

L11

4800

10.18.2016

abhi shelat

# Huffman

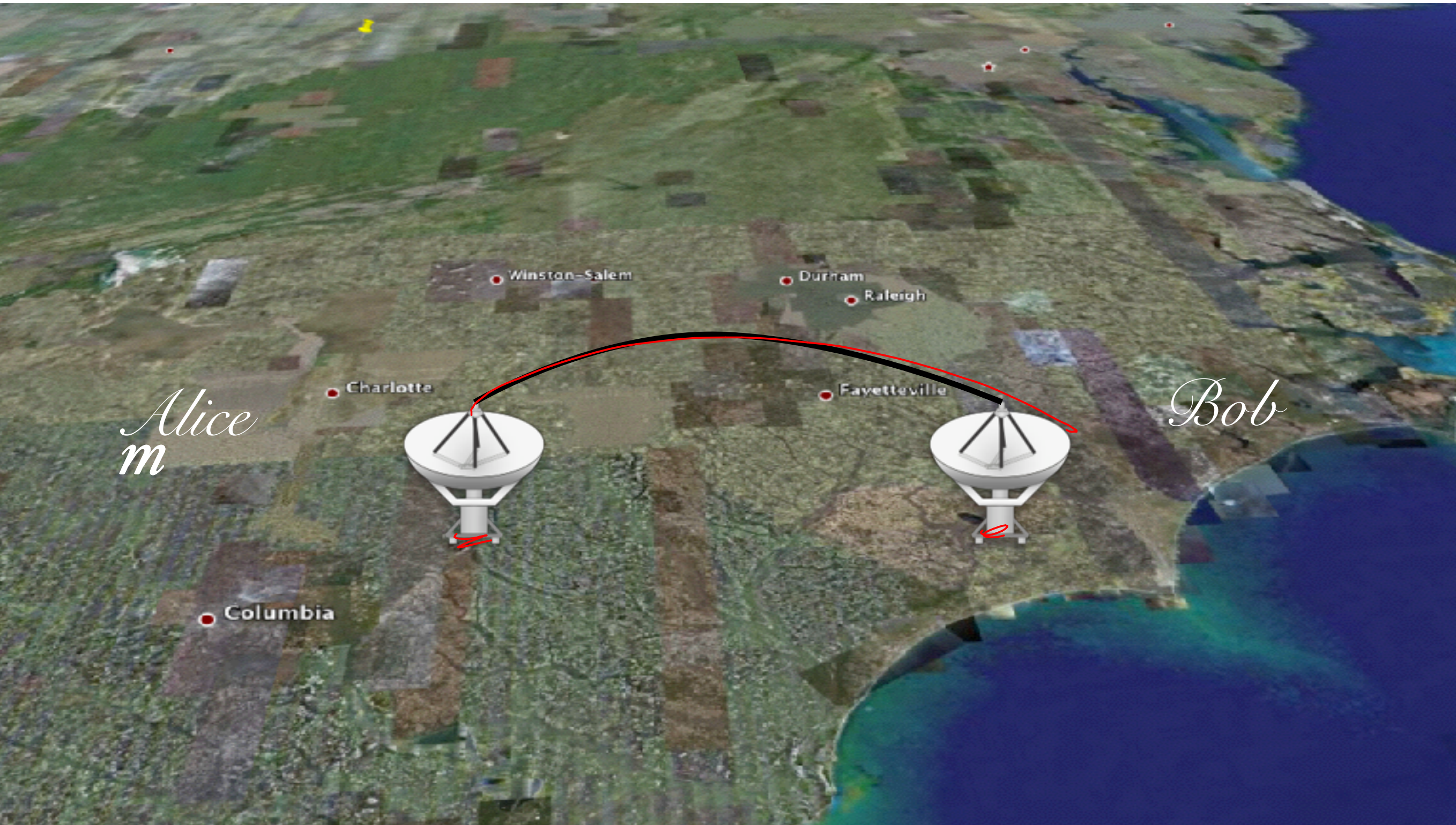
## CS4800

abhi shelat



image: wikimedia





*Alice*  
*m*

*Bob*

Columbia

Charlotte

Winston-Salem

Durham

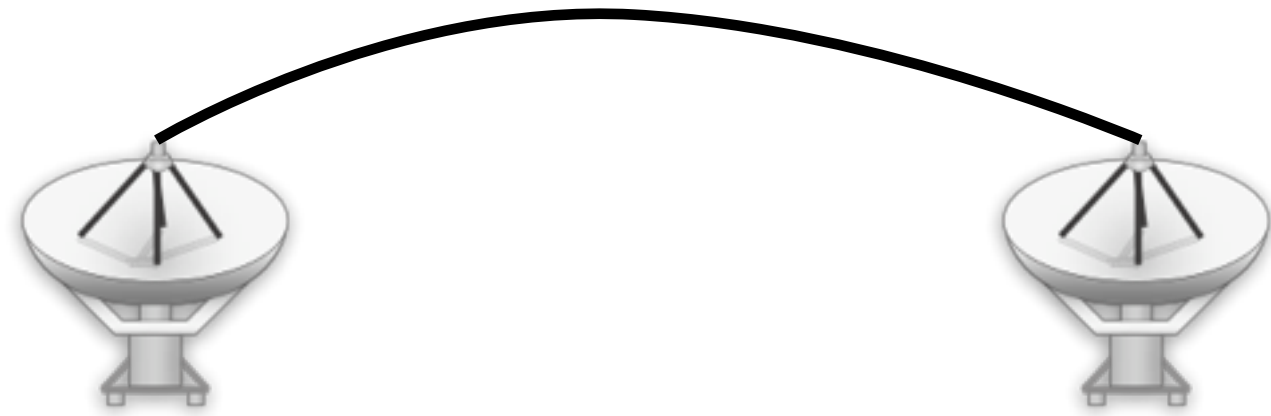
Raleigh

Fayetteville



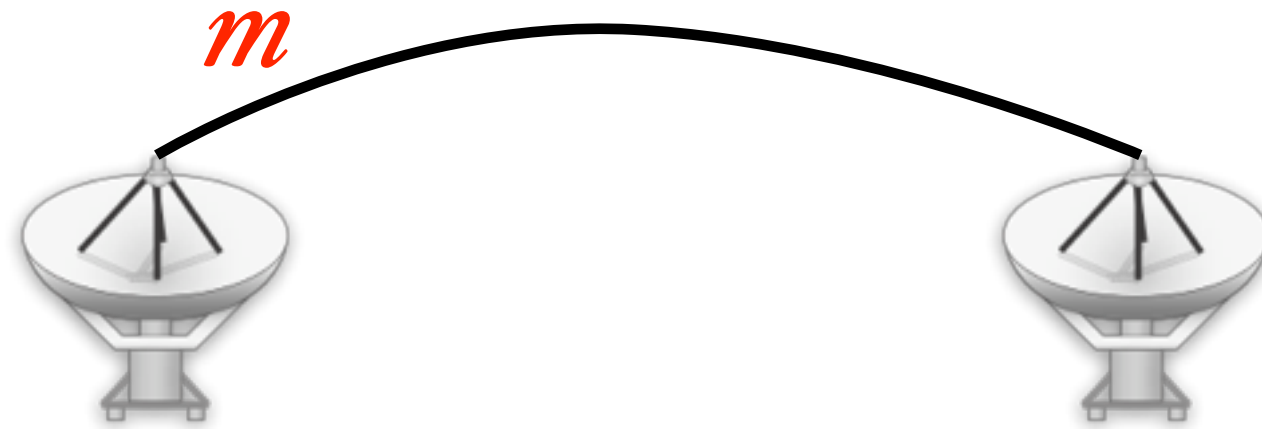






MOSCOW — President Vladimir V. Putin's typically theatrical order to withdraw the bulk of Russian forces from Syria, a process that the Defense Ministry said it began on Tuesday, seemingly caught Washington, Damascus and everybody in between off guard — just the way the Russian leader likes it.

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 world events.



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 world events.

$c \in C$   $f_c$   $T$

e: 235

i: 200

o: 170

u: 87

p: 78

g: 47

b: 40

f: 24

881



=

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

881

-> 2643

# def: cost of an encoding

*code* ↓  
*frequencies of your messages* ↗  
*length of code for character c in code T* ↖

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

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

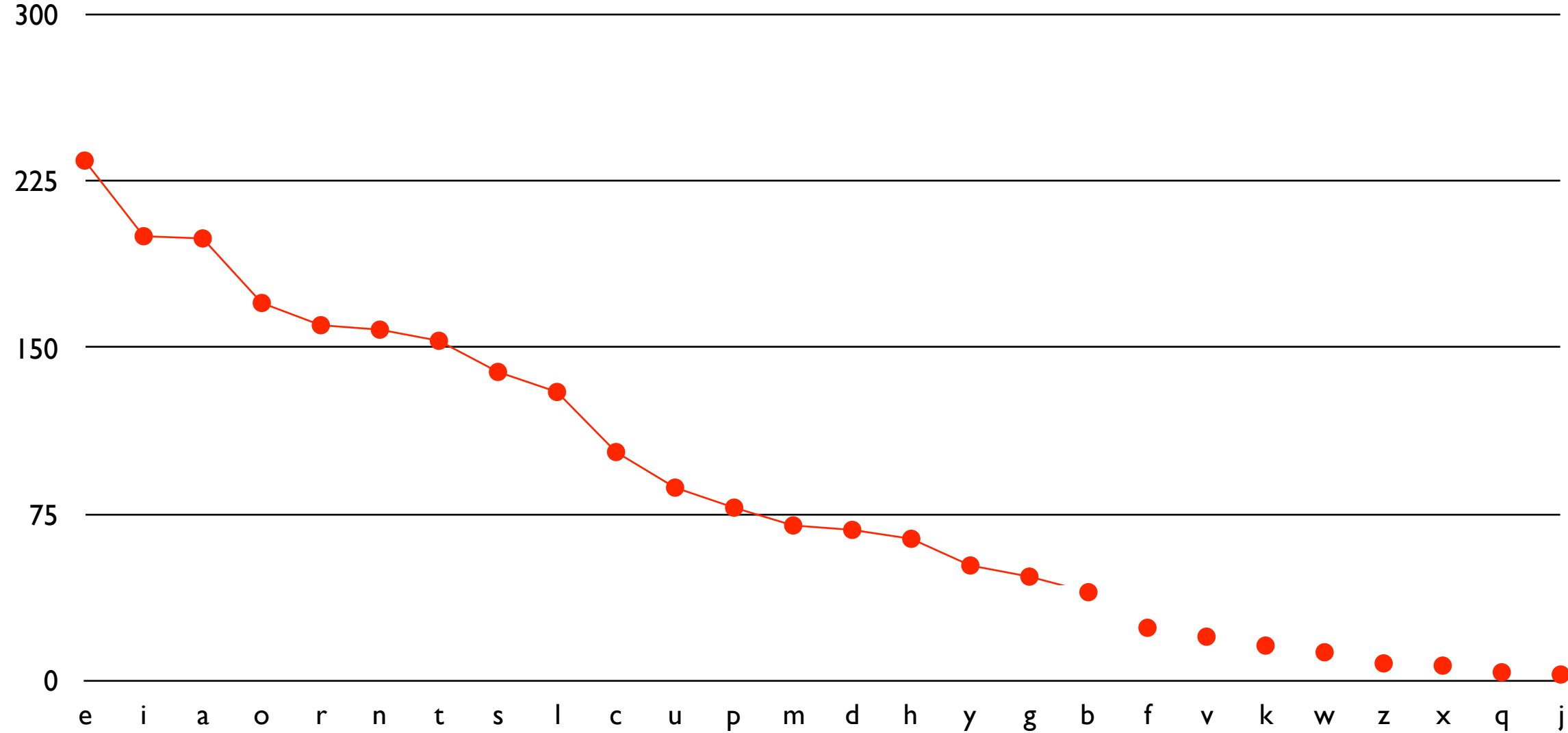
881

$\cdot 3 = 2643$



# 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



# 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](http://en.wikipedia.org/wiki/Morse_code)

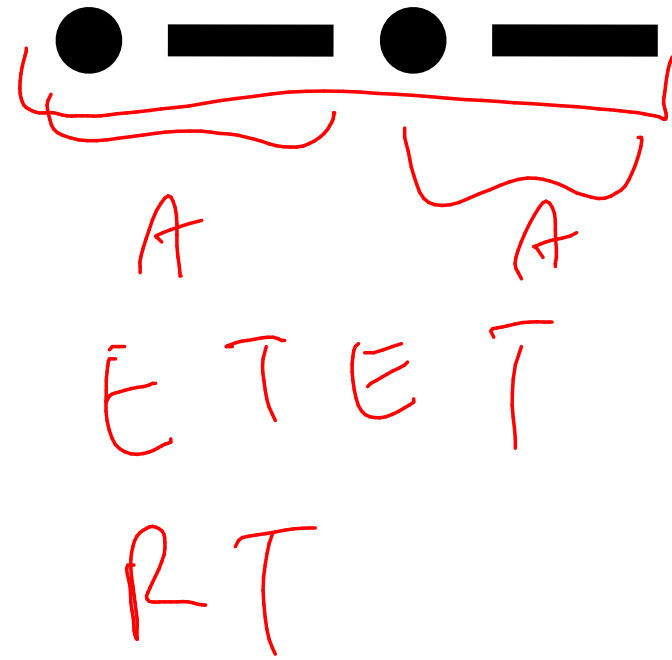


# morse code

## 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	• • —		



# def: prefix-free code

Code such that for any two symbols  $x, y \in C$   
 $x \neq y$  )  $\text{code}(x)$  is not a prefix of  $\text{code}(y)$



# 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

$\forall x, y \in C, x \neq y \implies \text{CODE}(x)$  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

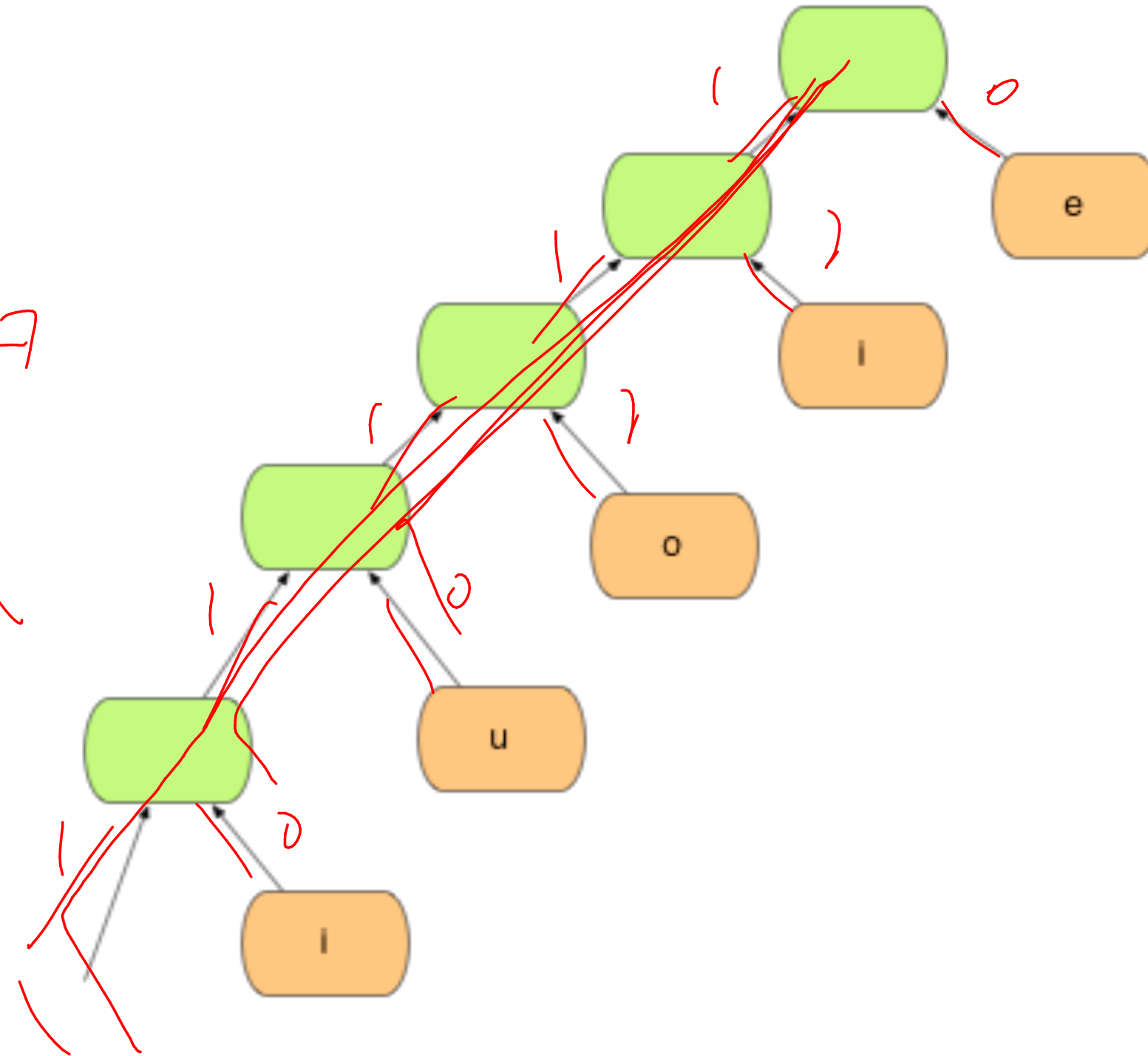
e:	235	0
i:	200	10
o:	170	110
u:	87	1110
p:	78	11110
g:	47	111110
b:	40	<u>1111110</u>
f:	24	11111110

111111010111110

# code to binary tree

e: 235  
i: 200  
o: 170  
u: 87  
p: 78  
g: 47  
b: 40  
f: 24

0  
10  
110  
1110  
11110  
111110  
1111110  
11111110

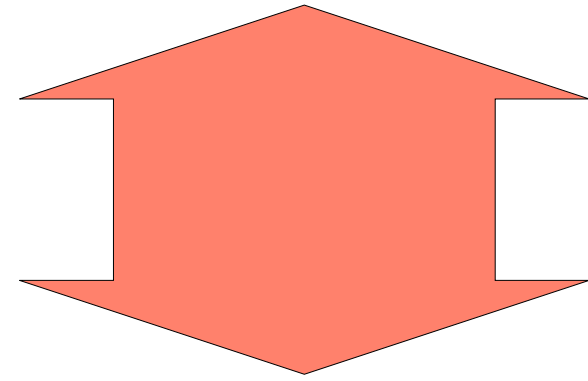


1111110101 / 11110

B I 9

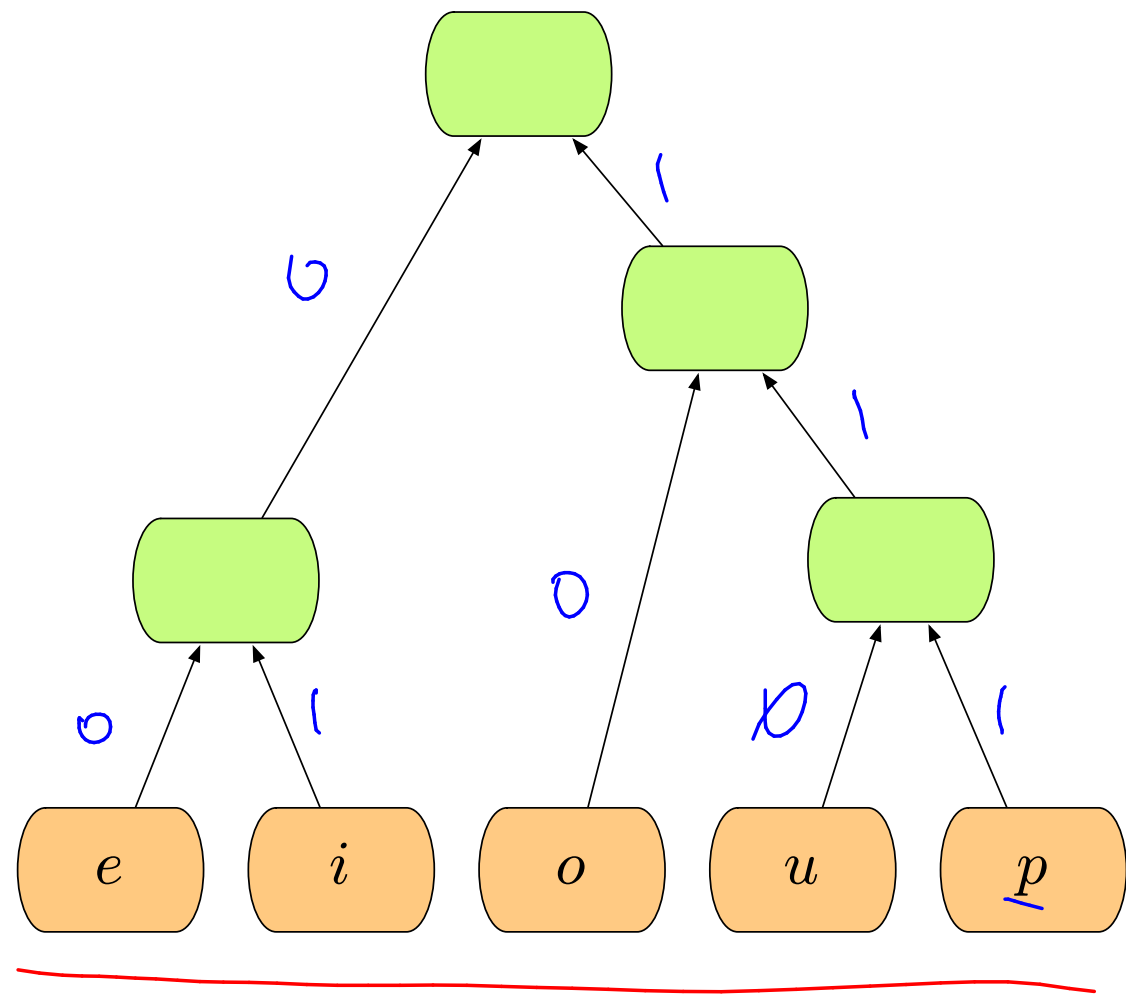


prefix code



binary tree

# use tree to encode



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

# goal

given the frequencies for some message space  $\{f_c\}_{c \in C}$   
design an optimal prefix free code.

all frequencies  
are  $> 0$ .

# goal

(all frequencies are  $> 0$ )

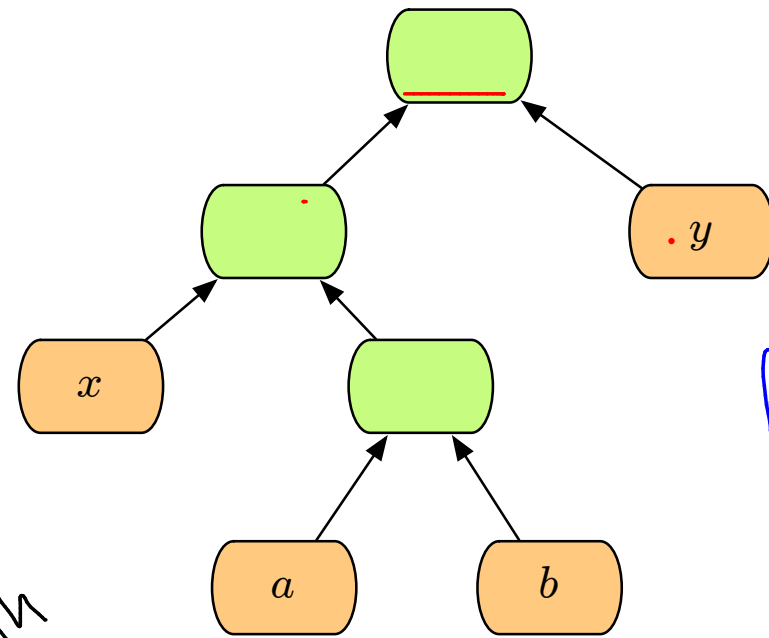
given the character frequencies  $\{f_c\}_{c \in C}$

produce a prefix code  $T$  with smallest cost

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



# property

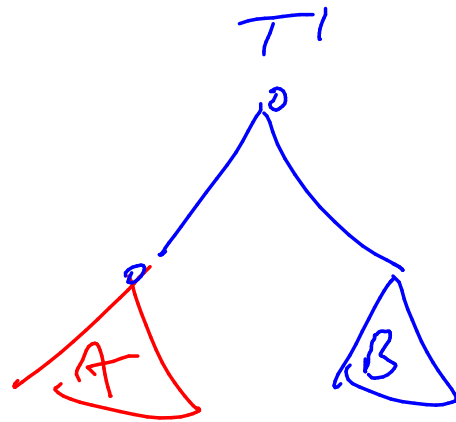
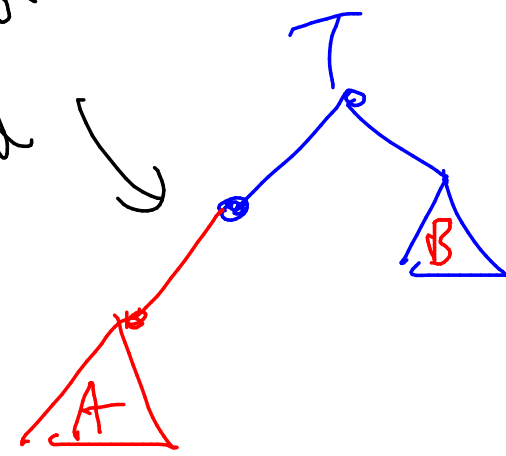


lemma: optimal tree must be full.

each node has either 0 children or 2 children

Proof: Suppose a code had a node with only 1 child. One could remove this node to produce a code with shorter codewords for several symbols

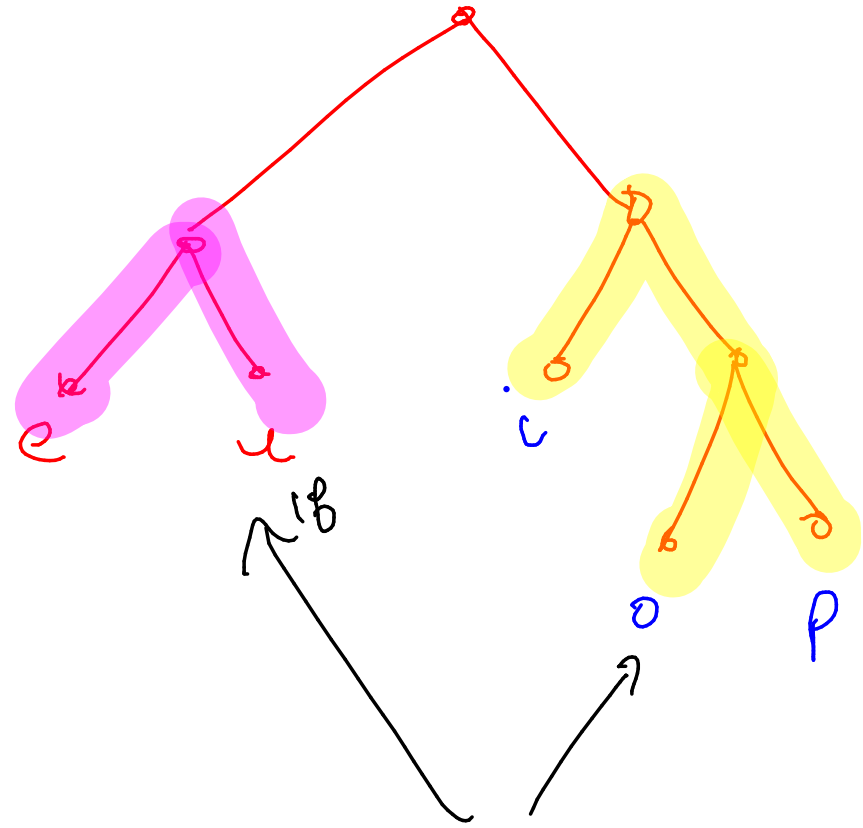
node with 1 child



divide & conquer?

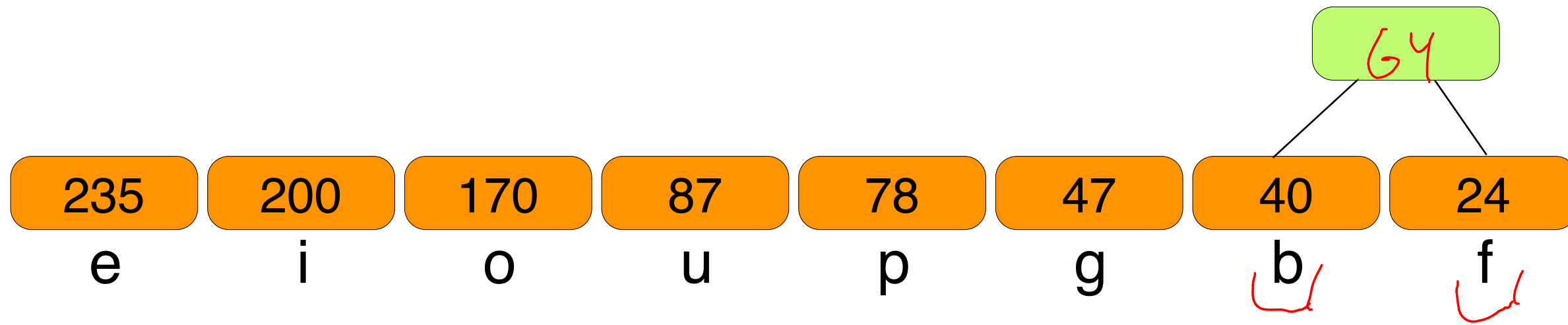
# counter-example

e: 32  
i: 25  
o: 20  
u: 18  
p: 5

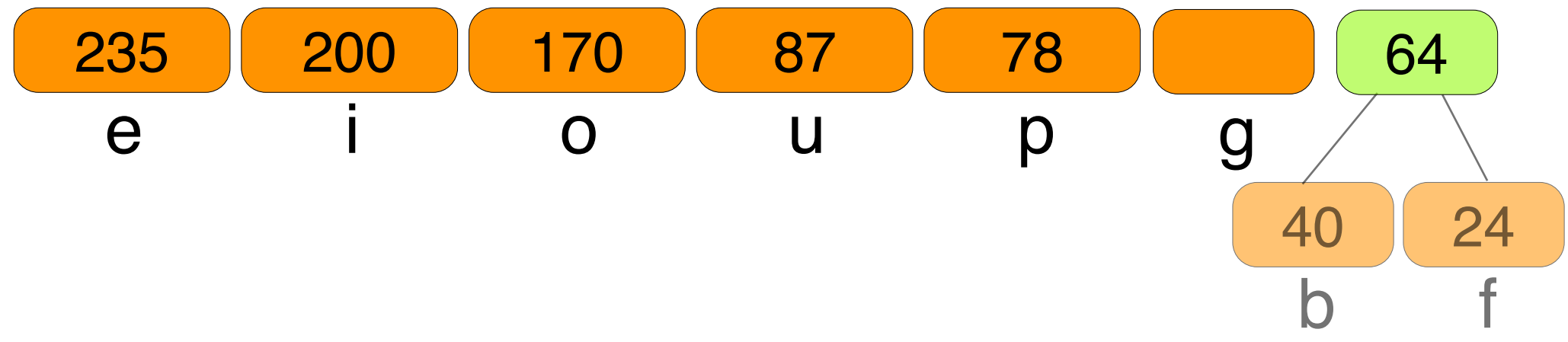


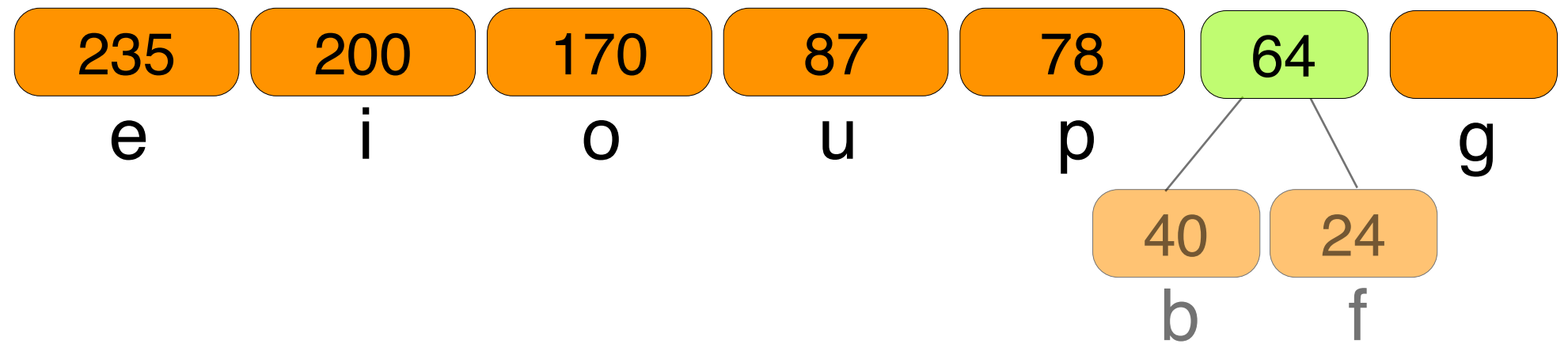
e: 2  
i: 2  
o: 3  
u: 2  
p: 3

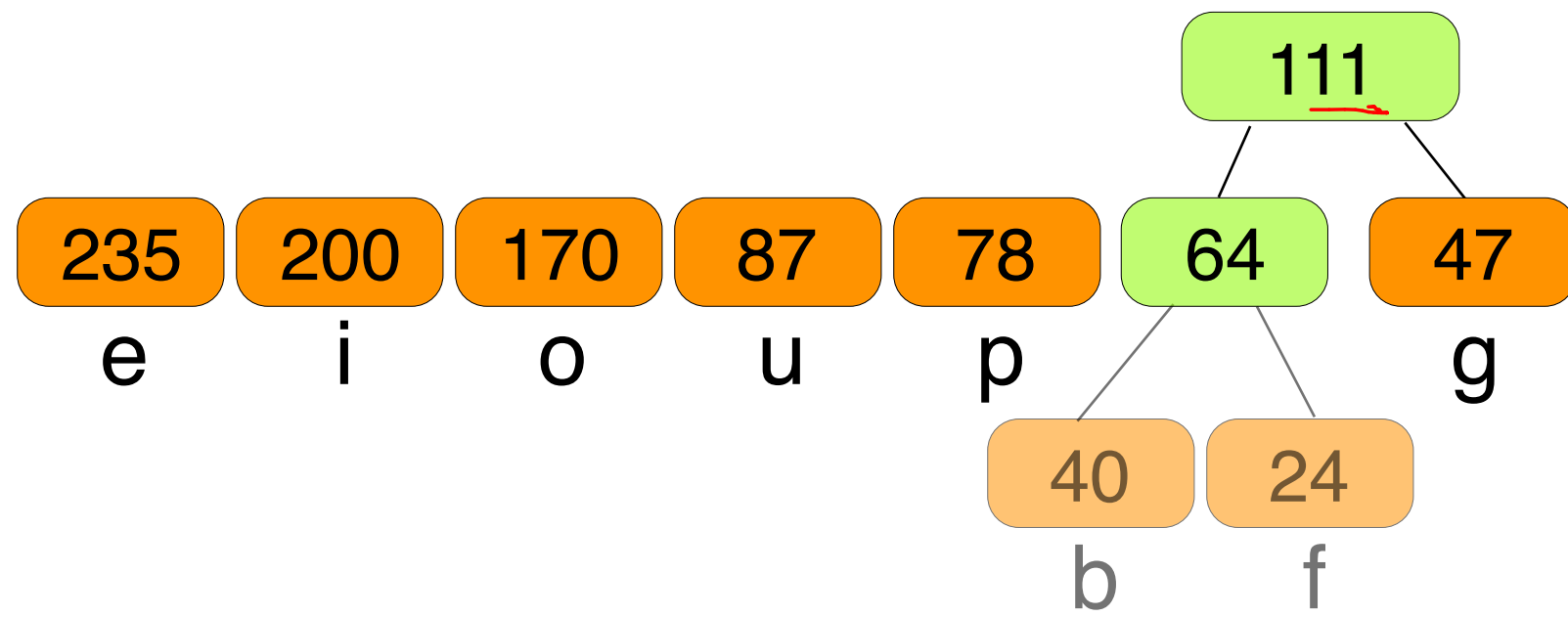
$(18 \cdot 2 + 20 \cdot 3) = 96$   
 $\hookrightarrow 20 \cdot 2 + 18 \cdot 3 = 94$

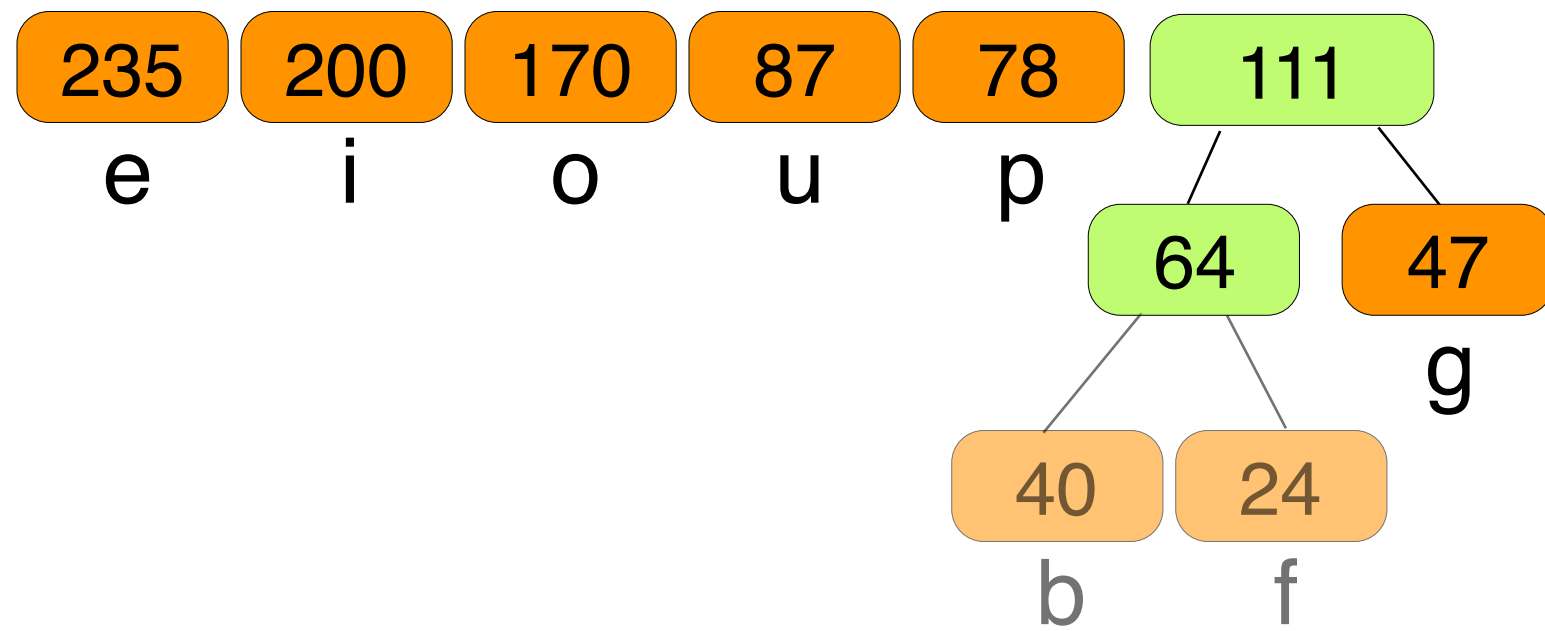


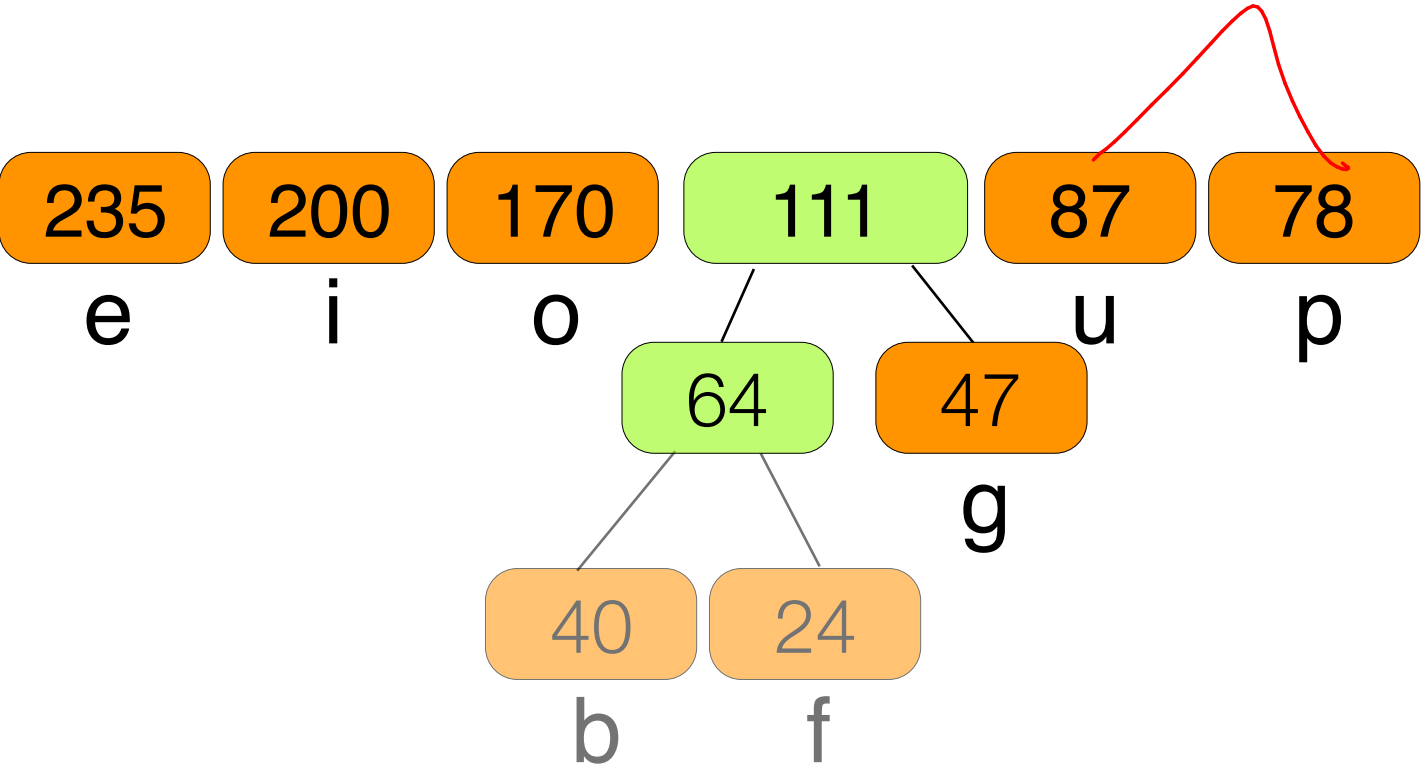


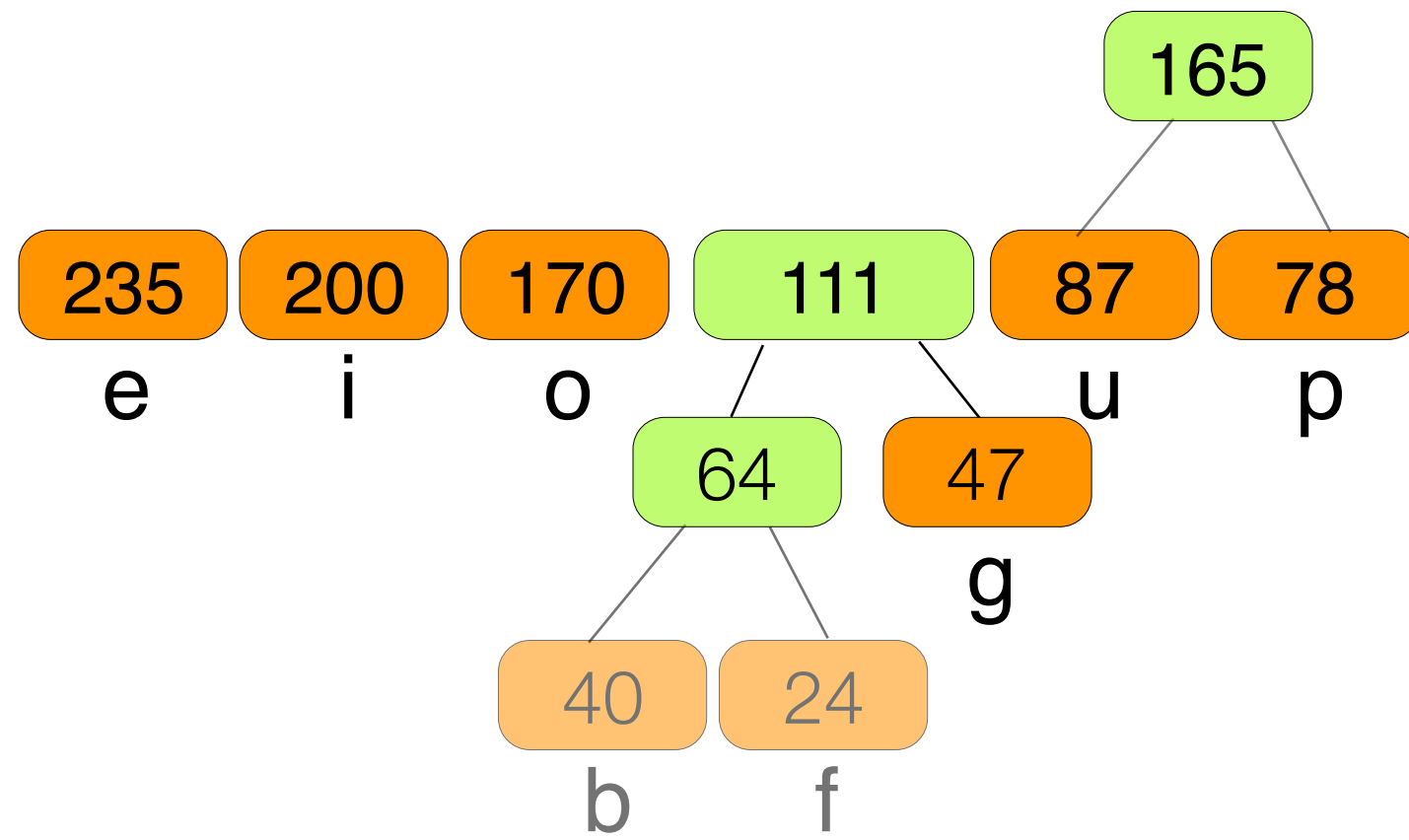




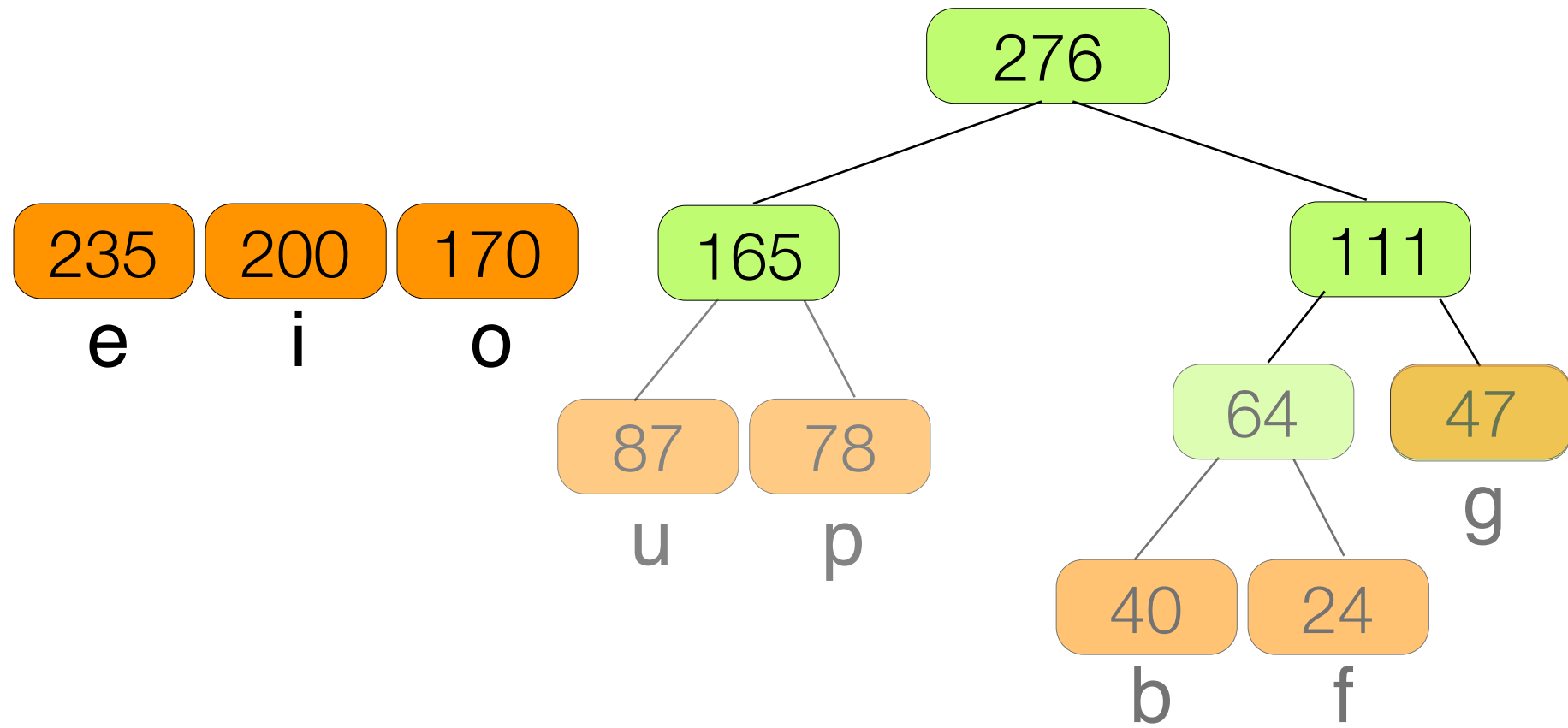


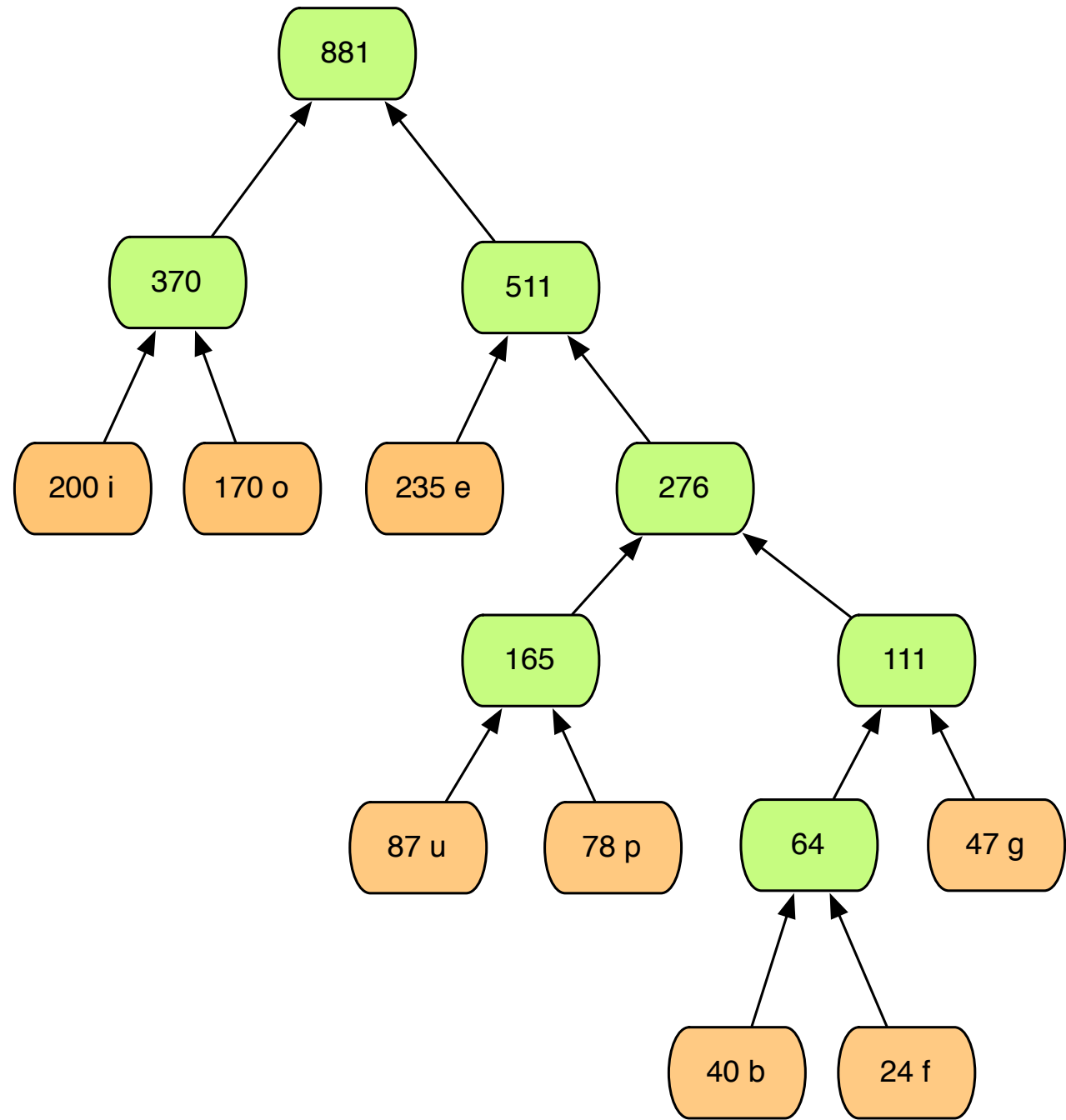


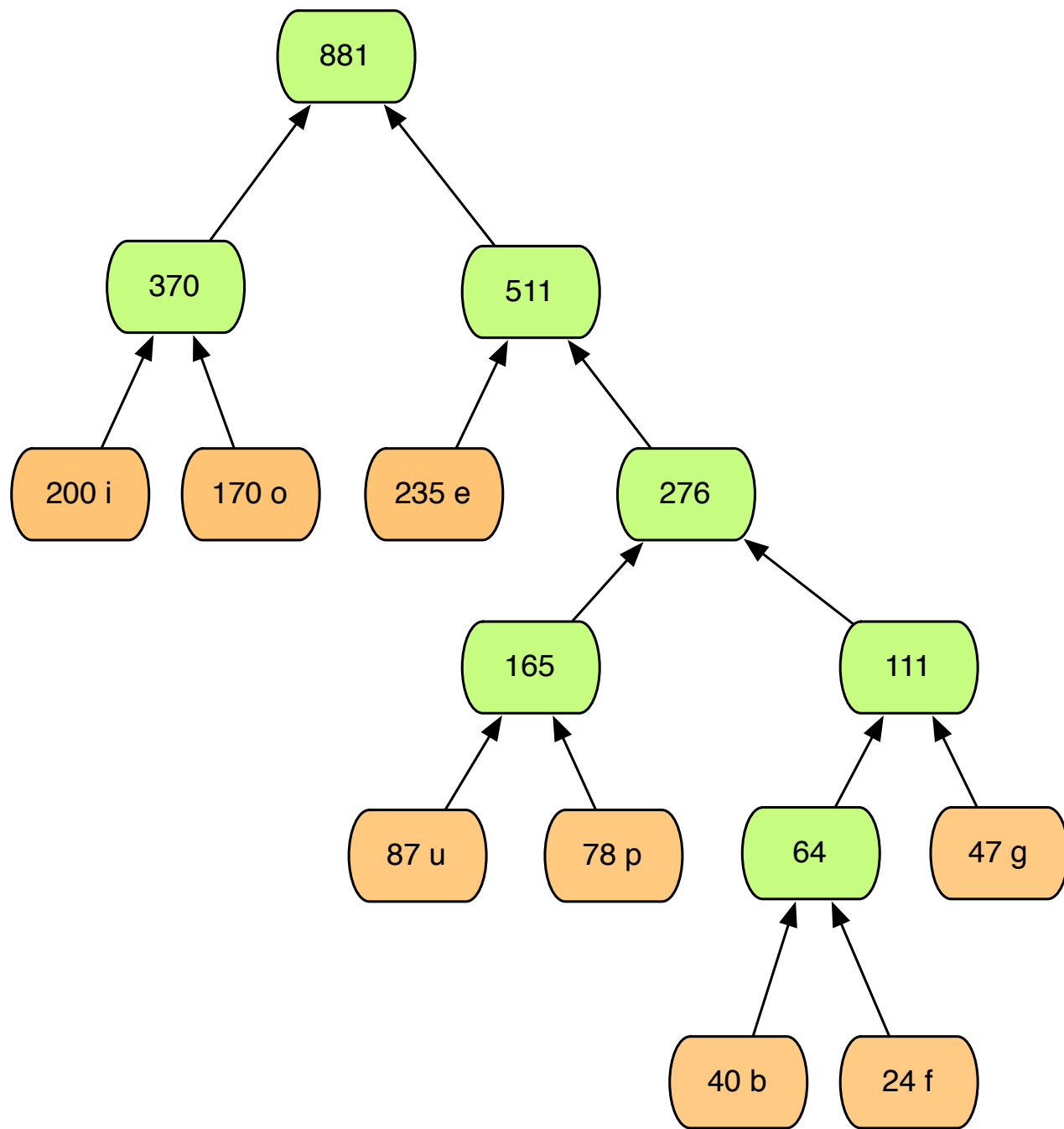




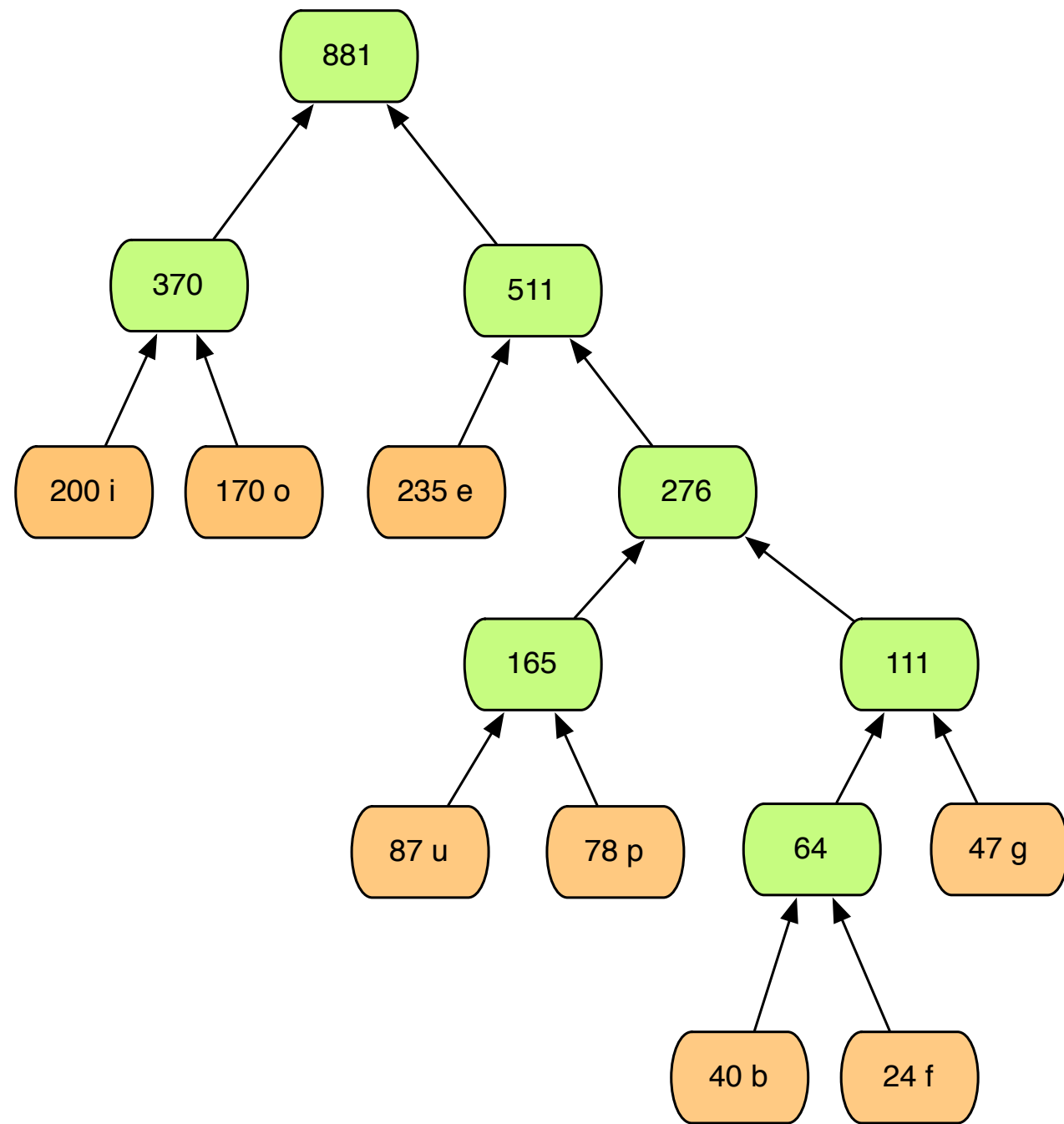








e: 235 01  
 i: 200 11  
 o: 170 10  
 u: 87 0011  
 p: 78 0010  
 g: 47 0000  
 b: 40 00011  
 f: 24 00010



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

2378

2643

10%

# objective

Prove that the Huffman algorithm produces the optimal prefix free code.

In search of an exchange argument.

# Exchange argument

lemma: Let  $f_x$  and  $f_y$  be the 2 smallest frequencies in  $\{f_c\}$ .

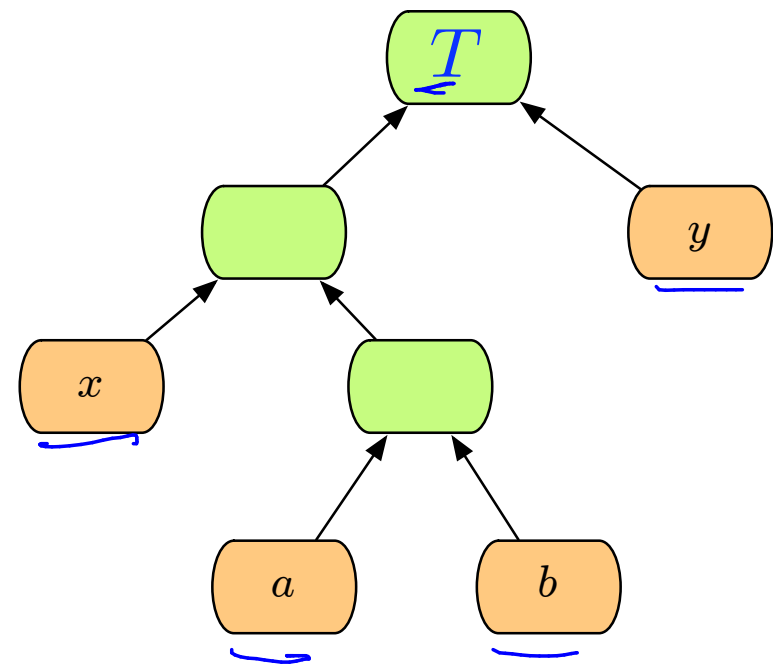
There exists an optimal code such that  $x$  and  $y$  are siblings.



# 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 optimal code for  $\{f_c\}$ .

If  $x$  and  $y$  are siblings in  $T$ , the claim holds.

Otherwise, let  $a$  and  $b$  be the 2 symbols at greatest depth in  $T$ . (which are siblings)

① Why do  $a$  &  $b$  exist??

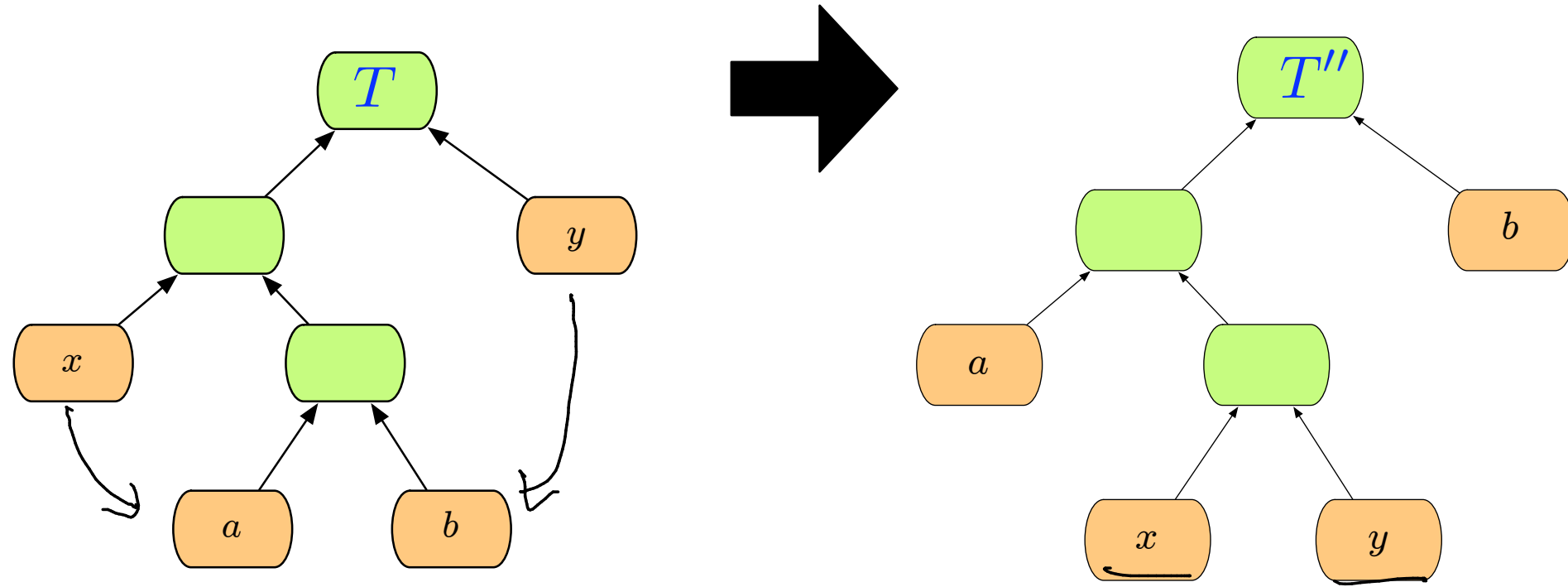
B/c  $T$  is optimal, & as we argued before,

nodes only have 0 or 2 children in optimal codes.

# 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.



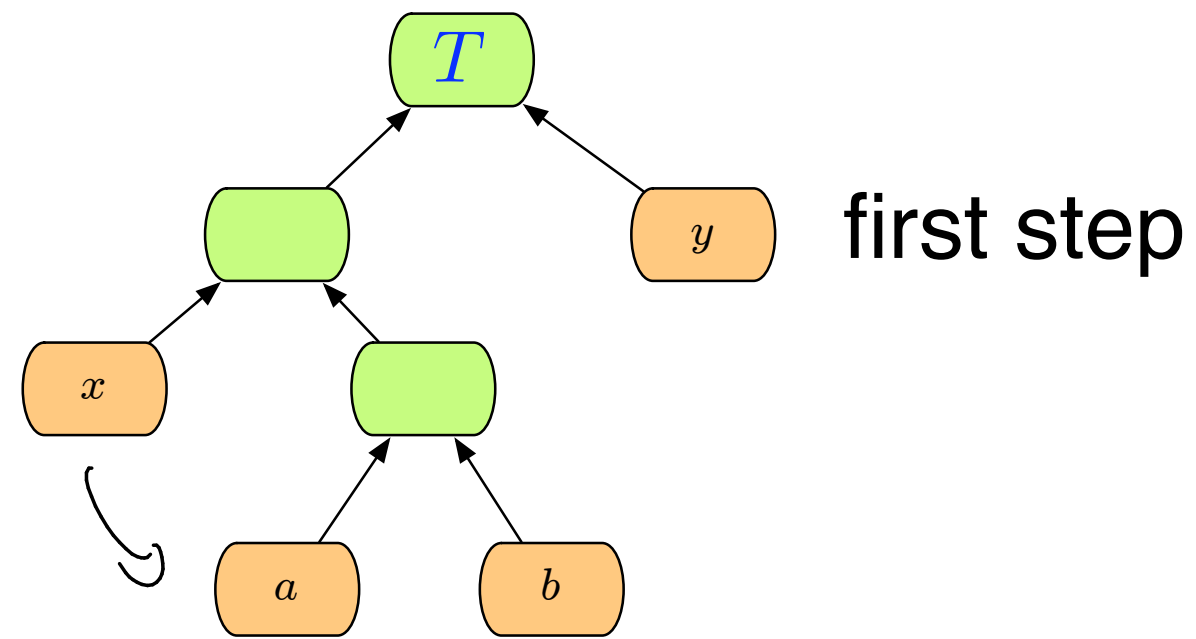
$$f_x \leq f_a \quad f_y \leq f_b$$

# exchange argument

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.

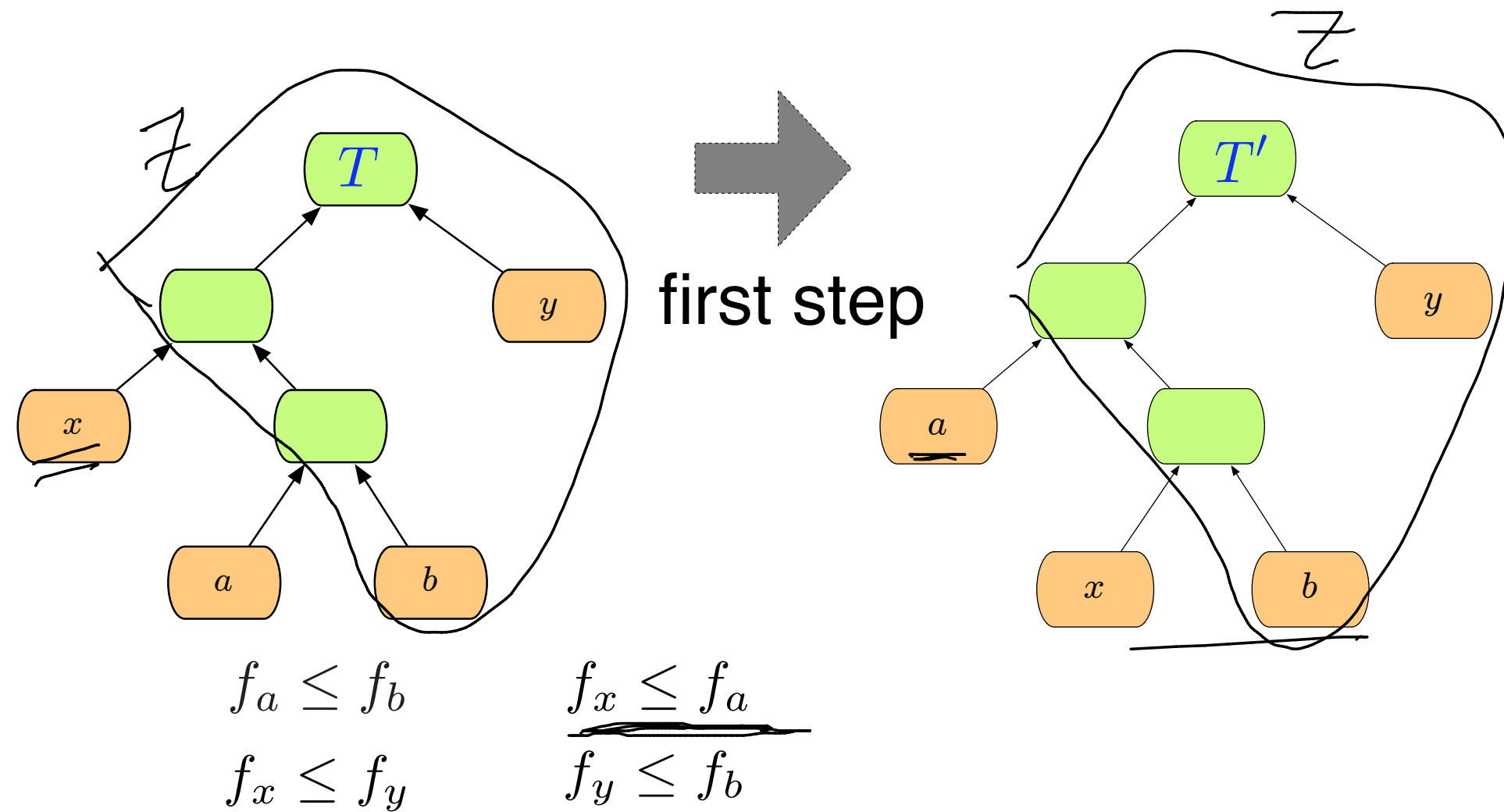
proof:

# exchange argument



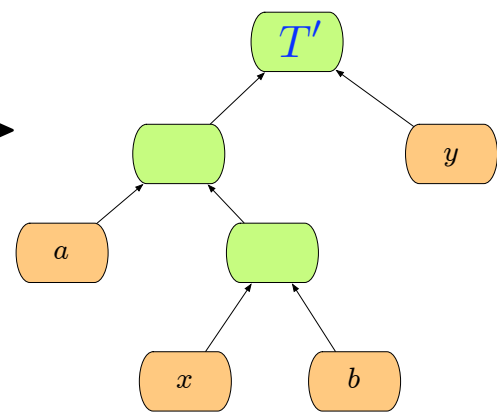
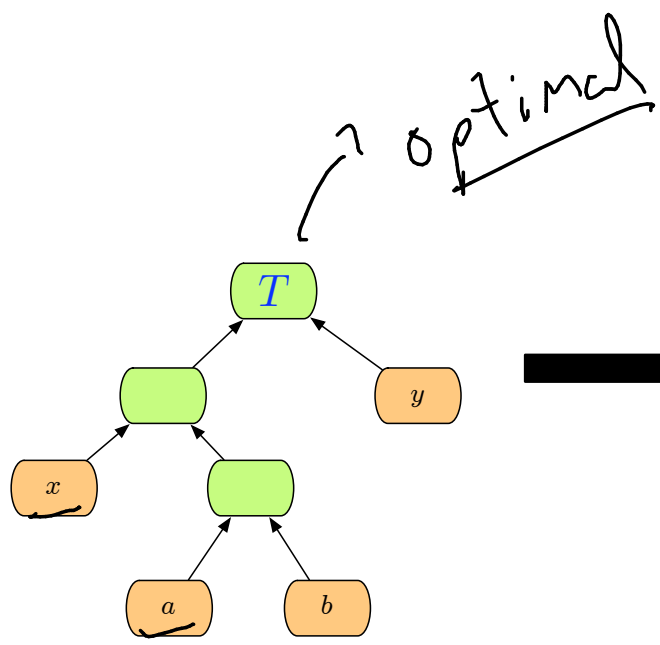
Exchange  $a$  with  $x$  in the tree to construct a new tree  $T'$ .

# exchange argument



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

$$B(T') = Z + f_a \cdot l_x + f_x \cdot l_a$$



$$B(T) = \cancel{Z} + f_x \cdot l_x + f_a \cdot l_a$$

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

$$B(T) - B(T') = f_x(l_x - l_a) - f_a(l_x - l_a) \geq 0$$

$$= \underbrace{(f_x - f_a)}_{\leq 0} \underbrace{(l_x - l_a)}_{\leq 0}$$

$$f_x \leq f_a$$

$$l_x \leq l_a$$

But  $T$  is optimal, and so  $B(T) = B(T')$

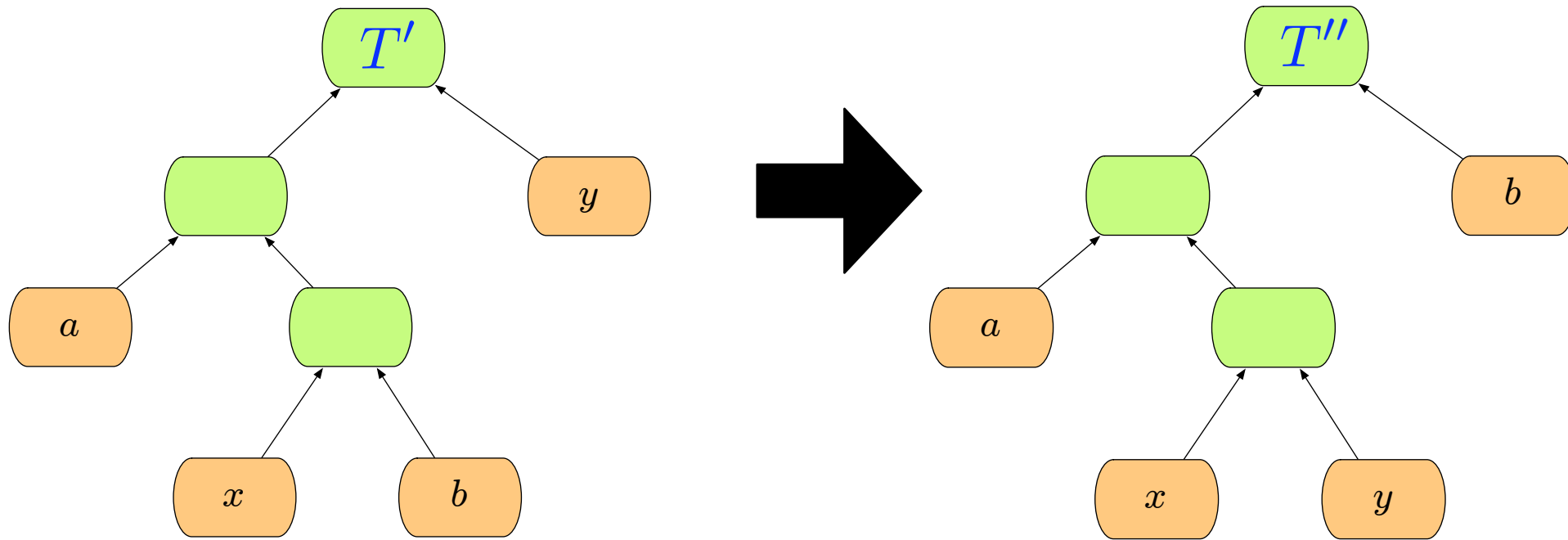




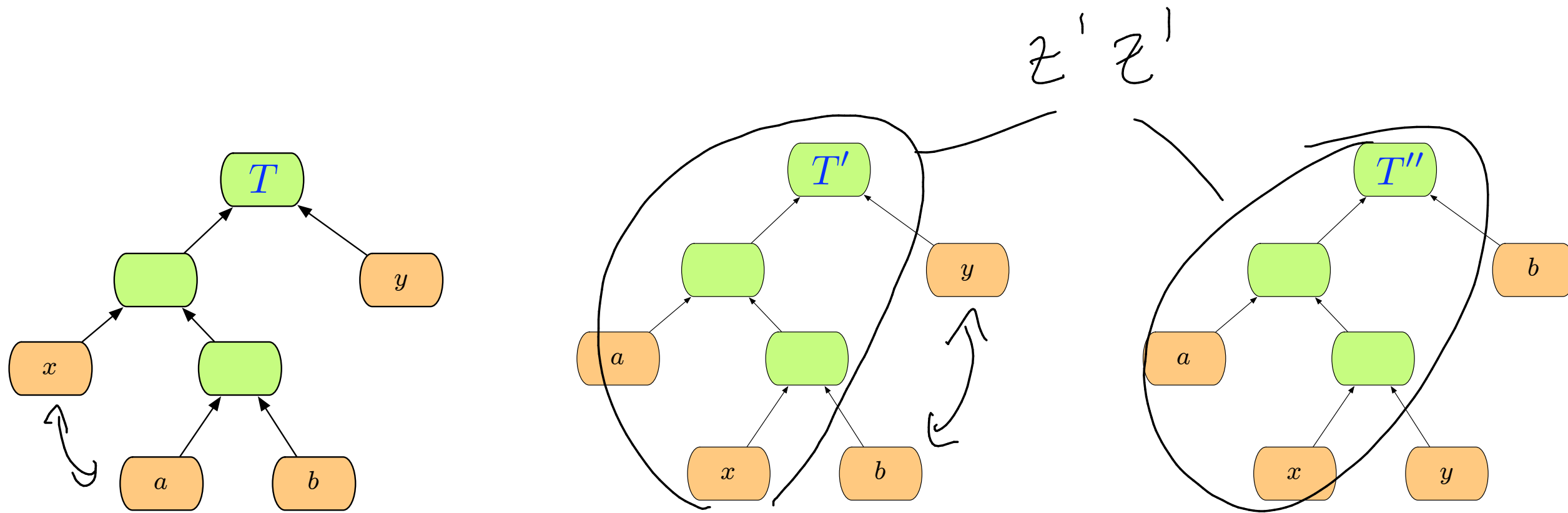
$$B(T) = \sum_c f_c l_c + f_x l_x + f_a l_a \quad B(T') = \sum_c f_c l'_c + f_x l'_x + f_a l'_a$$

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

# exchange argument



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

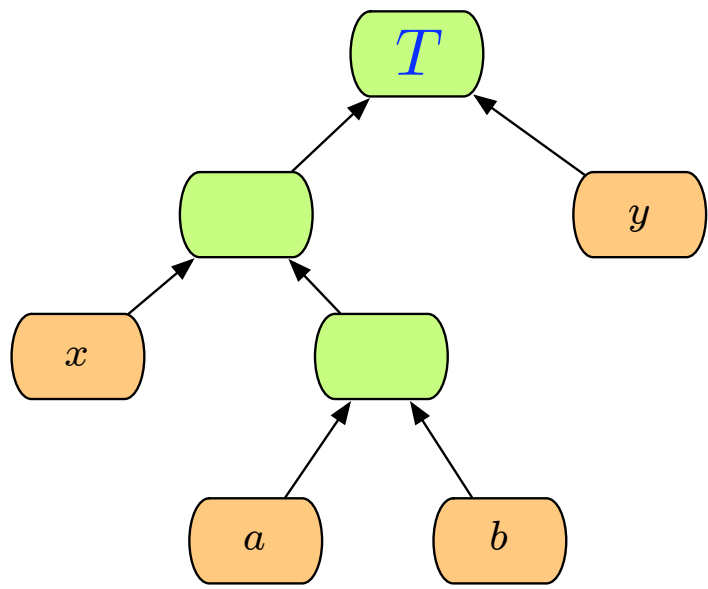


$$\underline{B(T) - B(T') \geq 0}$$

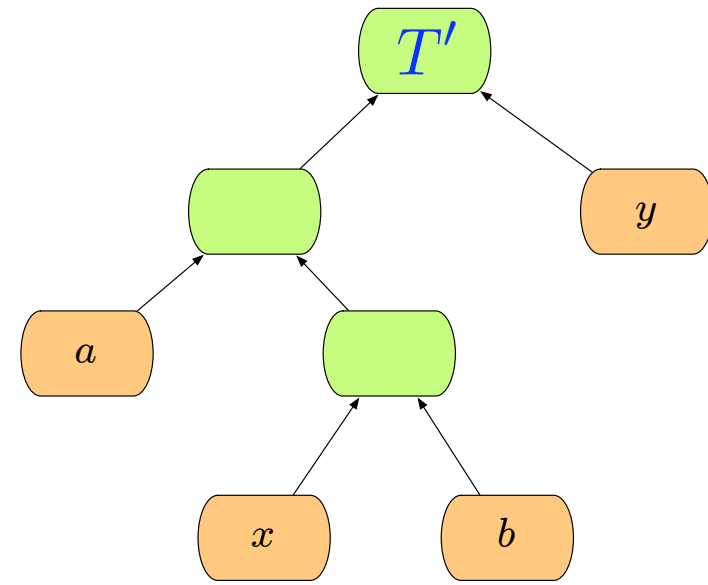
$$\underline{B(T') - B(T'') \geq 0}$$

$B(T) - B(T'') \geq 0$ , but again, b/c  $T$  is optimal,

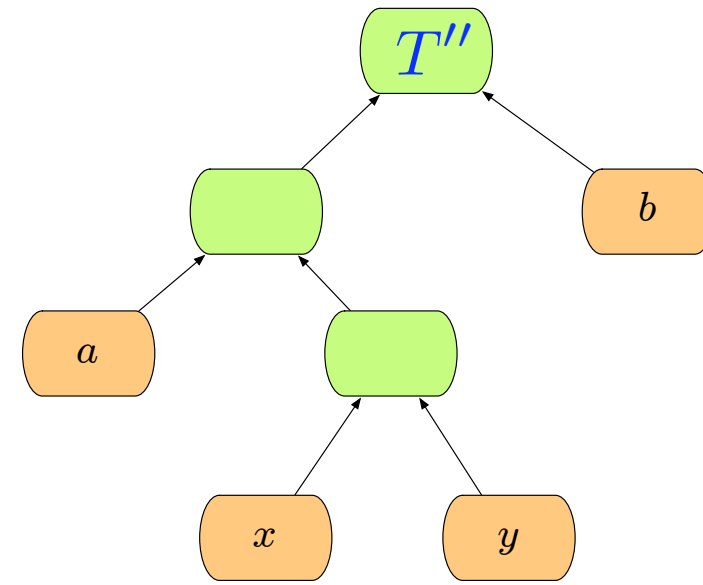
they must be equal,  $\Rightarrow T''$  is optimal.



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



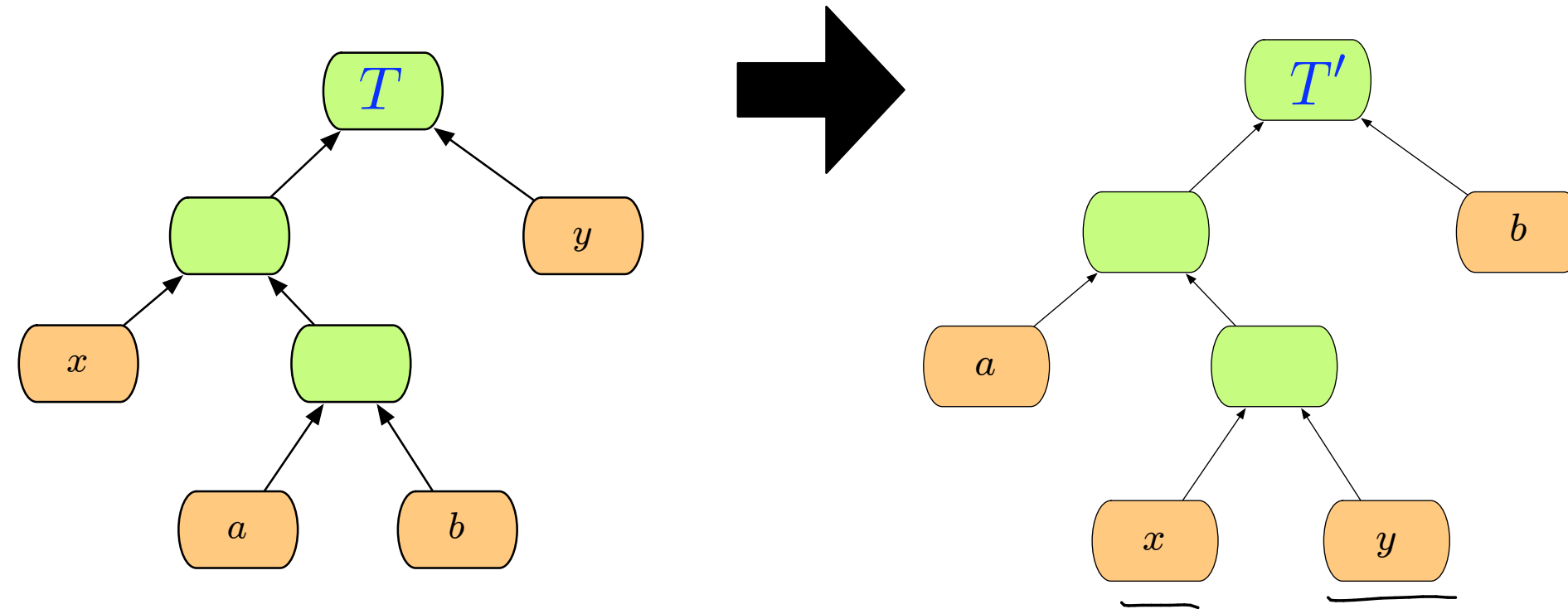
$$B(T') - B(T'') \geq 0$$



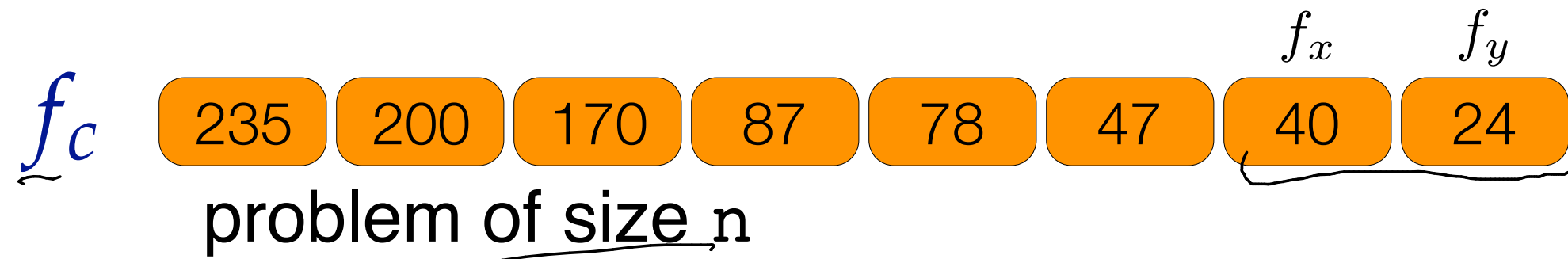
$T''$  is also optimal

# 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.

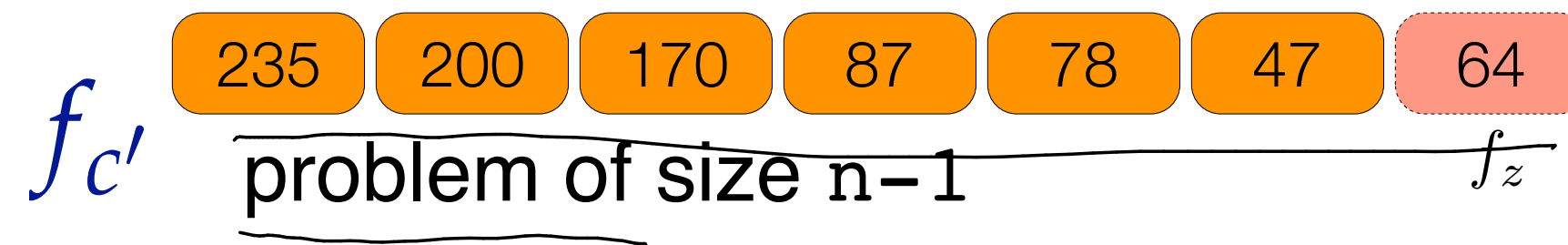
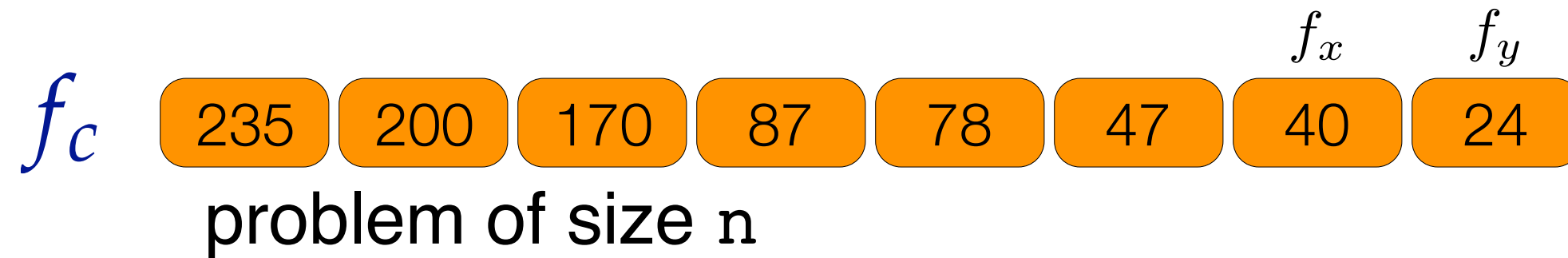


# optimal sub-structure

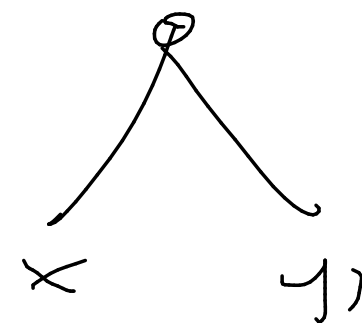




# optimal sub-structure

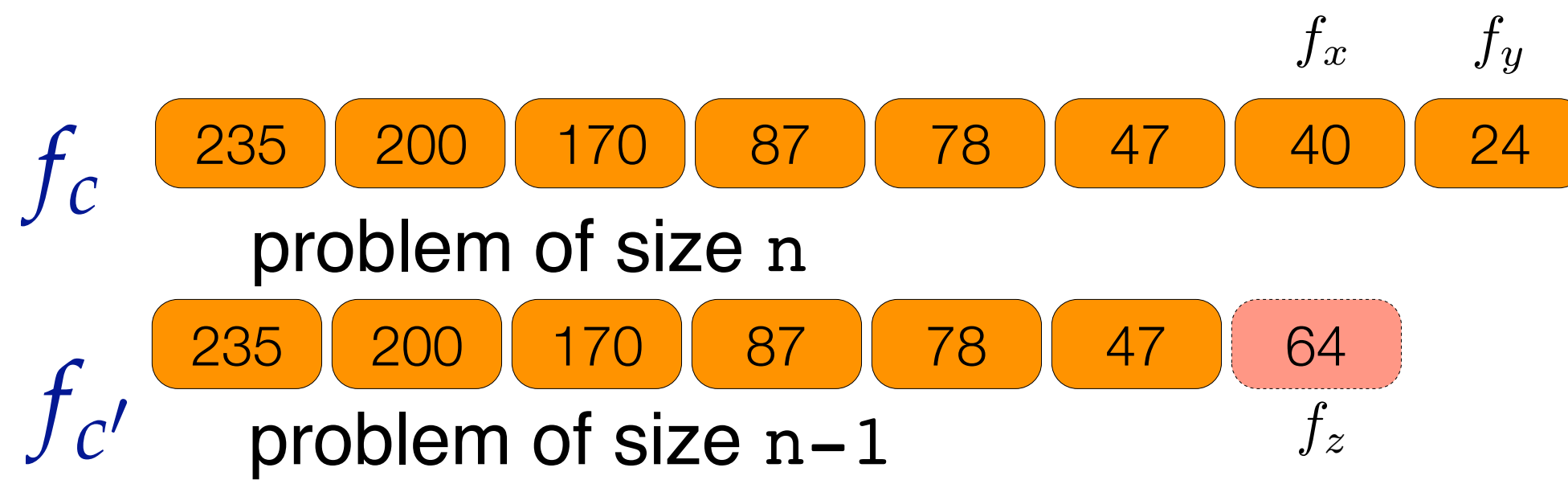


optimal solution to  $f_{c'}$ , and then replace  $f_z$  with



the resulting tree is optimal for  $f_c$ .

# optimal sub-structure



Lemma:

# optimal sub-structure

$f_c$  235 200 170 87 78 47 40 24

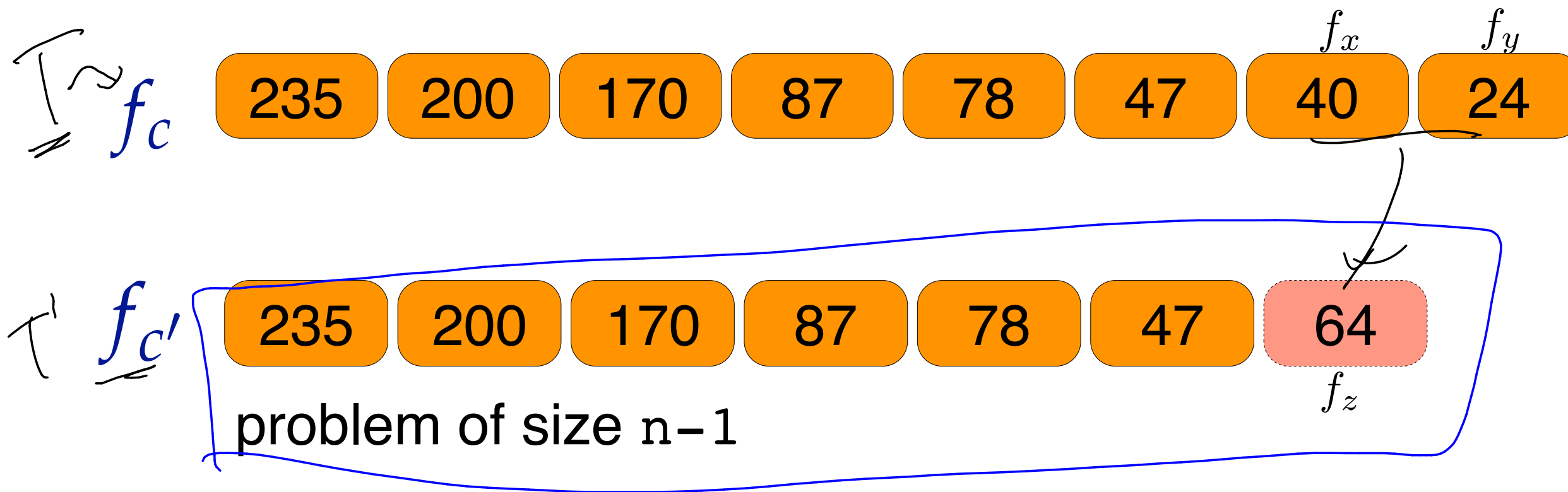
$f_x$   $f_y$

$f_{c'}$  235 200 170 87 78 47 64

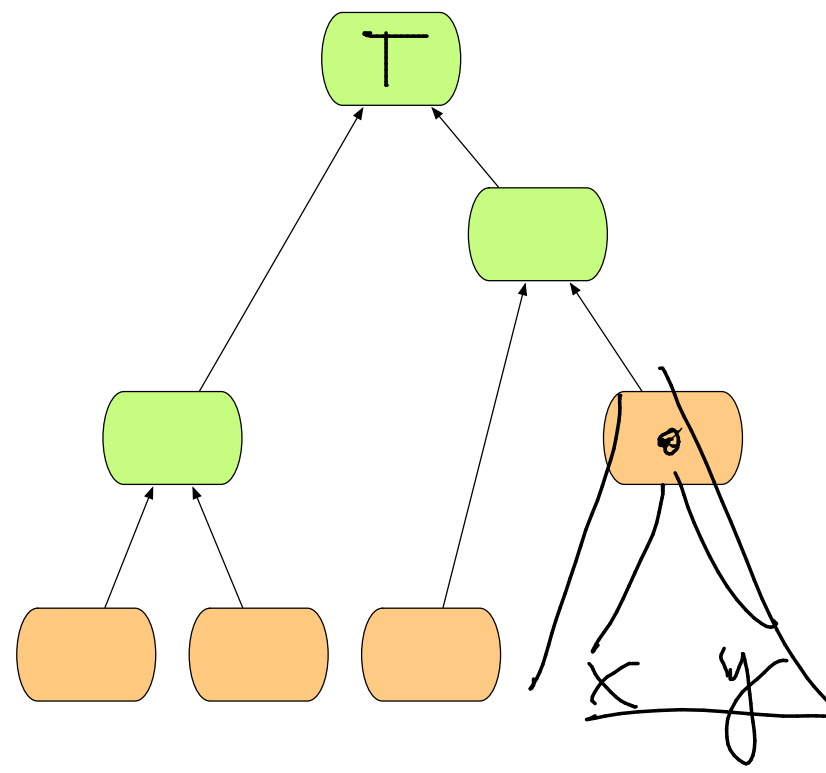
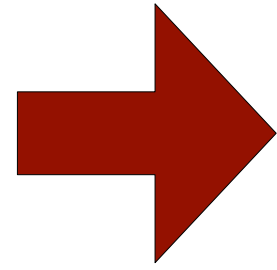
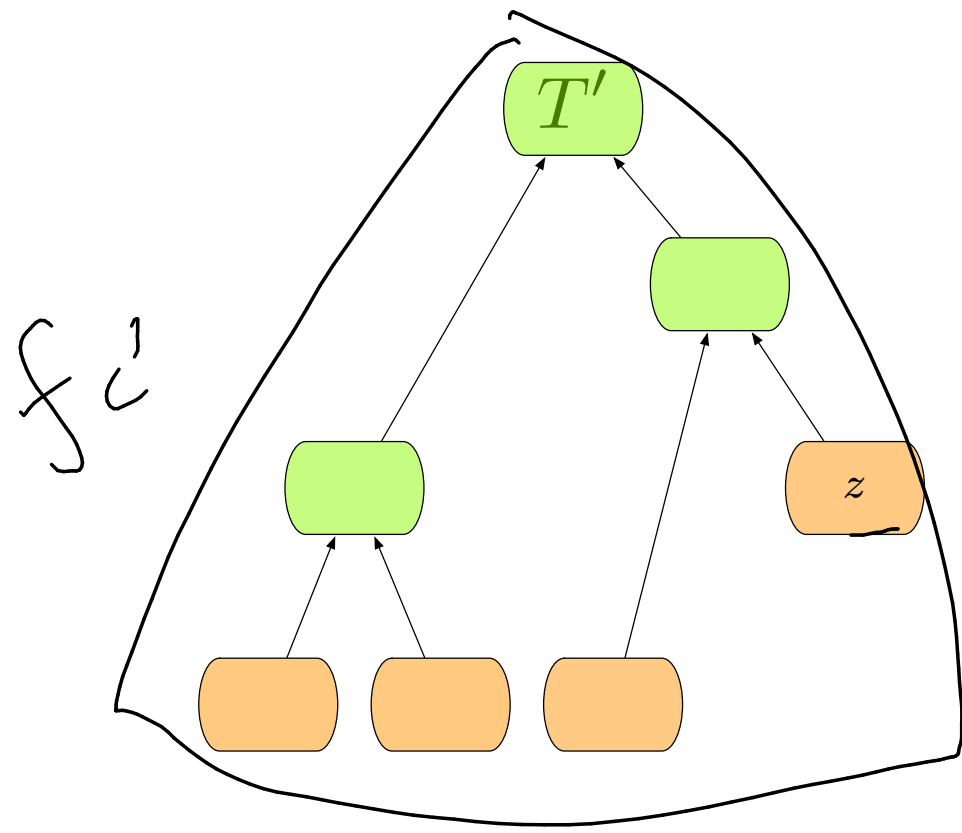
$f_z$

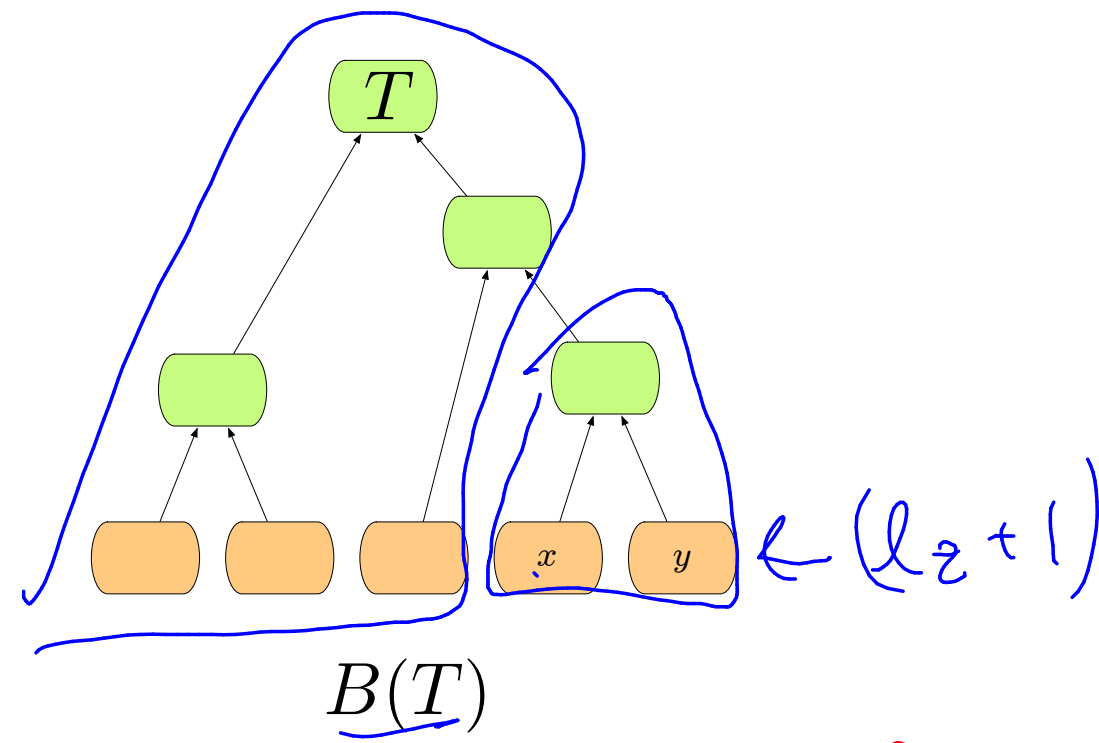
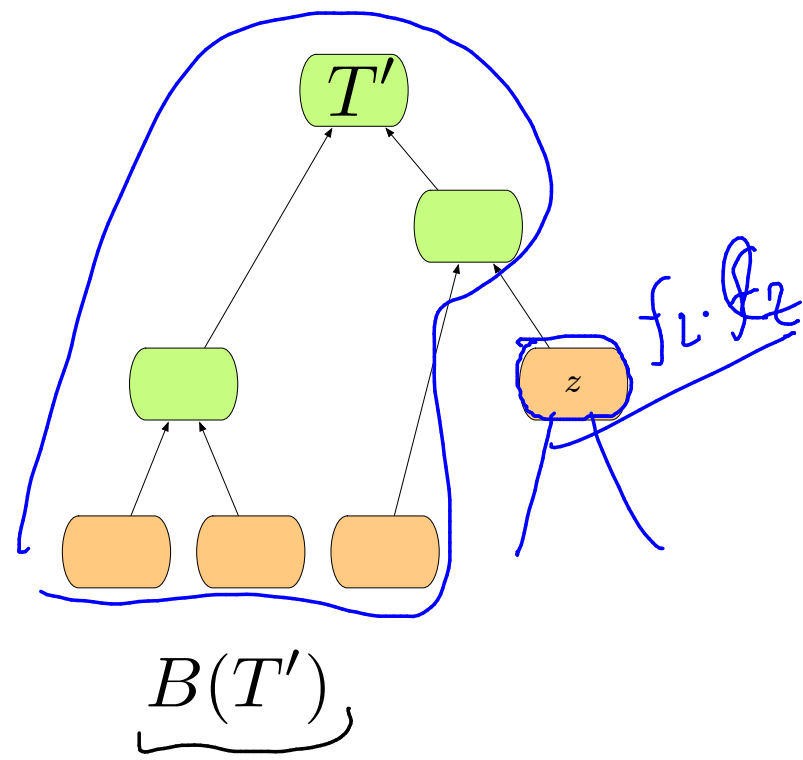
problem of size  $n-1$

# optimal sub-structure

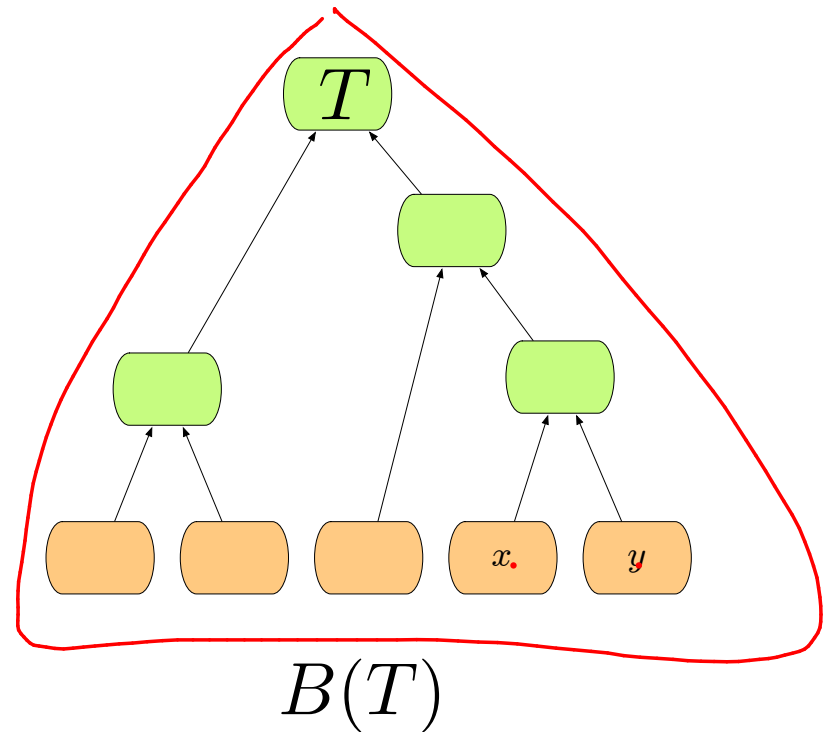
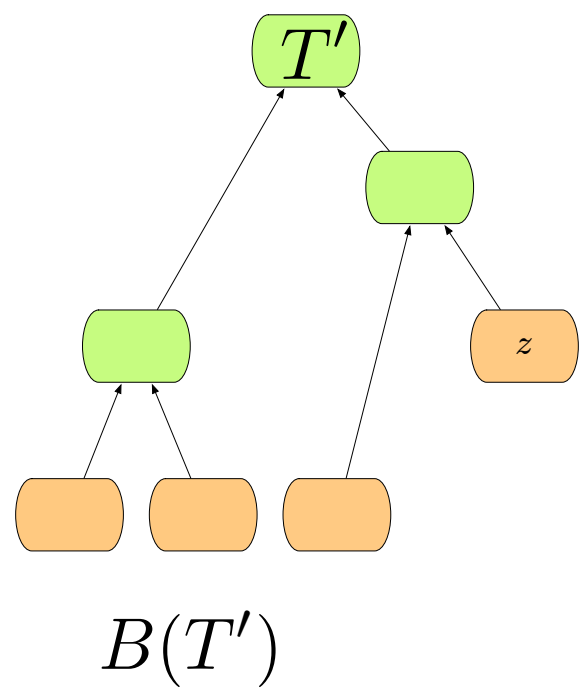


**Lemma:** The optimal solution for  $T$  consists of computing an optimal solution for  $T'$  and replacing the left  $z$  with a node having children  $x, y$ .





$$\begin{aligned}
 \underline{B(T)} &= B(T') - \underline{f_z \cdot l_z} + \underline{(l_z + 1)} \left( \overbrace{f_{x+1} f_y}^{f_z} \right) \\
 &= \underline{B(T') + \underline{f_{x+1} f_y}}
 \end{aligned}$$



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

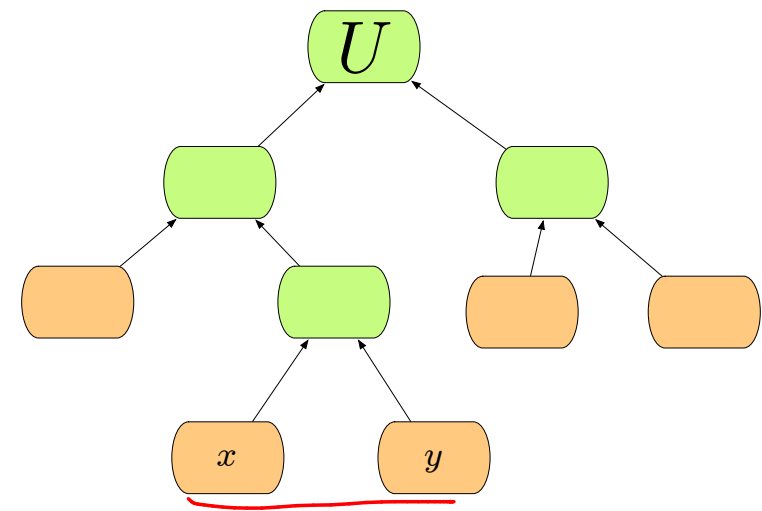
→

Suppose  $T$  is not optimal.



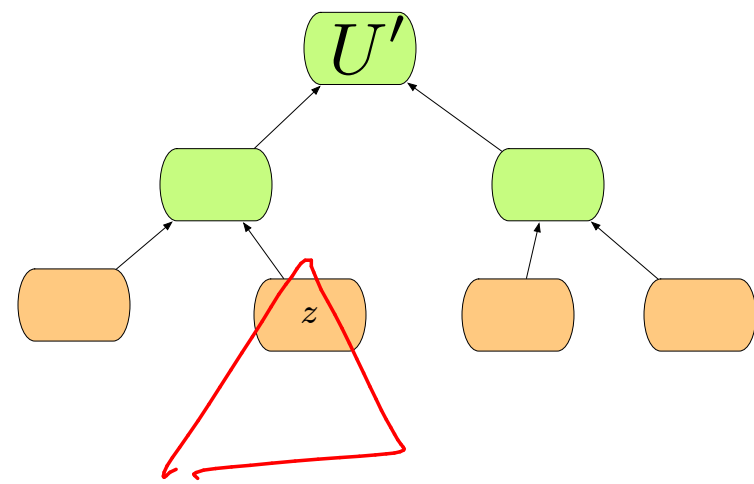
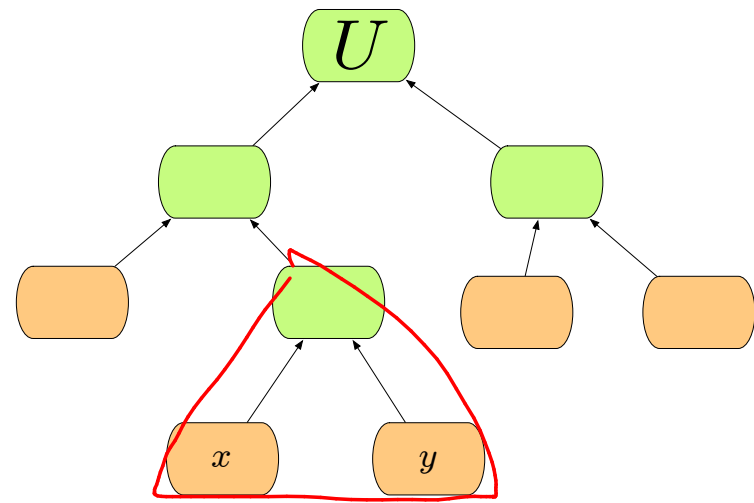
Suppose  $T$  is not optimal.

There is some other tree  $U$  such that



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

Suppose  $T$  is not optimal.



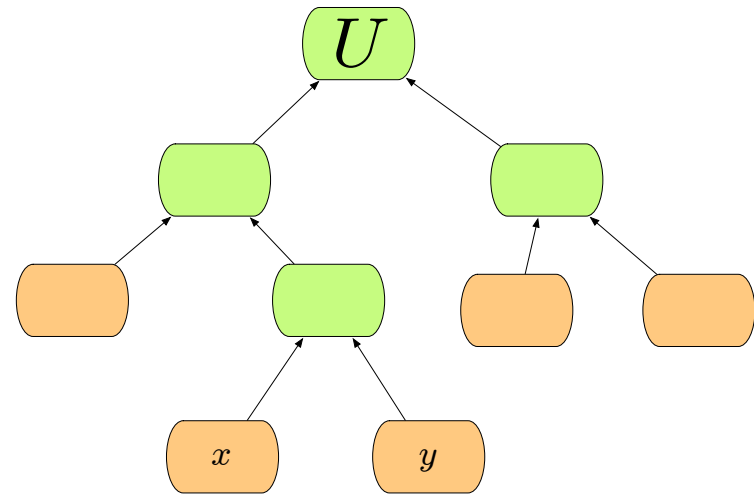
$$\underline{B(U)} < B(T) = \underline{B(T') + f_x + f_y}$$

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

$$\underline{B(u')} < B(T')$$

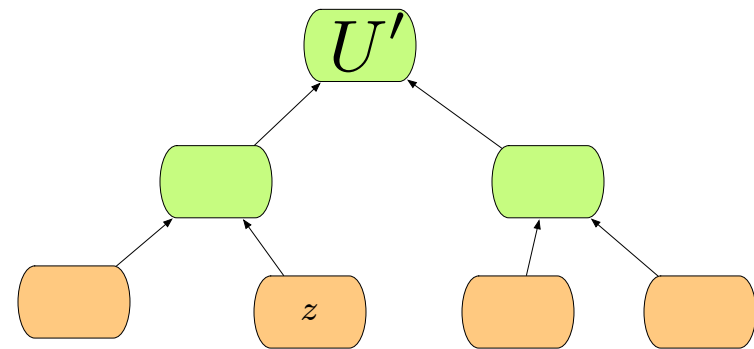
This suggests that  $T'$  was not an optimal solution. This is a contradiction.

# 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

