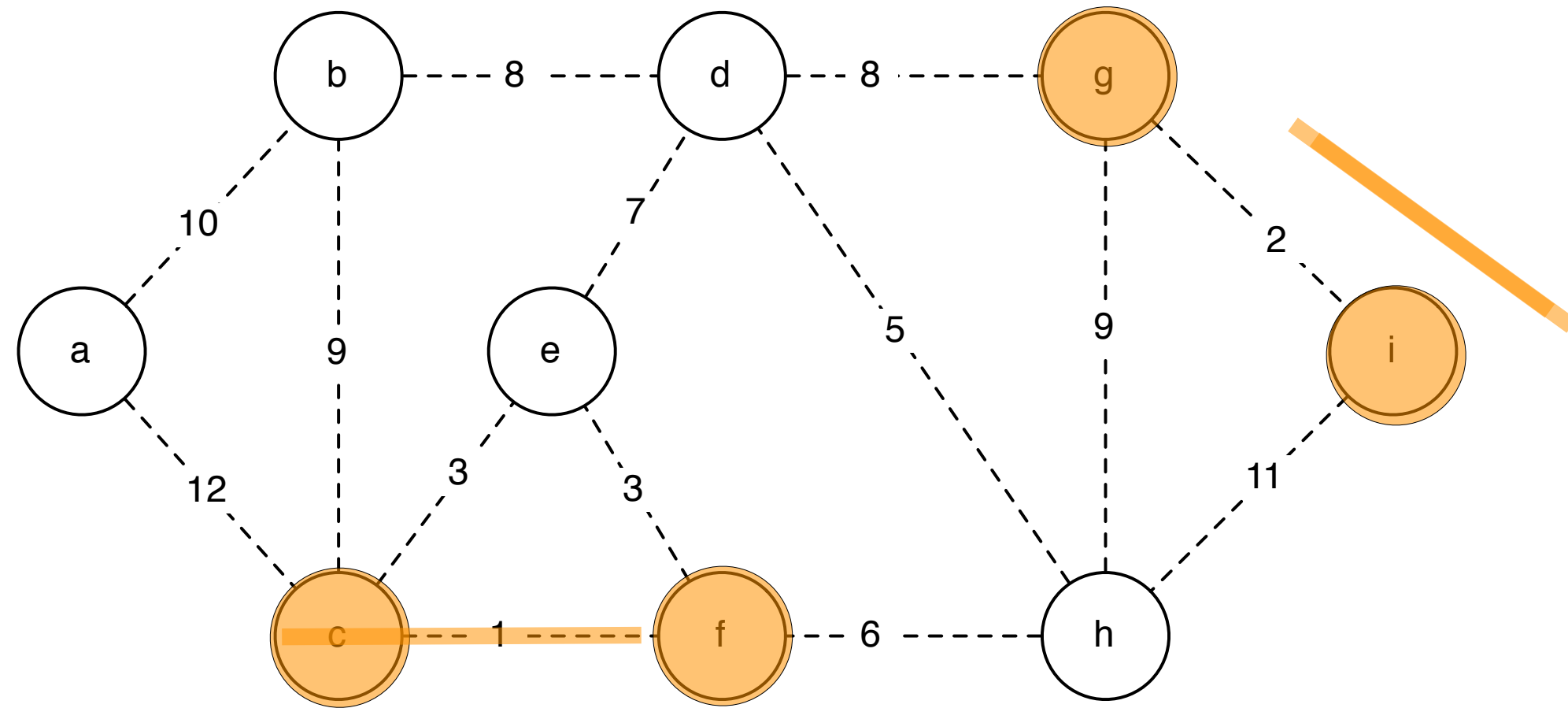
example



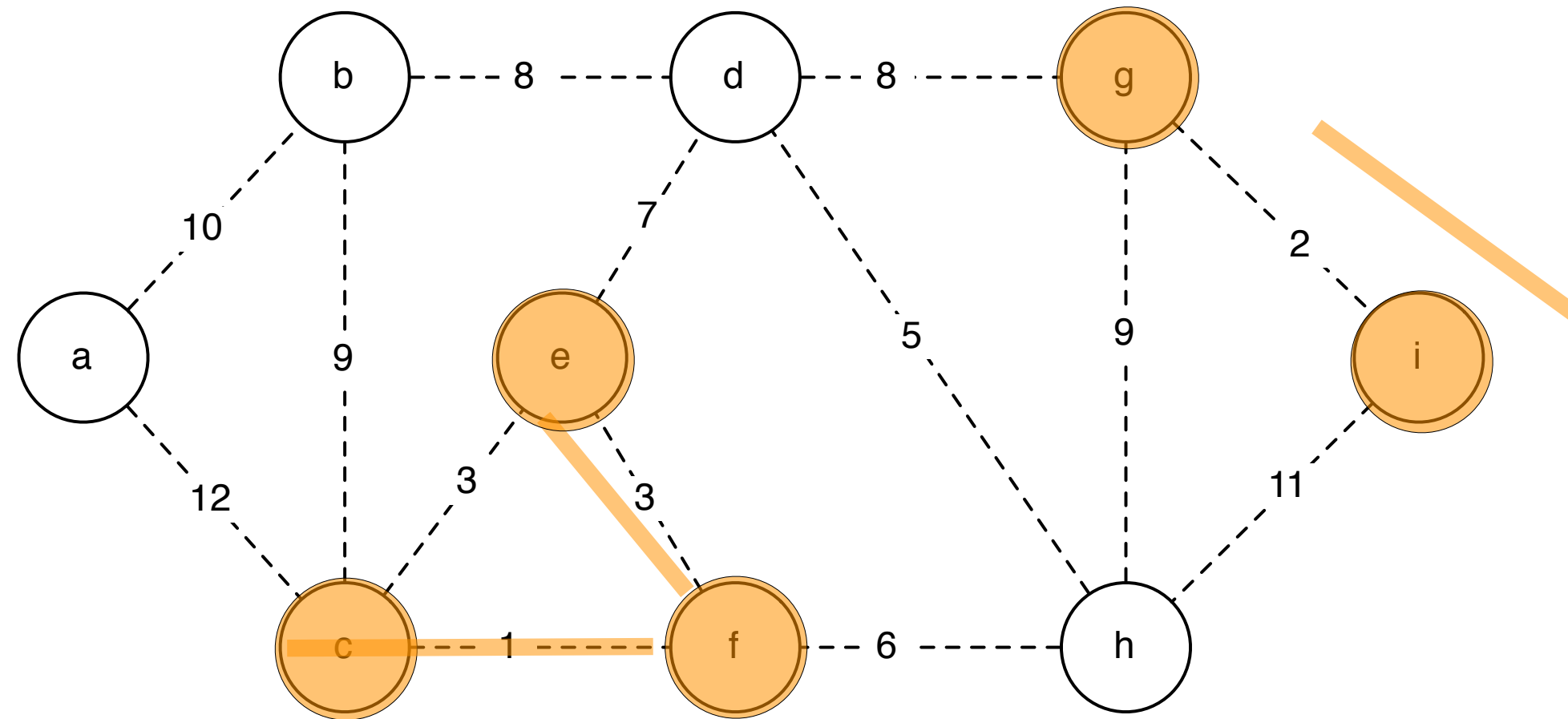
kruskal



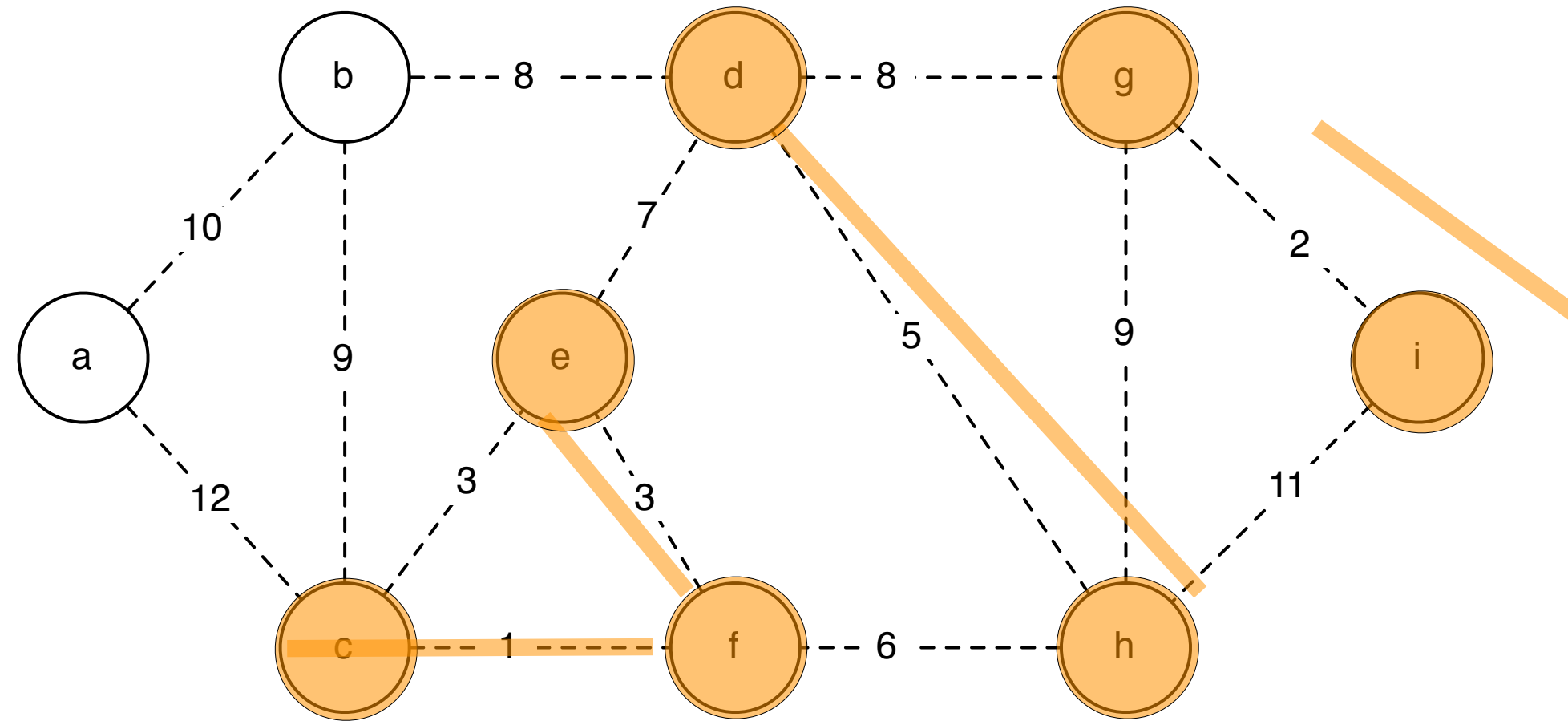
kruskal



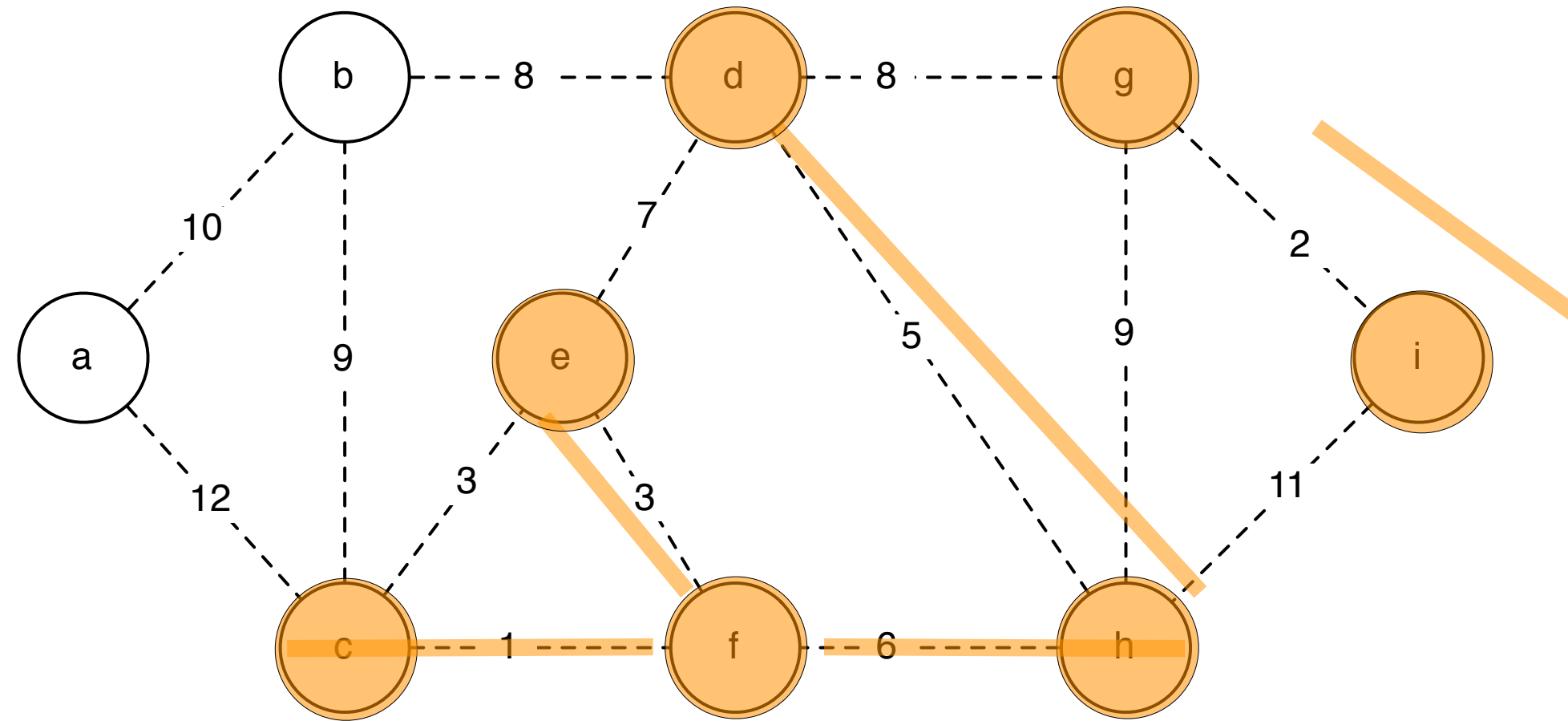
kruskal



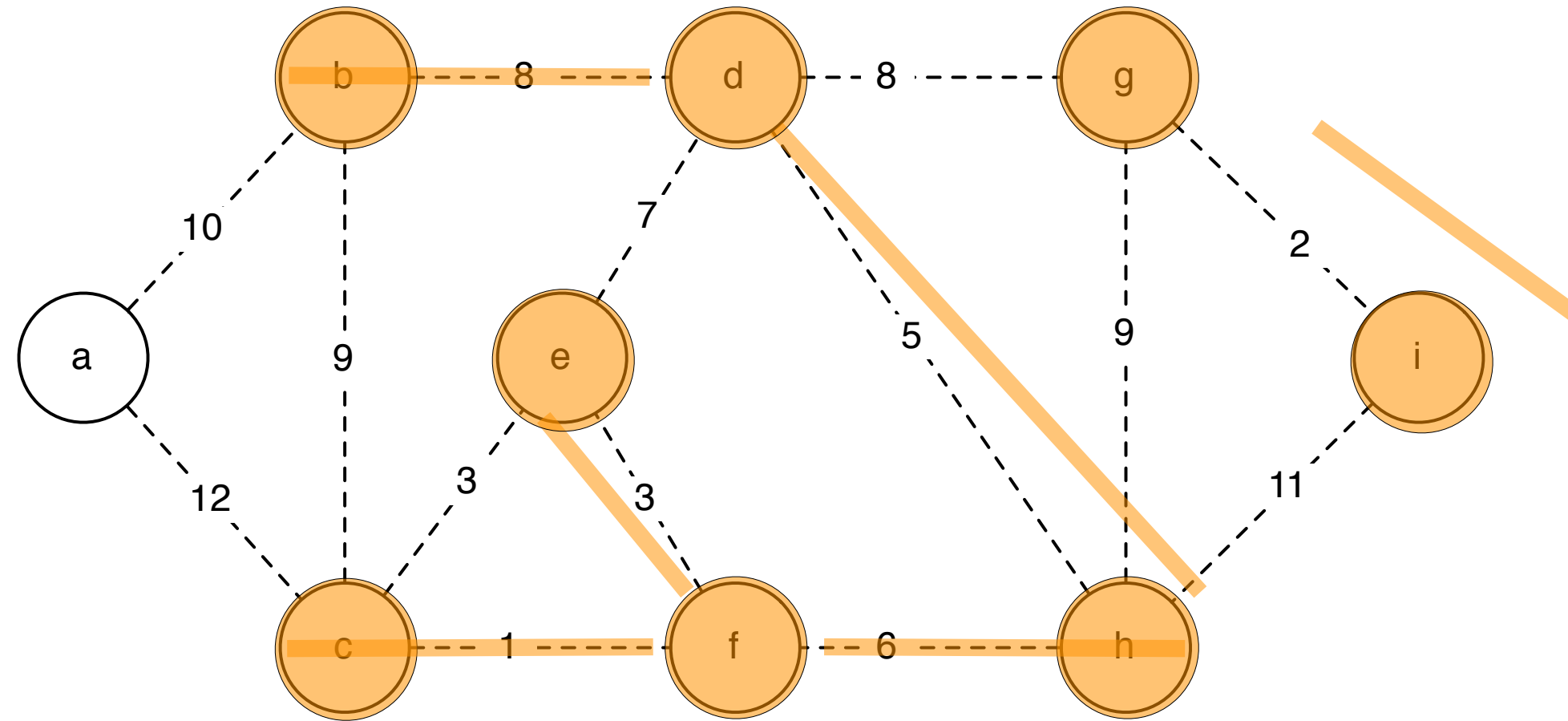
kruskal



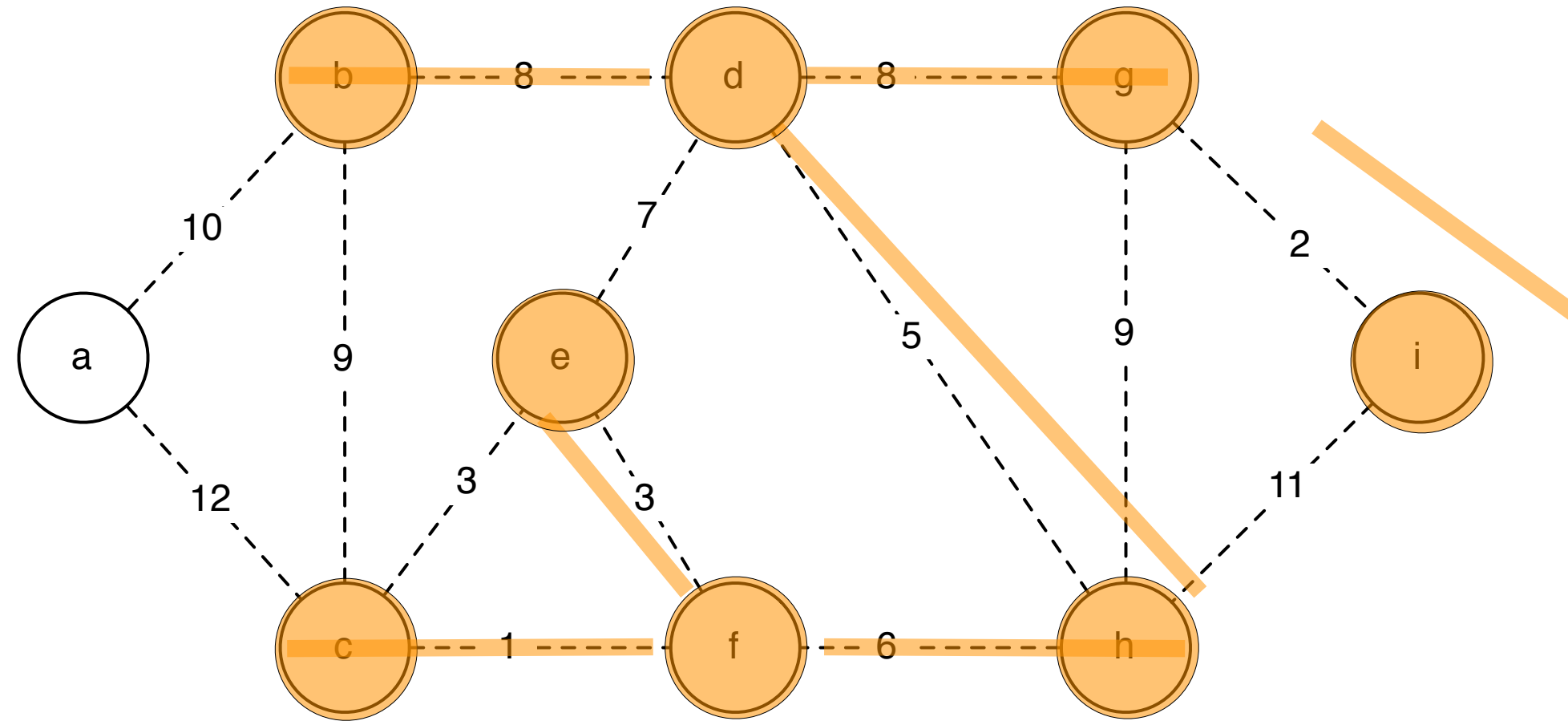
kruskal



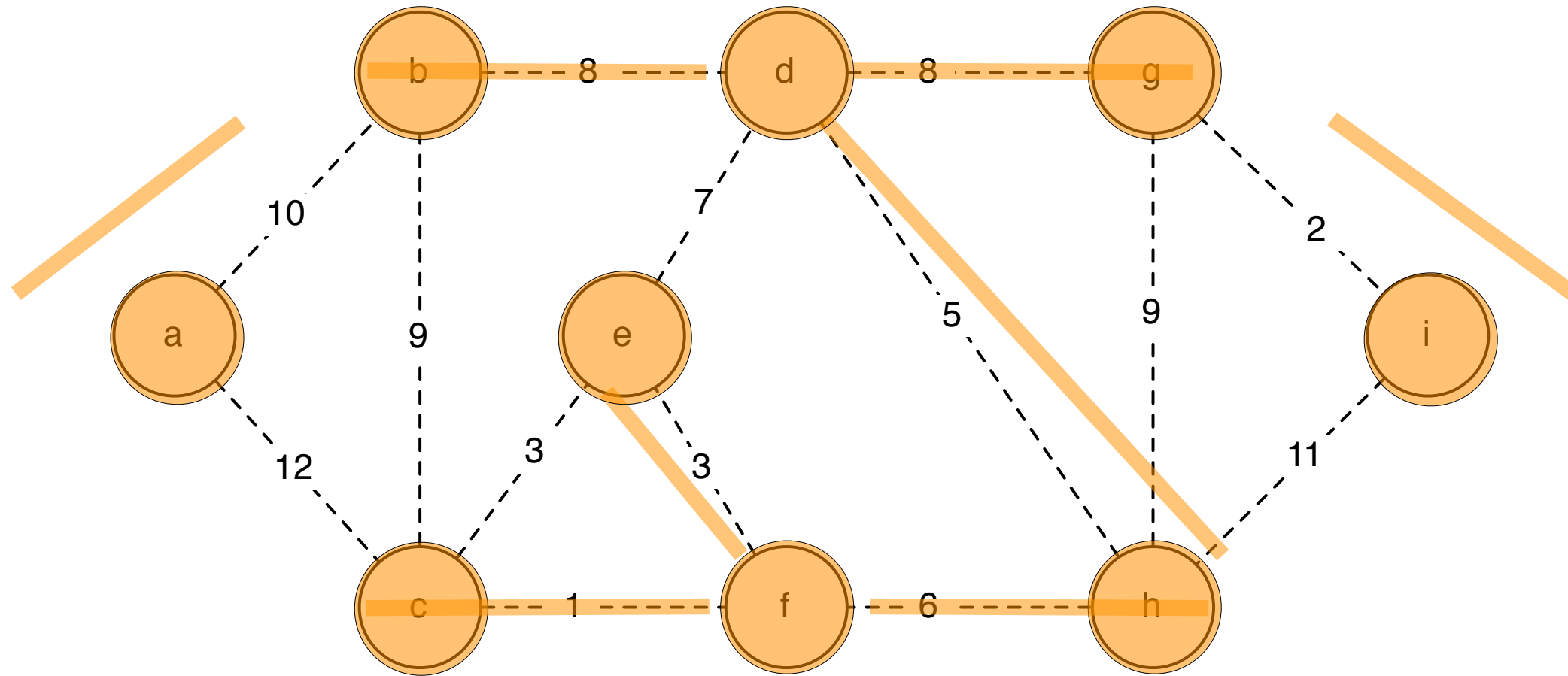
kruskal



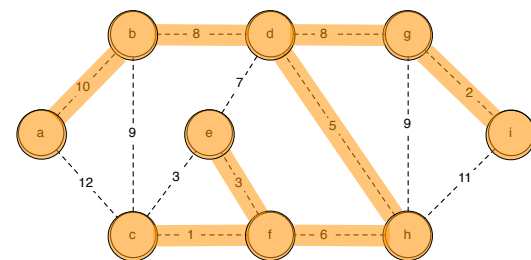
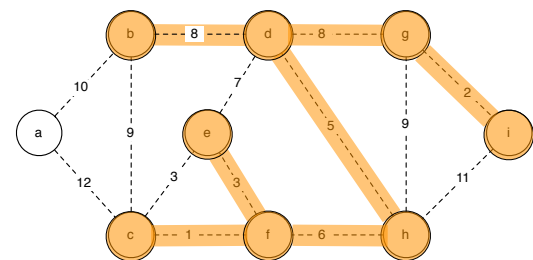
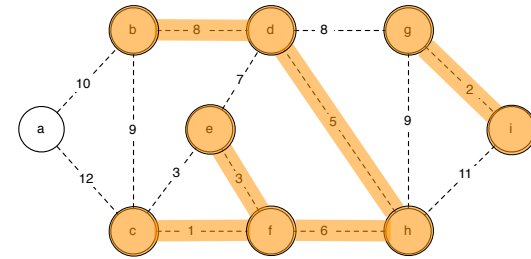
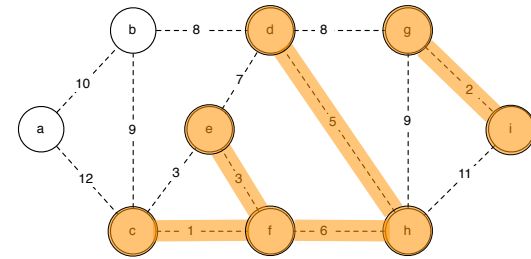
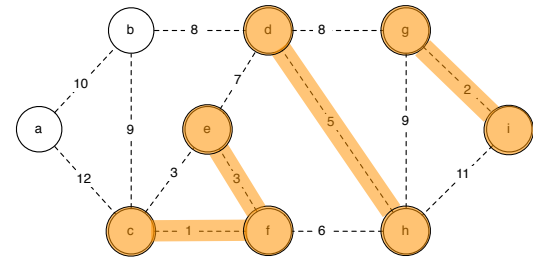
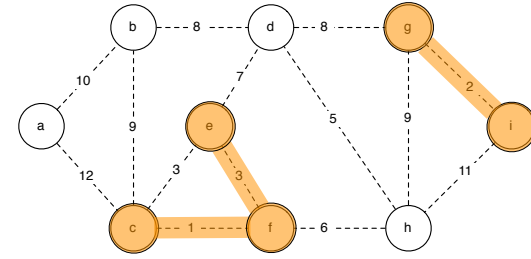
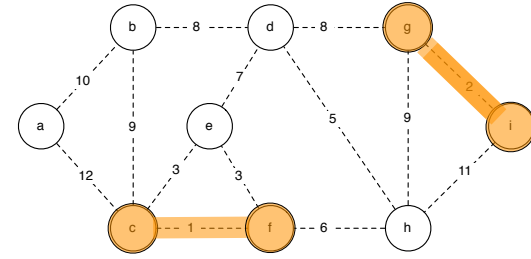
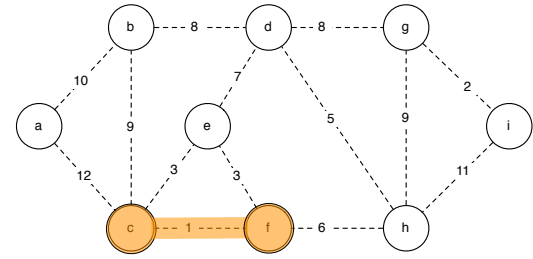
kruskal



kruskal



kruskal



why does this work?

- 1 $T \leftarrow \emptyset$
- 2 **repeat** $V - 1$ times:
- 3 add to T the lightest edge $e \in E$ that does not create a cycle