

Huffman character encoding

CMSC 420

You bigoted UTF-8 you!

- “Archimedes” will take 10 bytes to store...
- “Αρχιμήδης” will take >> 10 bytes to store!

You bigoted UTF-8 you!

- “Archimedes” will take 10 bytes to store...
- “Αρχιμήδης” will take >> 10 bytes to store!
- But with the vast majority of text being in English... people have adapted!
 - And the Earth keeps on turning.



You bigoted UTF-8 you!

- “Archimedes” will take 10 bytes to store...
- “Αρχιμήδης” will take >> 10 bytes to store!
- But with the vast majority of text being in English... people have adapted!
 - And the Earth keeps on turning.
- UTF8 is a variable-length encoding that has taken advantage of the frequency of English communications on the Web to present an overall very efficient way to store character data.



Quiz

- Is **ASCII** a variable – length encoding?

Yes

No

Quiz

- Is **ASCII** a variable – length encoding?

Yes

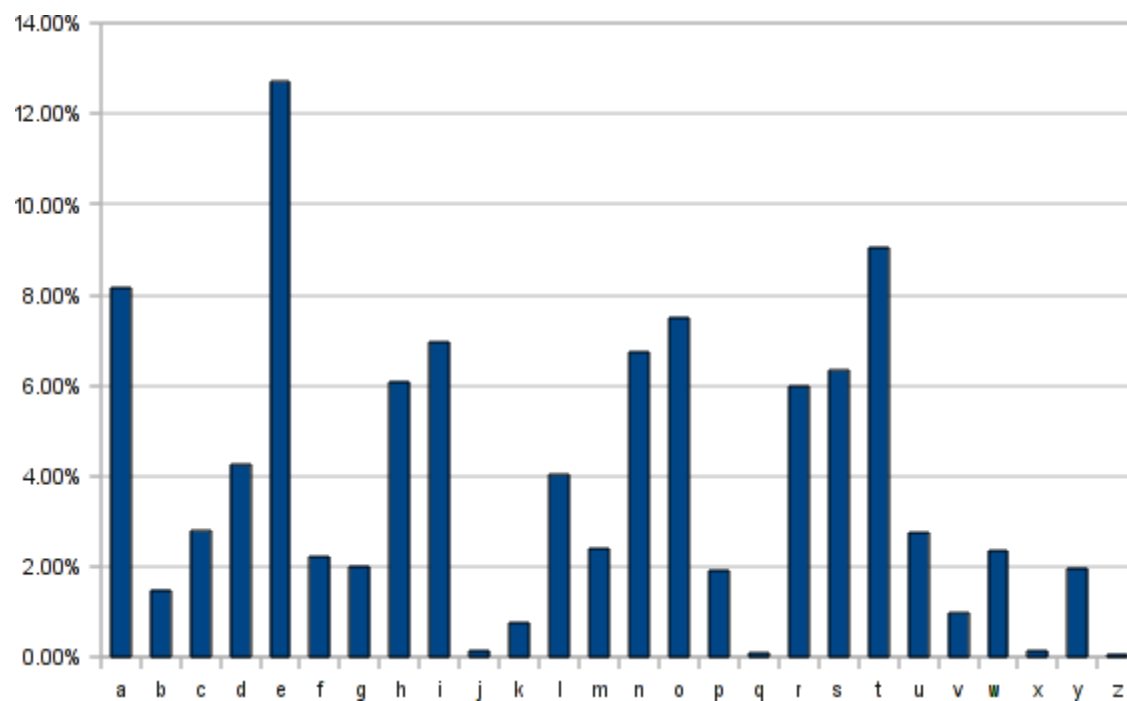
No

Dec	Hex	Oct	Char	Dec	Hex	Oct	Char	Dec	Hex	Oct	Char	Dec	Hex	Oct	Char
0	0	000	NUL (null)	32	20	040	#32; Space	64	40	100	#64; @	96	60	140	#96; P
1	1	001	SOH (start of heading)	33	21	041	#33; !	65	41	101	#65; A	97	61	141	#97; Q
2	2	002	STX (start of text)	34	22	042	#34; "	66	42	102	#66; B	98	62	142	#98; R
3	3	003	ETX (end of text)	35	23	043	#35; #	67	43	103	#67; C	99	63	143	#99; S
4	4	004	EOF (end of transmission)	36	24	044	#36; \$	68	44	104	#68; D	100	64	144	#100; T
5	5	005	ENQ (enquiry)	37	25	045	#37; %	69	45	105	#69; E	101	65	145	#101; U
6	6	006	ACK (acknowledge)	38	26	046	#38; &	70	46	106	#70; F	102	66	146	#102; V
7	7	007	BEL (bell)	39	27	047	#39; '	71	47	107	#71; G	103	67	147	#103; W
8	8	010	BS (backspace)	40	28	050	#40; (72	48	110	#72; H	104	68	150	#104; X
9	9	011	TAB (horizontal tab)	41	29	051	#41;)	73	49	111	#73; I	105	69	151	#105; Y
10	A	012	LF (NL line feed, new line)	42	2A	052	#42; *	74	4A	112	#74; J	106	6A	152	#106; Z
11	B	013	VT (vertical tab)	43	2B	053	#43; +	75	4B	113	#75; K	107	6B	153	#107; [
12	C	014	FF (NP form feed, new page)	44	2C	054	#44; ,	76	4C	114	#76; L	108	6C	154	#108; \
13	D	015	CR (carriage return)	45	2D	055	#45; -	77	4D	115	#77; M	109	6D	155	#109;]
14	E	016	SO (shift out)	46	2E	056	#46; .	78	4E	116	#78; N	110	6E	156	#110; ^
15	F	017	SI (shift in)	47	2F	057	#47; /	79	4F	117	#79; O	111	6F	157	#111; _
16	10	020	DLE (data link escape)	48	30	060	#48; 0	80	50	120	#80; P	112	70	160	#112; `
17	11	021	DC1 (device control 1)	49	31	061	#49; 1	81	51	121	#81; Q	113	71	161	#113; a
18	12	022	DC2 (device control 2)	50	32	062	#50; 2	82	52	122	#82; R	114	72	162	#114; b
19	13	023	DC3 (device control 3)	51	33	063	#51; 3	83	53	123	#83; S	115	73	163	#115; c
20	14	024	DC4 (device control 4)	52	34	064	#52; 4	84	54	124	#84; T	116	74	164	#116; d
21	15	025	NAK (negative acknowledge)	53	35	065	#53; 5	85	55	125	#85; U	117	75	165	#117; e
22	16	026	SYN (synchronous idle)	54	36	066	#54; 6	86	56	126	#86; V	118	76	166	#118; f
23	17	027	ETB (end of trans. block)	55	37	067	#55; 7	87	57	127	#87; W	119	77	167	#119; g
24	18	030	CAN (cancel)	56	38	070	#56; 8	88	58	130	#88; X	120	78	170	#120; h
25	19	031	EM (end of medium)	57	39	071	#57; 9	89	59	131	#89; Y	121	79	171	#121; i
26	1A	032	SUB (substitute)	58	3A	072	#58; :	90	5A	132	#90; Z	122	7A	172	#122; j
27	1B	033	ESC (escape)	59	3B	073	#59; ;	91	5B	133	#91; [123	7B	173	#123; k
28	1C	034	FS (file separator)	60	3C	074	#60; <	92	5C	134	#92; \	124	7C	174	#124; l
29	1D	035	GS (group separator)	61	3D	075	#61; =	93	5D	135	#93;]	125	7D	175	#125; m
30	1E	036	RS (record separator)	62	3E	076	#62; >	94	5E	136	#94; ^	126	7E	176	#126; n
31	1F	037	US (unit separator)	63	3F	077	#63; ?	95	5F	137	#95; _	127	7F	177	#127; o

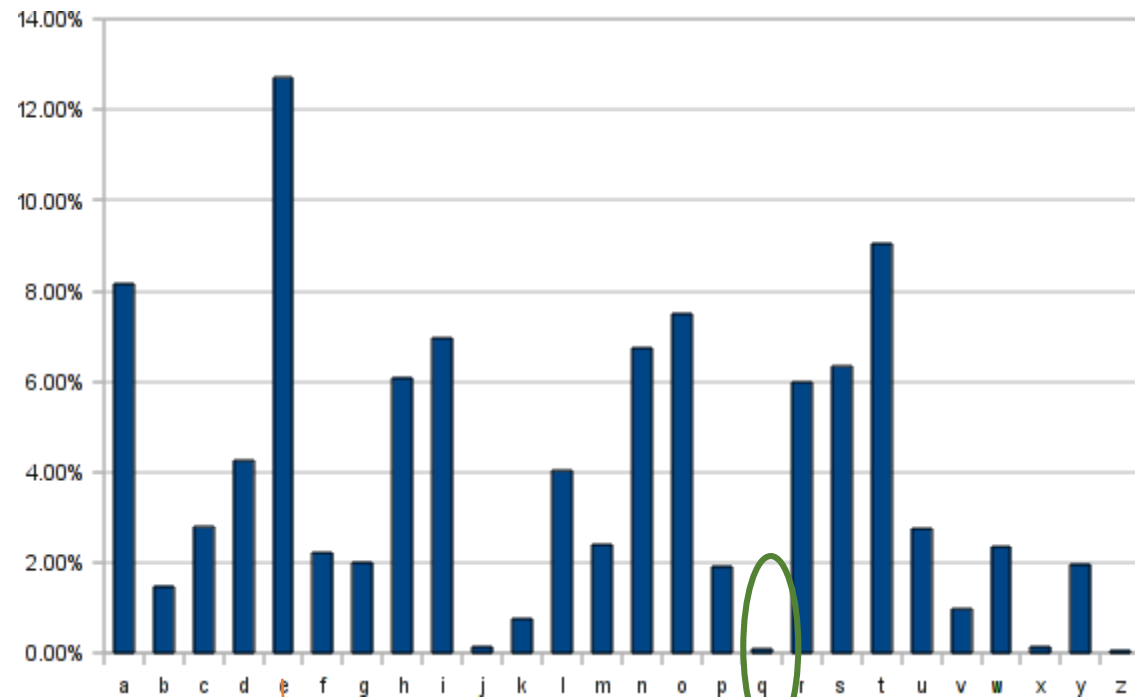
Source: www.LookupTables.com

- Remember: **7 bits per character**.
- Can represent 128 characters, with about 30 of them non-printable.
- But still, a remarkably useful subset, with 51.6% of the Web written in English, and lots of those characters including:
 - Punctuation (, , . , ! , ...) almost universally used
 - MANY characters that other Romance languages, like Spanish and French, use (a , s , t , r , ...)
 - ALL different kinds of whitespace ('\0' , tabs , space , CRLF , ...)!

A different kind of redundancy...



A different kind of redundancy...

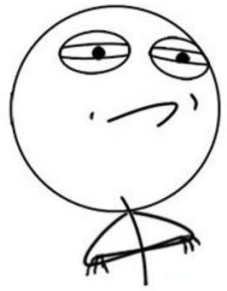


Encoding
cost: 7 bits

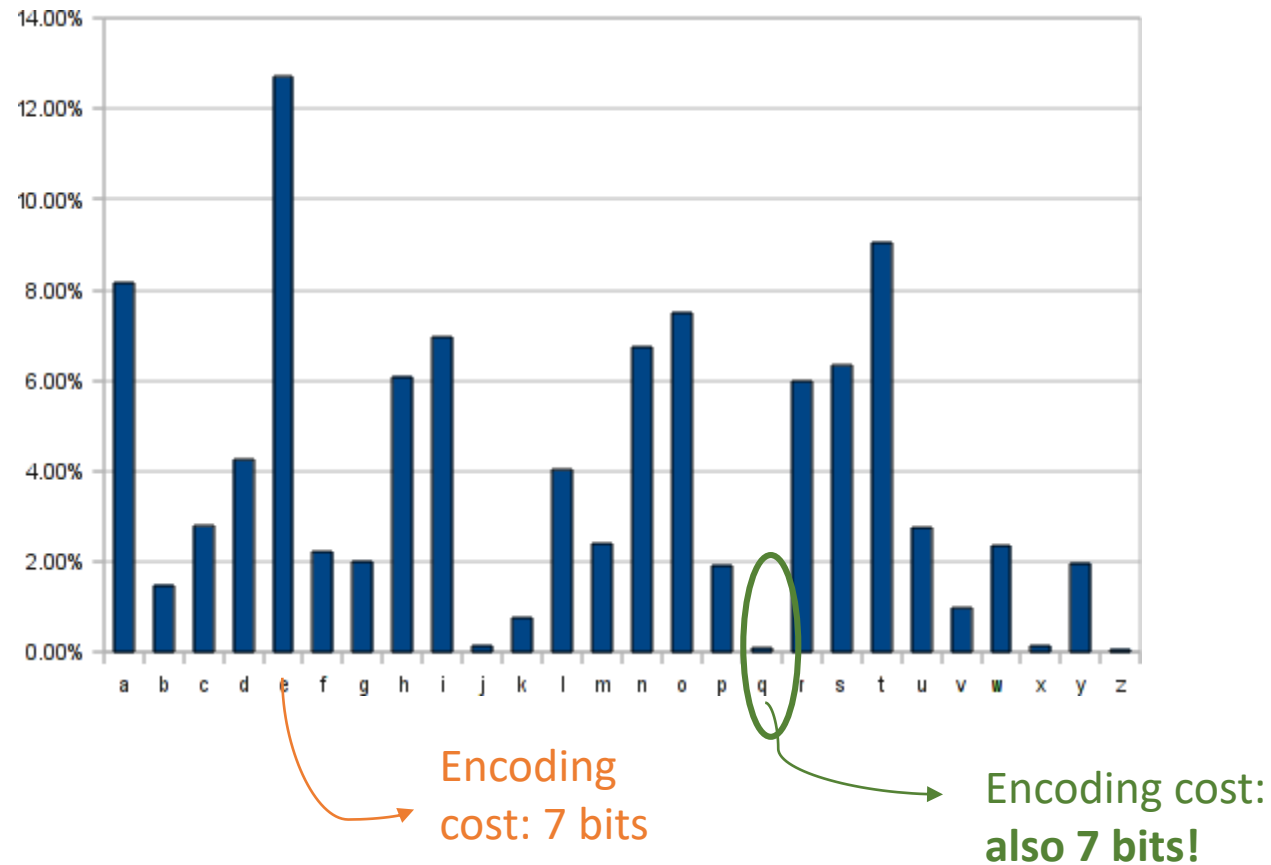
Encoding cost:
also 7 bits!

This doesn't
seem fair!

A different kind of redundancy...



We should do better.



This doesn't seem fair!

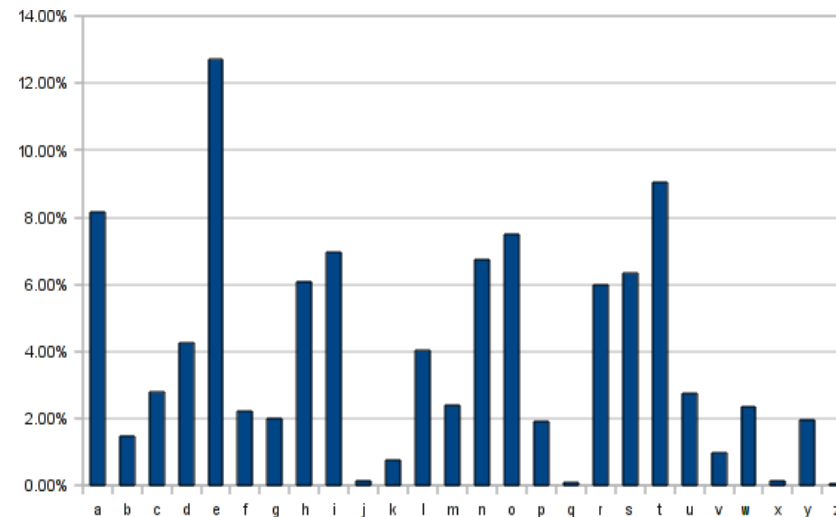
Oh, by the way, bits matter now

- When we move from memory storage to transmitting information through a network, every little bit counts!
- The bit is the most basic unit information in Claude Shannon's *Information Theory*.
- We essentially forget any kind of packet padding that might exist at the packet – level representation of data, and count the transmitted bits from source to sink.



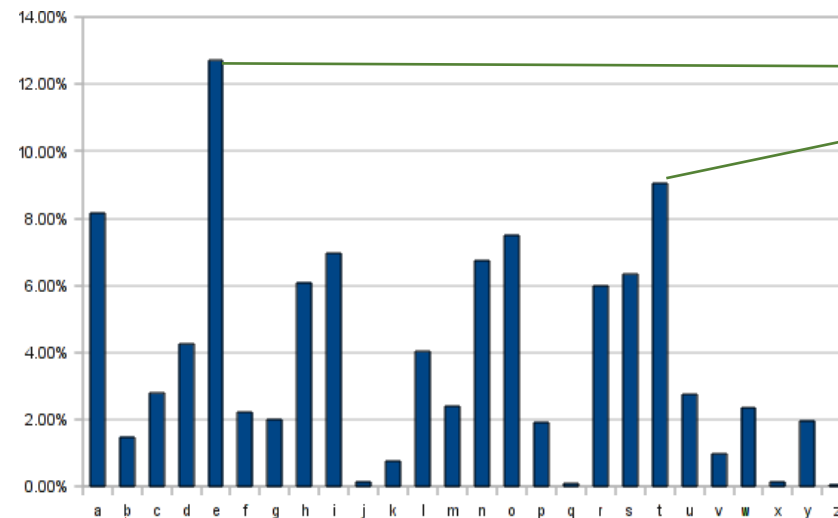
What we want

- For unigrams (characters) with the largest frequencies in this histogram to receive economical storage, since we expect to use them a lot, given their frequency.
- For those that we expect to rarely use, accept a long bit representation.



What we want

- For unigrams (characters) with the largest frequencies in this histogram to receive economical storage, since we expect to use them a lot, given their frequency.
- For those that we expect to rarely use, accept a long bit representation.



Frequent? Make them cheap.

Rare? I don't care!

Careful!

- A crucial detail: **We want our encodings to not be confusable!** For example, if our encoding for e is 01, that of a is 00 and that of h is 0100, **what should the sink interpret if it sees 0100?**



- So any encoding algorithm we come up **with has this constraint as well!**

Careful!

- A crucial detail: **We want our encodings to not be confusable!** For example, if our encoding for e is 01, that of a is 00 and that of h is 0100, **what should the sink interpret if it sees 0100?**



- So any encoding algorithm we come up **with has this constraint as well!**
- We call this property the **prefix property**.

Huffman Coding

- We will describe a **simple greedy algorithm** that is guaranteed to maintain that property.
 - This algorithm is called **Huffman Coding**, pioneered by [David A. Huffman](#)
- The algorithm will **build a binary trie, bottom-up**.

Huffman Coding

- We will describe a **simple greedy algorithm** that is guaranteed to maintain that property.
 - This algorithm is called **Huffman Coding**, pioneered by [David A. Huffman](#)
- The algorithm will **build a binary trie, bottom-up**.
- Assumption: **A loop has been run over the entire text, and a histogram of character frequencies has been constructed**.
 - This can be **parallelized across $k \geq 2$ threads** so that the time is reduced from $|T|$ to $\frac{|T|}{k} + c$, where c is the small cost of adding up and normalizing the individual frequencies of the characters, as they were returned by the threads.

Huffman Coding

- We insert all the `<character, frequency>` pairs in a **Priority Queue**
- Our goal will be, at every point in time, to have **fast access** to **the two nodes with the smallest frequencies**

Building the binary tree....

- Running example string:

“Jason loves chocolate.”

Building the binary tree....

- Running example string:

“Jason loves chocolate.”



Don't forget spaces and punctuation!

Building the binary tree....

- Running example string:

"Jason loves chocolate."

Don't forget spaces and punctuation!

1. Scan the text and **build nodes containing <character, #occurrences> pairs**

SPACE CHARACTER!

J, 1

a, 2

s, 2

o, 4

n, 1

l, 2

v, 1

e, 2

c, 2

h, 1

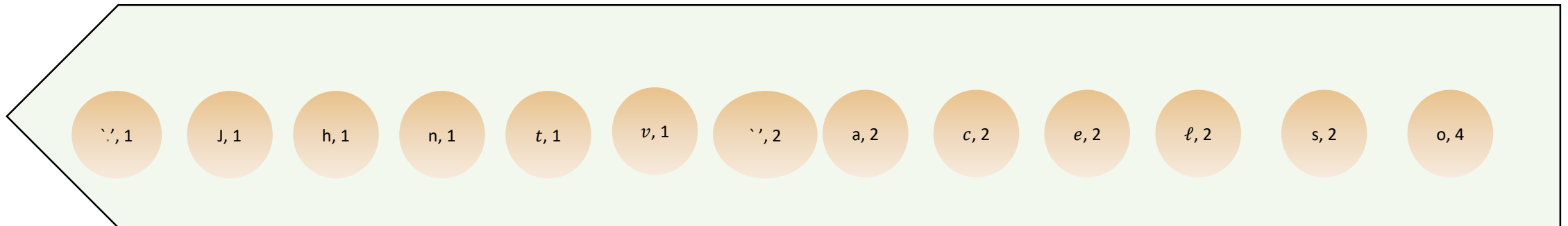
t, 1

`, 2

`, 1

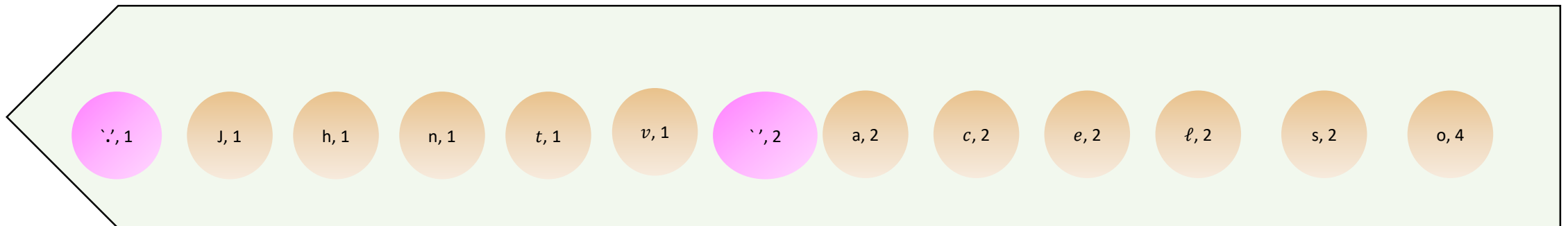
Building the binary tree....

2. Put all of the nodes in a **Priority Queue** sorted on frequencies, in ascending order. **Smaller is higher priority.**
 - Ties will be broken lexicographically.
 - Again, **smaller up front** 😊



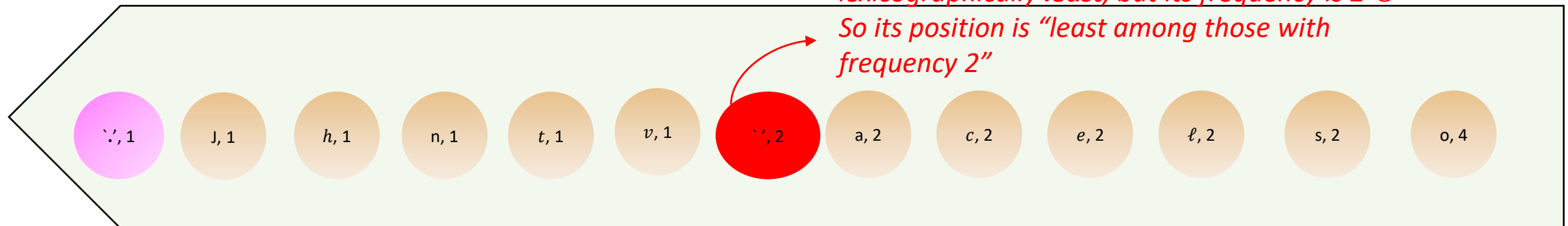
Building the binary tree....

- Put all of the nodes in a **Priority Queue** sorted on frequencies, in ascending order. **Smaller is higher priority.**
 - Ties will be broken lexicographically.
 - Again, **smaller up front** 😊
 - Lexicographically, **periods < characters** and **spaces < periods < characters**
 - Don't worry:** In exams, **you will be given an ASCII table** to deduce lexicographic order on the fly.



Building the binary tree....

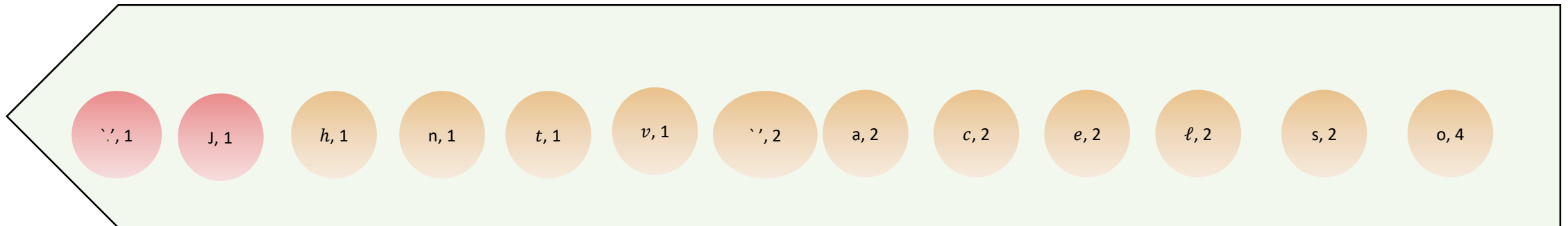
2. Put all of the nodes in a **Priority Queue** sorted on frequencies, in ascending order. **Smaller is higher priority.**
- Ties will be broken lexicographically.
 - Again, **smaller up front** 😊
 - Lexicographically, **periods < characters** and **spaces < periods < characters**
 - **Don't worry:** In exams, **you will be given an ASCII table** to deduce lexicographic order on the fly.



Building the binary tree....

3. As long as the queue has more than 1 elements in it

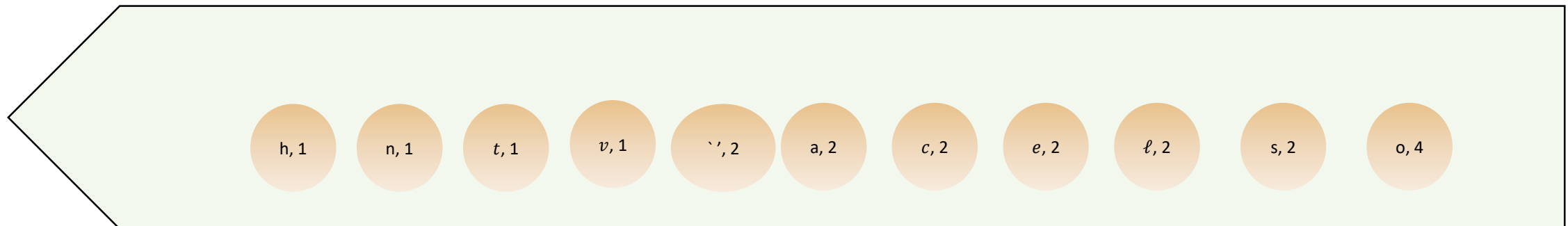
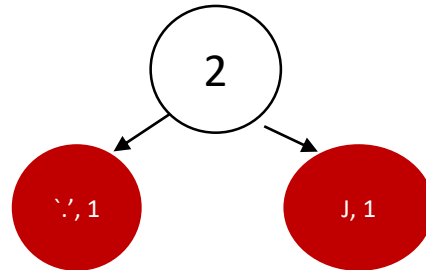
- Apply **getMin()** **twice** to dequeue the 2 **least** elements in the queue.
 - For a **MinHeap**-based PQ, this is **extremely fast**: $2 \cdot \log_2 |A|$ where A is your alphabet. Even for the subset of printable Unicode U with $|U| = 65535$ we have $2 \cdot \log_2 |A| = 32$ compares in the **worst case** for a **queue that will be pinned to cache because of temporal locality**! Recall: **heaps stored as arrays** (consecutive cache storage)



Building the binary tree....

3. As long as the queue has more than 1 elements in it

- Apply **getMin()** **twice** to dequeue the 2 **least** elements in the queue.
 - For a **MinHeap**-based PQ, this is **extremely fast**: $O(2 \cdot \log_2 |A|)$ where A is your alphabet. Even for the subset of printable Unicode U with $|U| = 65535 \Rightarrow 2 \cdot \log_2 |A| = 32$ compares in the **worst case** for a **queue that will be pinned to cache because of temporal locality**! Recall: **heaps stored as arrays** (consecutive cache storage)
- Create a **new node with those two nodes as children** and a frequency value equal to the **sum of the children's frequency values**.

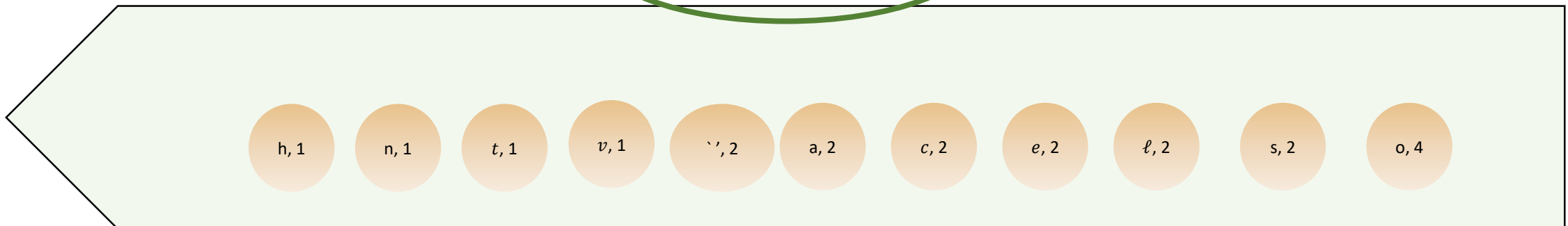
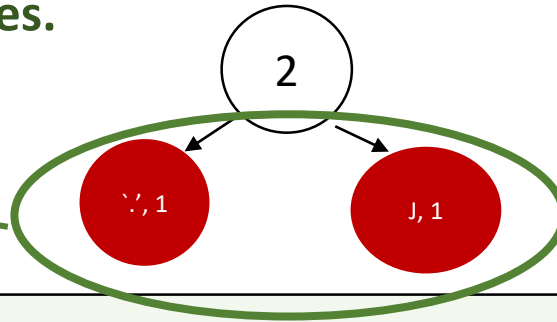


Building the binary tree....

3. As long as the queue has more than 1 elements in it

- Apply **getMin()** **twice** to dequeue the 2 **least** elements in the queue.
 - For a **MinHeap**-based PQ, this is **extremely fast**: $O(2 \cdot \log_2 |A|)$ where A is your alphabet. Even for the subset of printable Unicode U with $|U| = 65535 \Rightarrow 2 \cdot \log_2 |A| = 32$ compares in the **worst case** for a **queue that will be pinned to cache because of temporal locality**! Recall: **heaps stored as arrays** (consecutive cache storage)
- Create a **new node** with those two nodes as children and a frequency value equal to the **sum of the children's frequency values**.

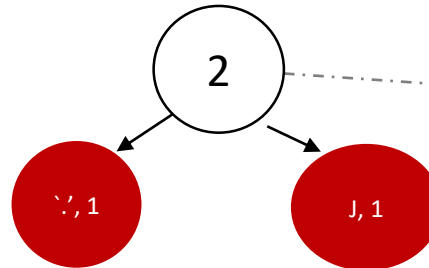
We follow the convention that **the first child pulled from the queue is on the left**, and **the second on the right**



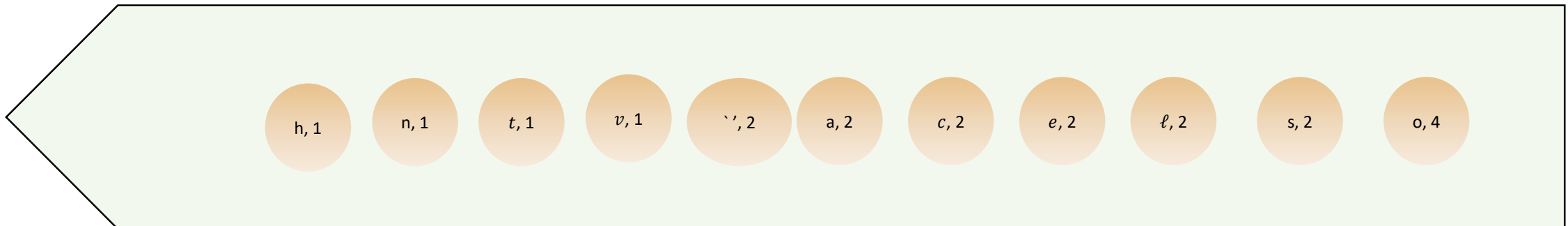
Building the binary tree....

3. As long as the queue has more than 1 elements in it

- Apply **getMin()** **twice** to dequeue the 2 **least** elements in the queue.
 - For a **MinHeap**-based PQ, this is **extremely fast**: $O(2 \cdot \log_2 |A|)$ where A is your alphabet. Even for the subset of printable Unicode U with $|U| = 65535 \Rightarrow 2 \cdot \log_2 |A| = 32$ compares in the **worst case** for a **queue that will be pinned to cache because of temporal locality**! Recall: **heaps stored as arrays** (consecutive cache storage)
- Create a **new node with those two nodes as children** and a frequency value equal to the **sum of the children's frequency values**.



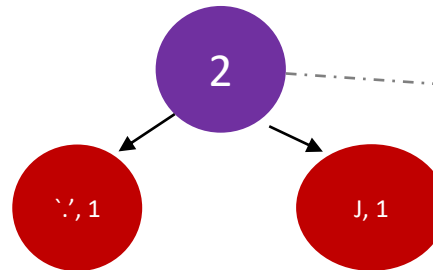
Important implementation detail; These nodes are assumed to hold the NULL character (first in ASCII table!)



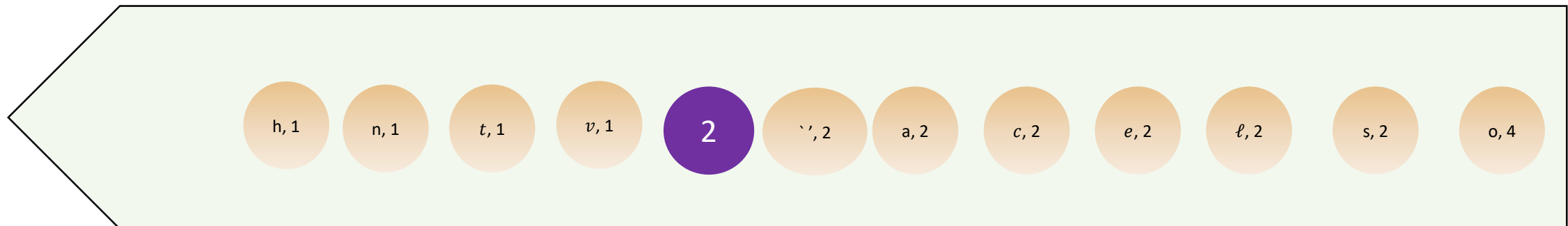
Building the binary tree....

3. As long as the queue has more than 1 elements in it

- Apply **getMin()** **twice** to dequeue the 2 **least** elements in the queue.
 - For a **MinHeap**-based PQ, this is **extremely fast**: $O(2 \cdot \log_2 |A|)$ where A is your alphabet. Even for the subset of printable Unicode U with $|U| = 65535 \Rightarrow 2 \cdot \log_2 |A| = 32$ compares in the **worst case** for a **queue that will be pinned to cache because of temporal locality**! Recall: **heaps stored as arrays** (consecutive cache storage)
- Create a **new node with those two nodes as children** and a frequency value equal to the **sum of the children's frequency values**.
- **Insert that node into the trie!**



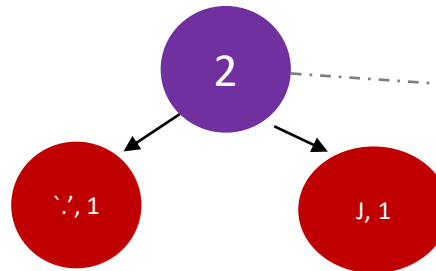
Important implementation detail; These nodes are assumed to hold the NULL character (first in ASCII table!)



Building the binary tree....

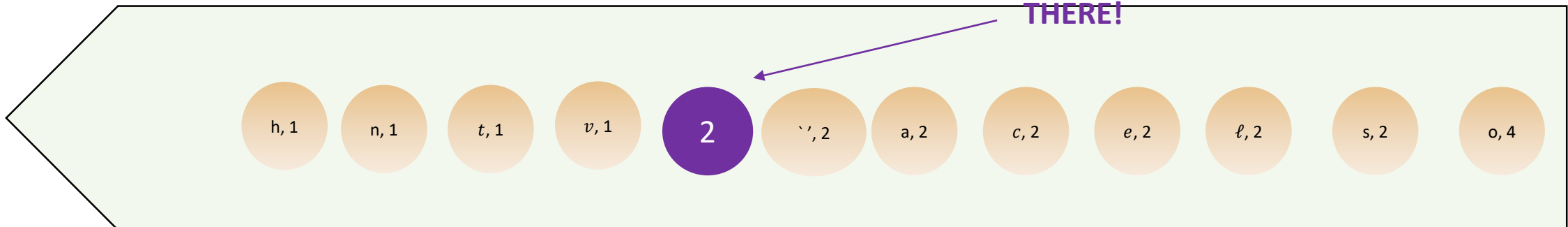
3. As long as the queue has more than 1 elements in it

- Apply **getMin()** **twice** to dequeue the 2 **least** elements in the queue.
 - For a **MinHeap**-based PQ, this is **extremely fast**: $O(2 \cdot \log_2 |A|)$ where A is your alphabet. Even for the subset of printable Unicode U with $|U| = 65535 \Rightarrow 2 \cdot \log_2 |A| = 32$ compares in the **worst case** for a **queue that will be pinned to cache because of temporal locality**! Recall: **heaps stored as arrays** (consecutive cache storage)
- Create a **new node with those two nodes as children** and a frequency value equal to the **sum of the children's frequency values**.
- **Insert that node into the trie!**

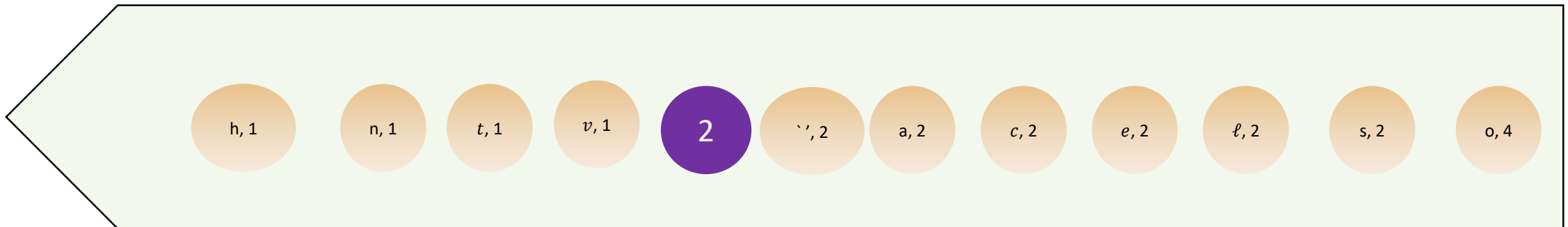
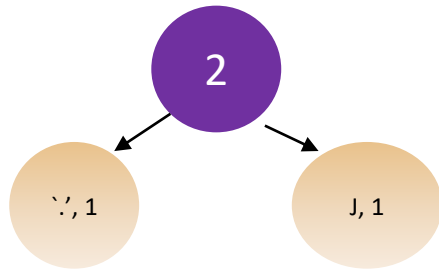


Important implementation detail; These nodes are assumed to hold the NULL character (first in ASCII table!)

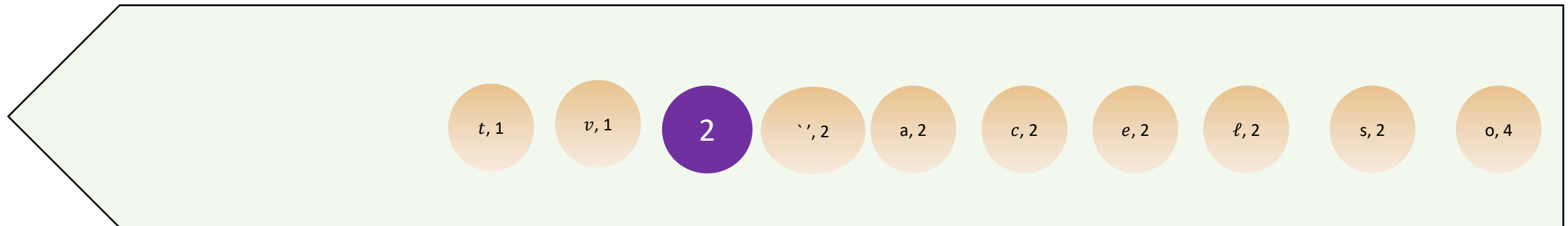
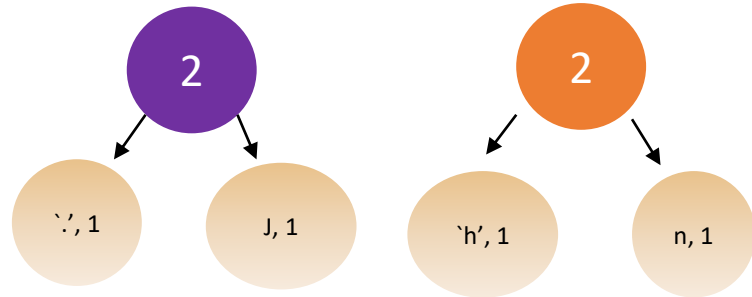
And that is why the node fits **EXACTLY THERE!**



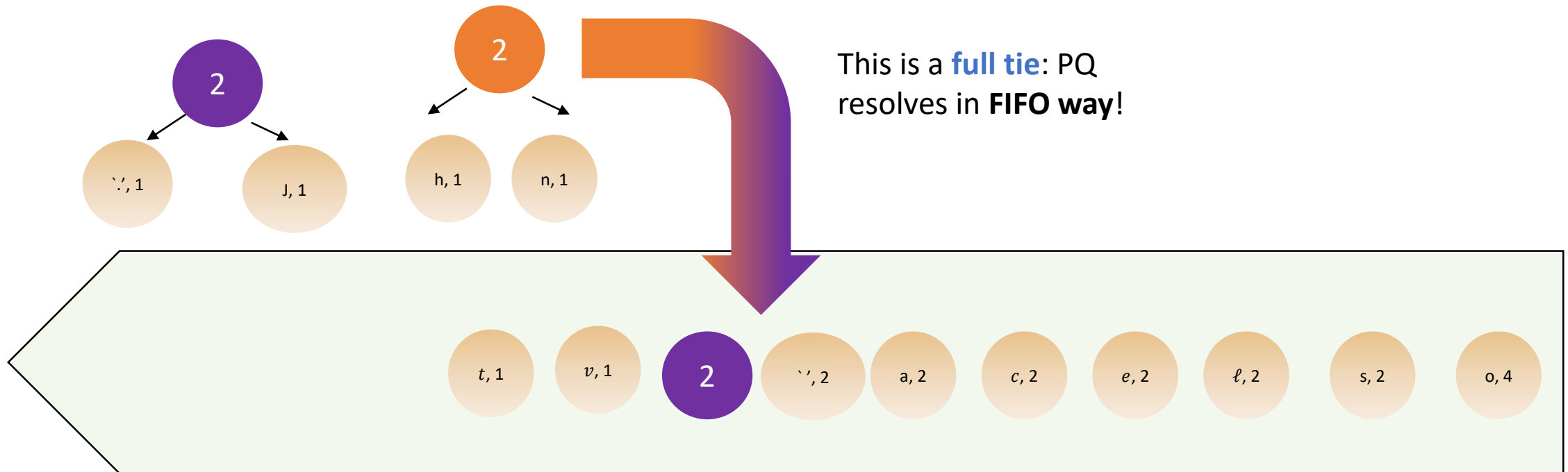
Building the binary tree....



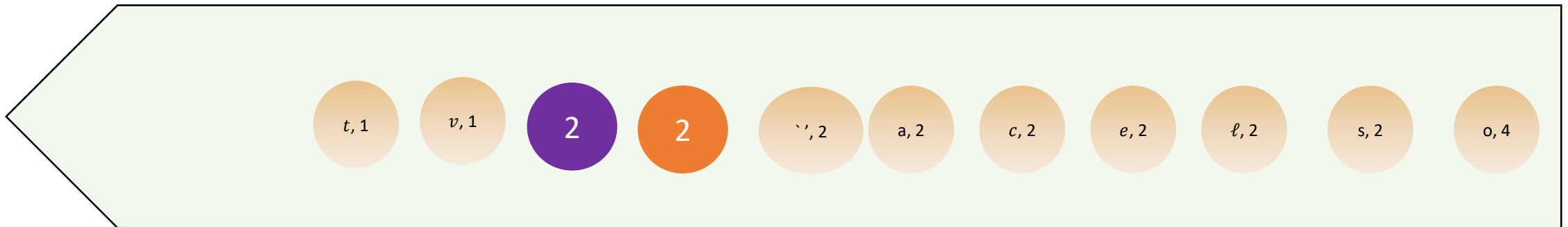
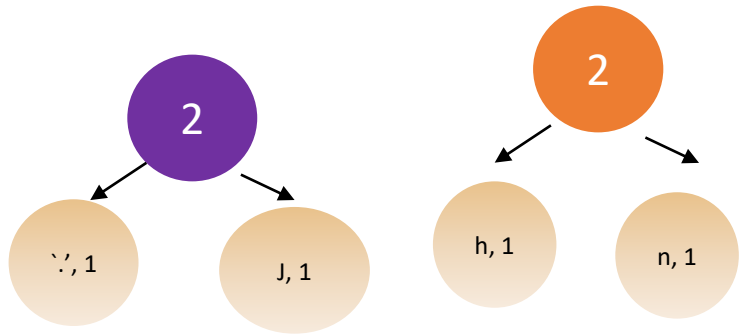
Building the binary tree....



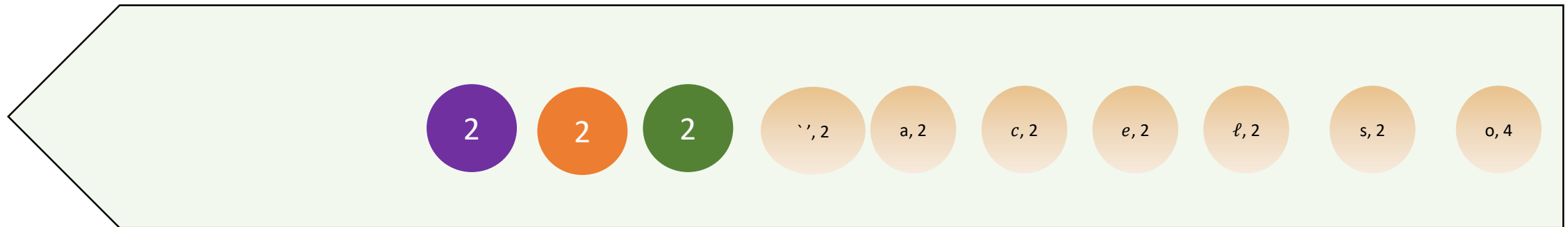
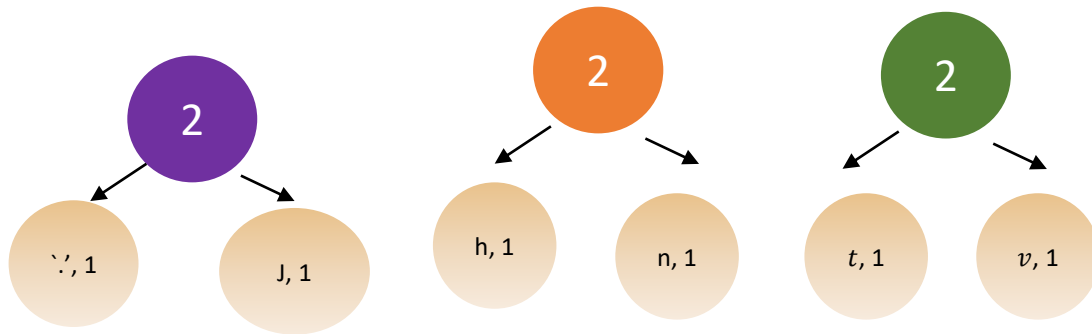
Building the binary tree....



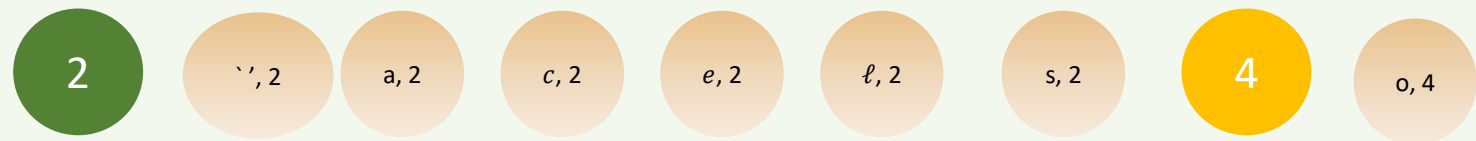
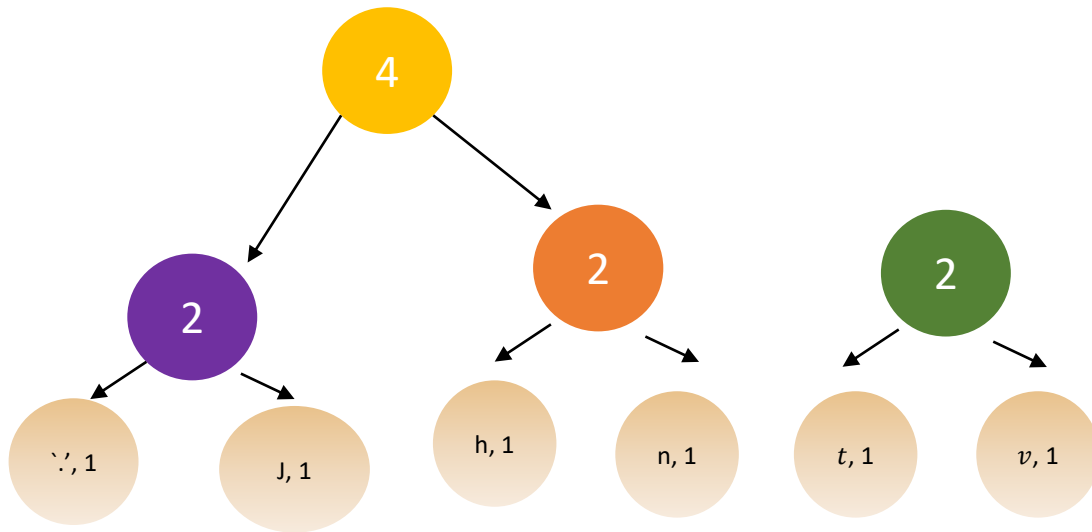
Building the binary tree....



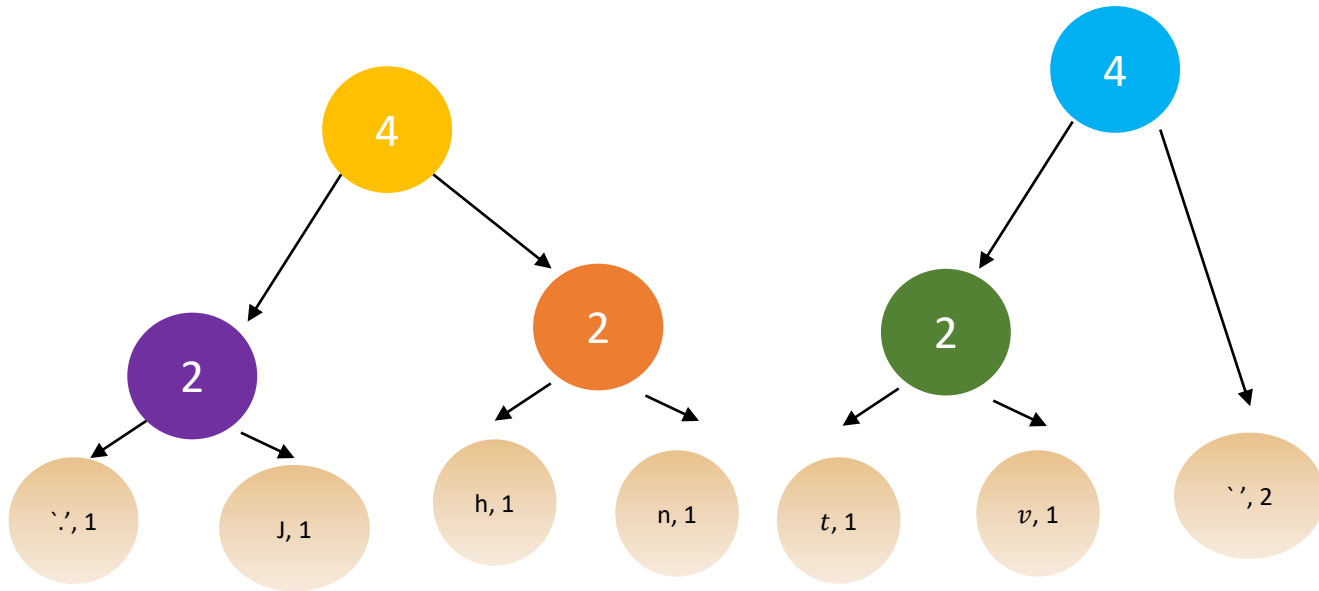
Building the binary tree....



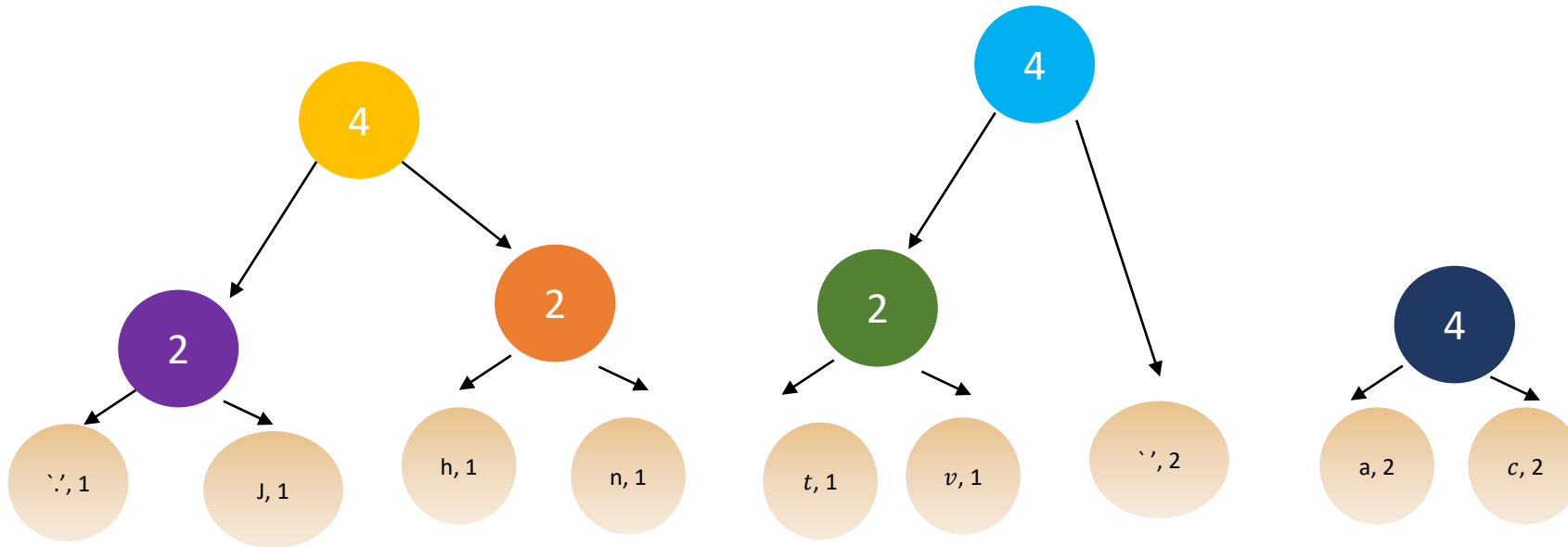
Building the binary tree....



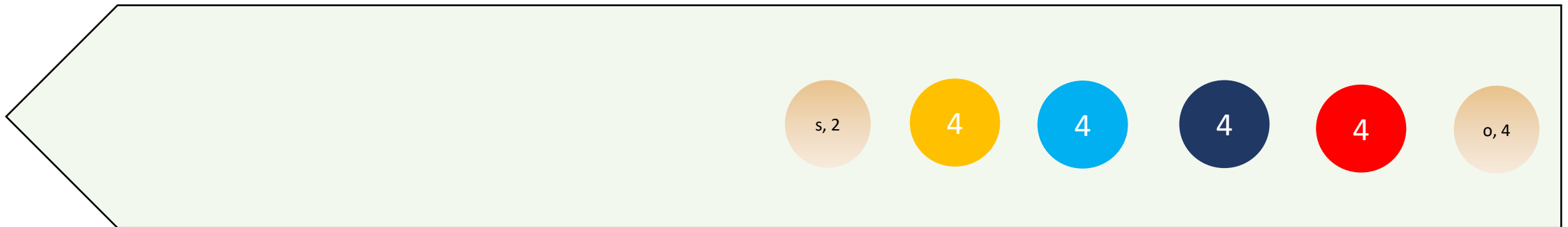
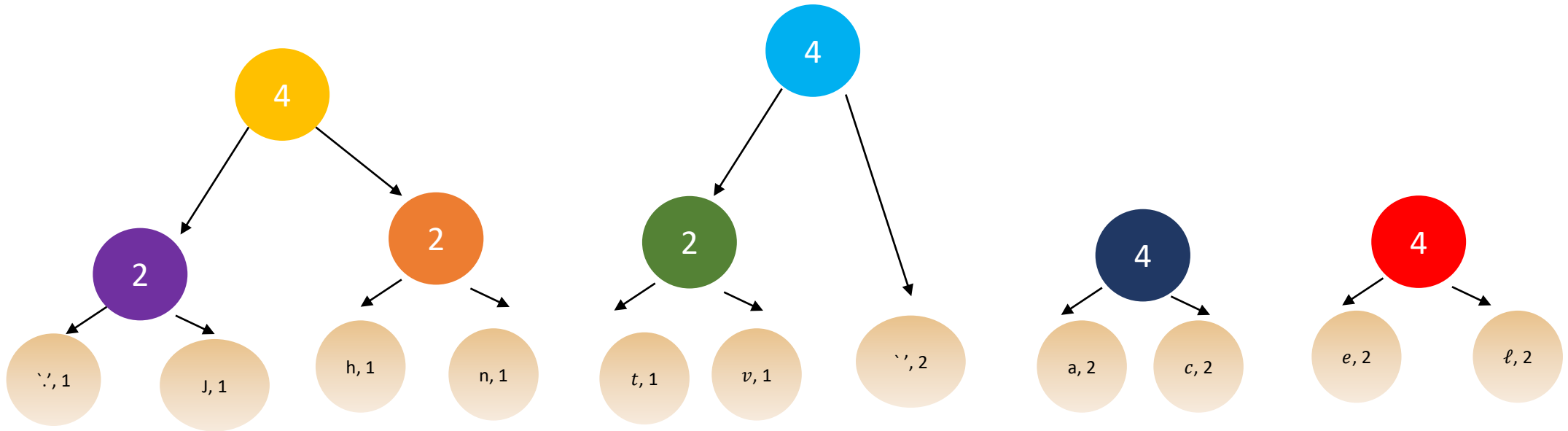
Building the binary tree....



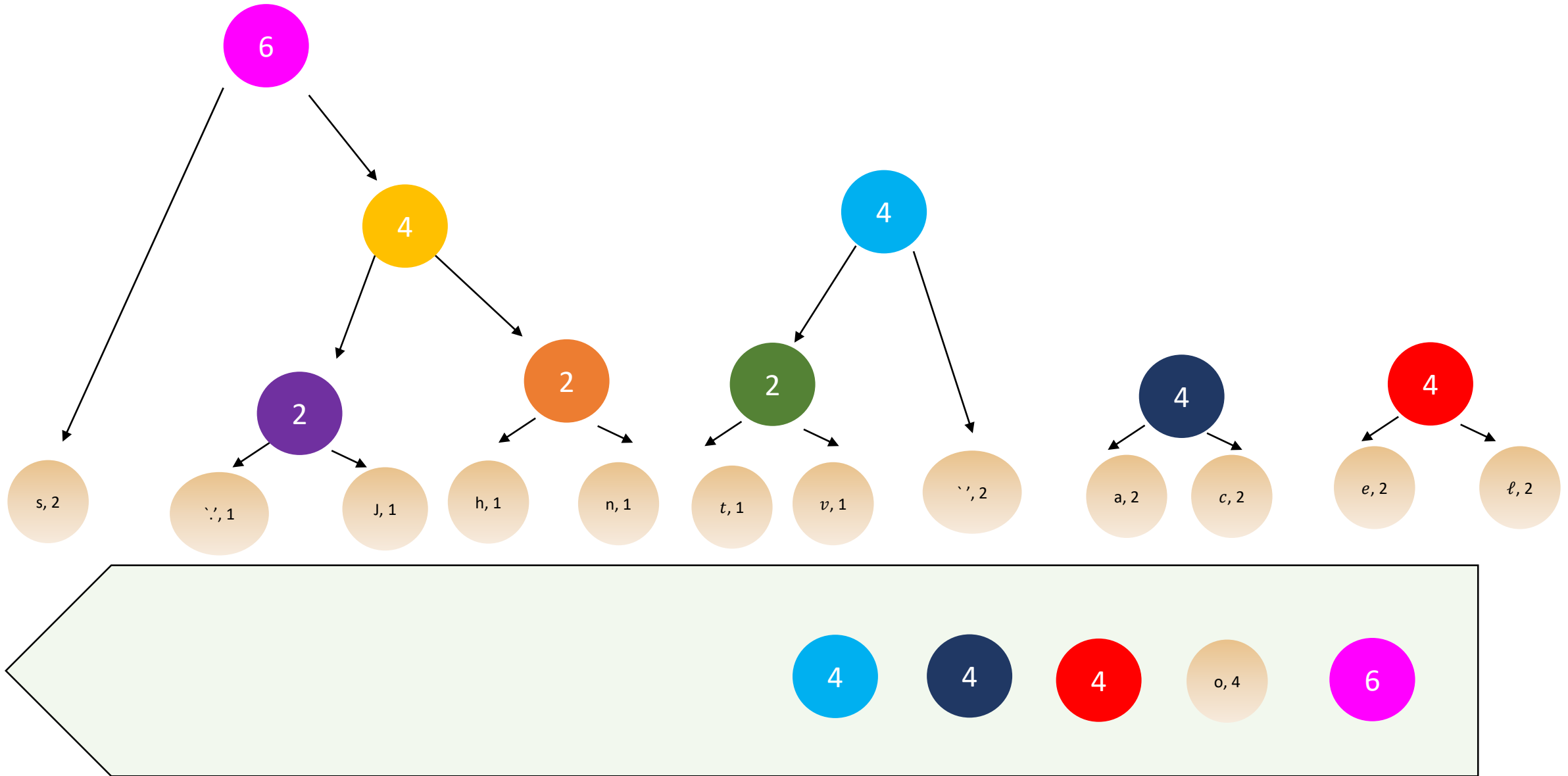
Building the binary tree....



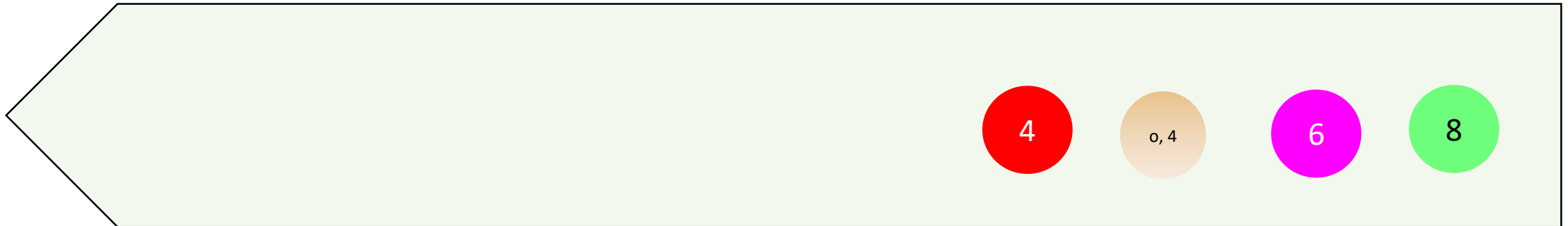
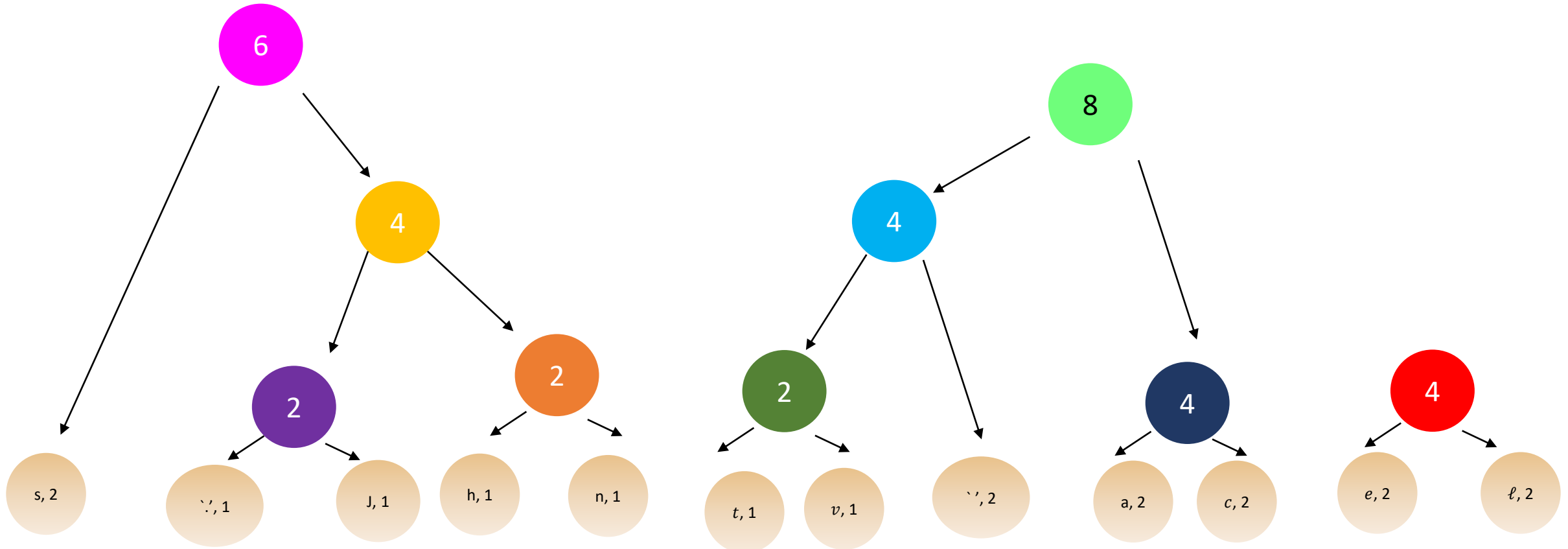
Building the binary tree....



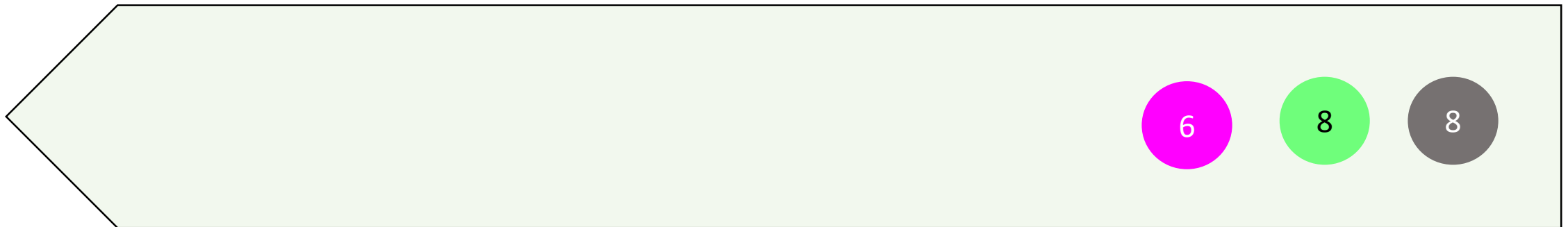
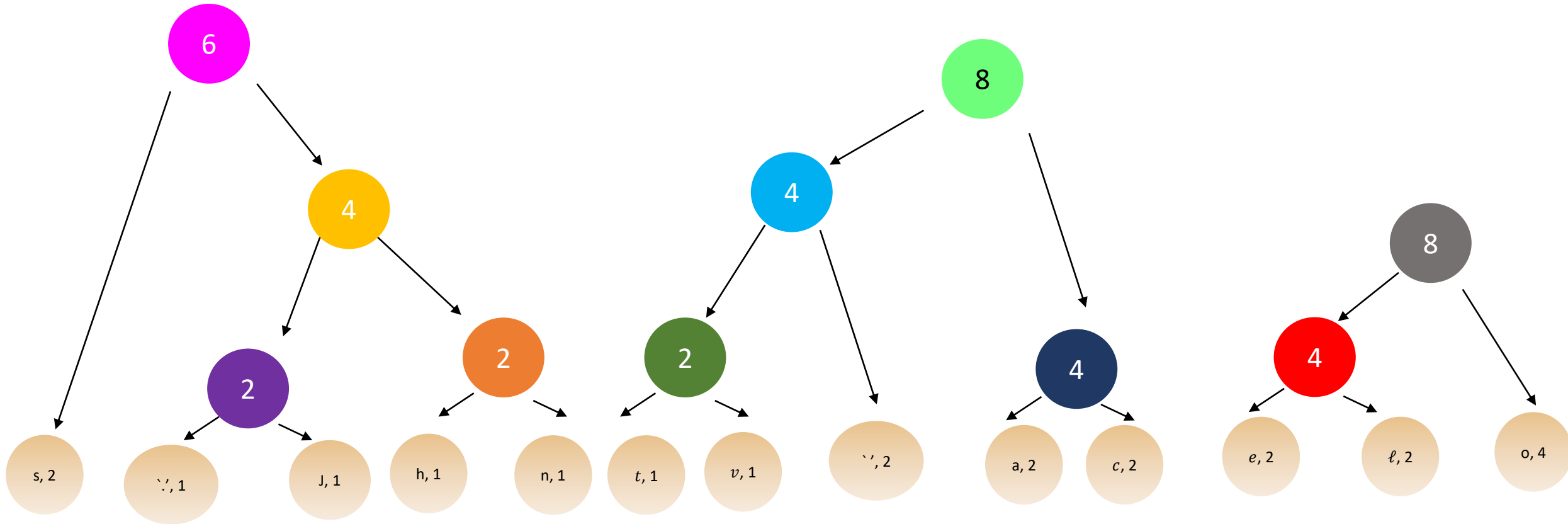
Building the binary tree....



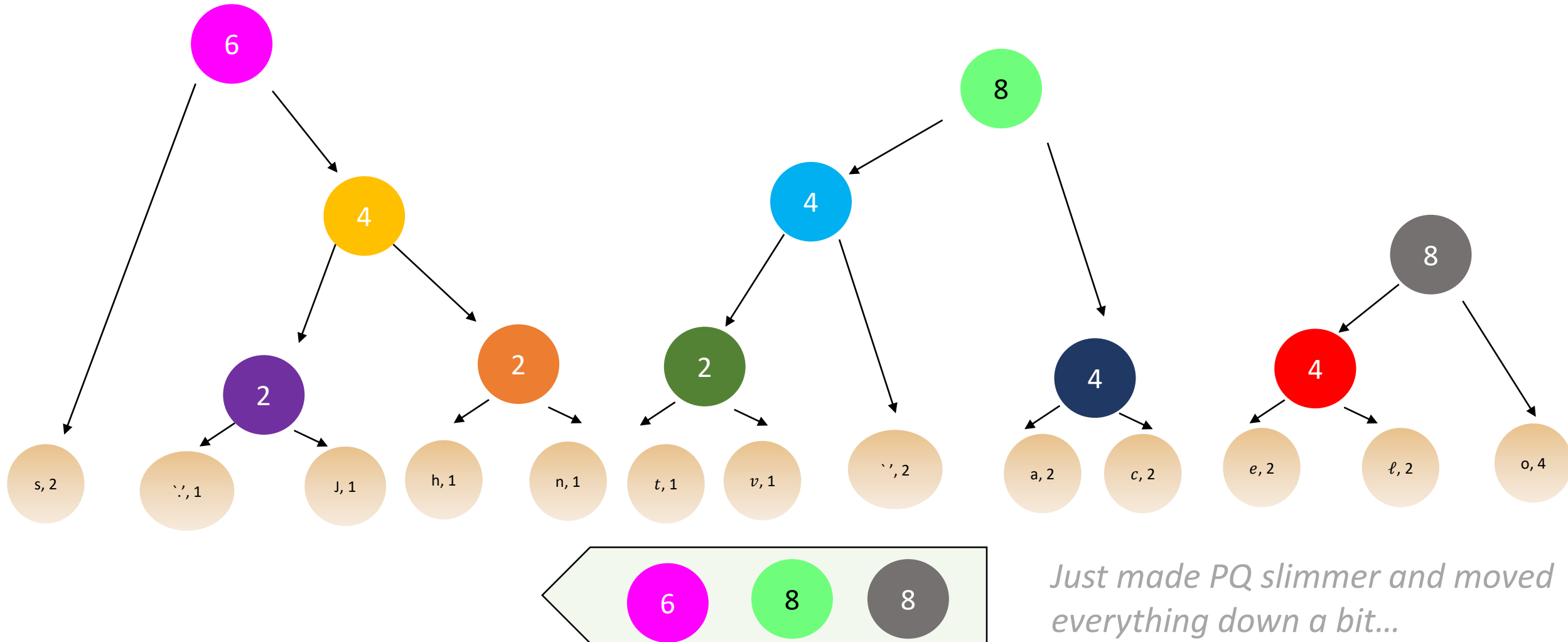
Building the binary tree....



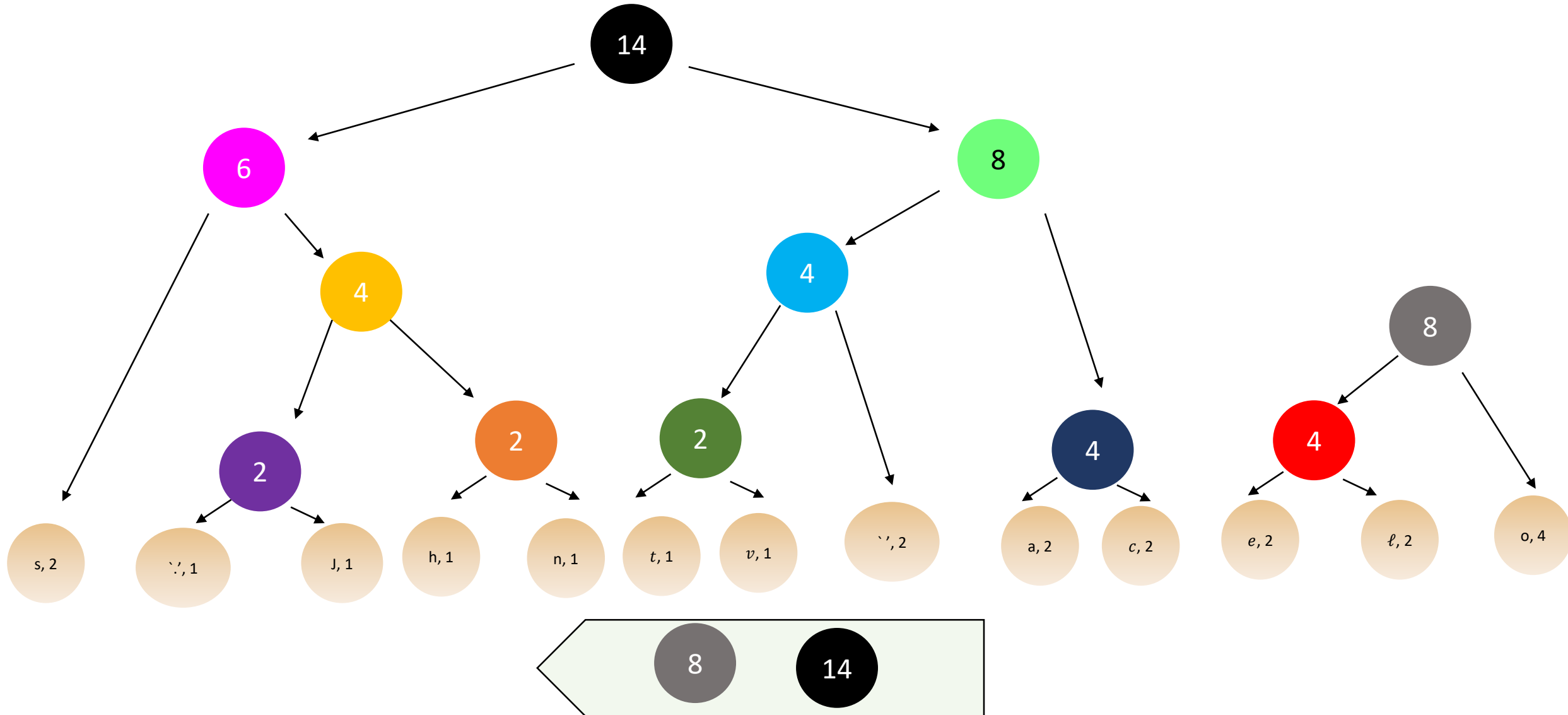
Building the binary tree....



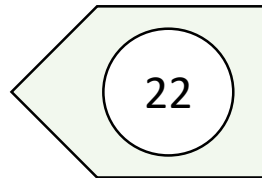
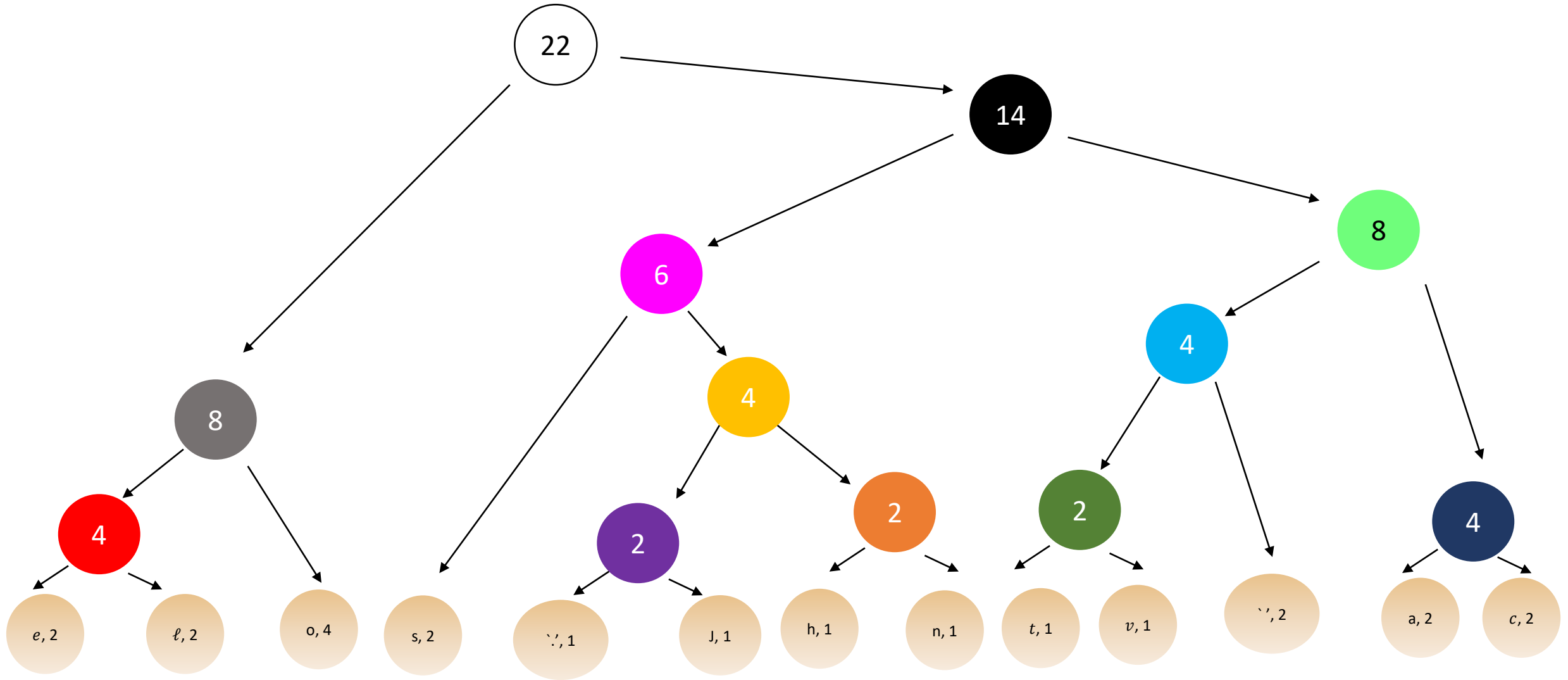
Building the binary tree....



Building the binary tree....



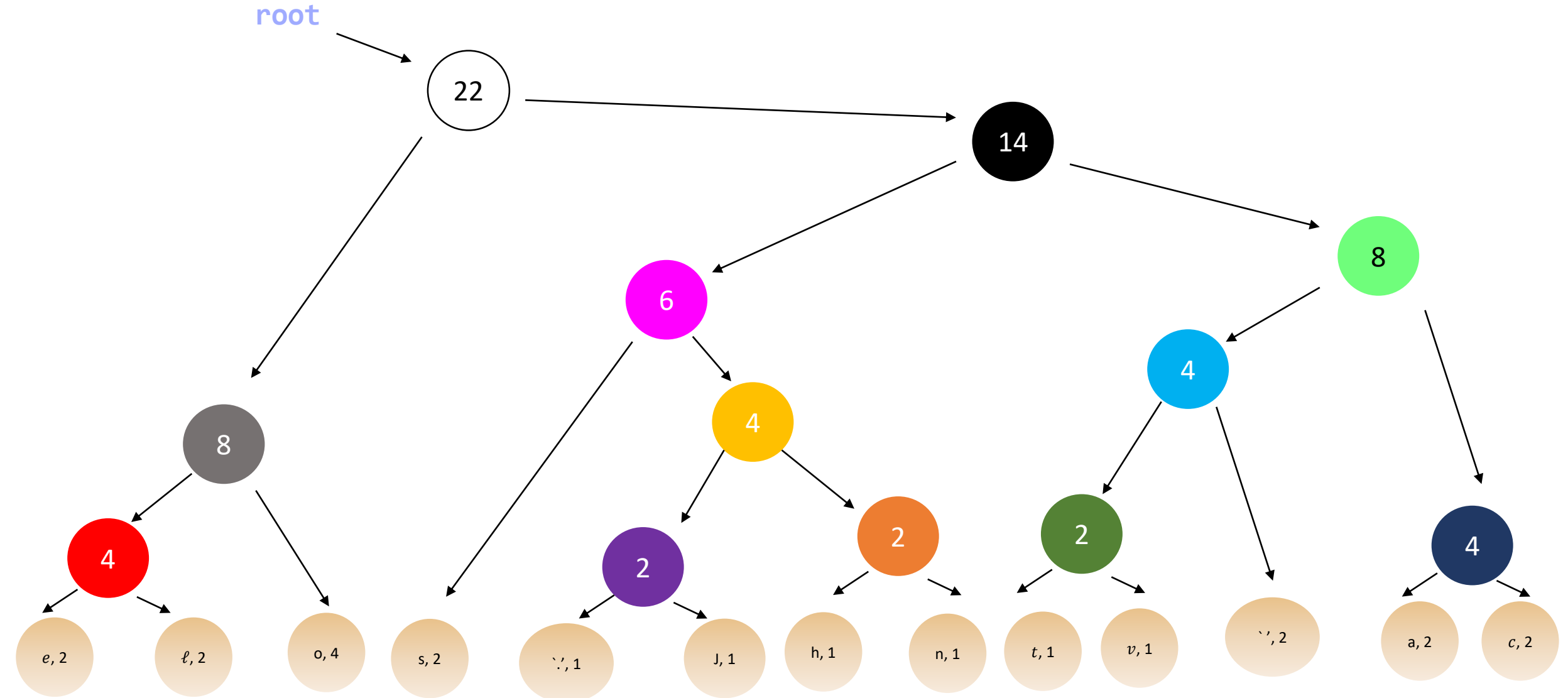
Building the binary tree....



Queue now has single element; loop ends!

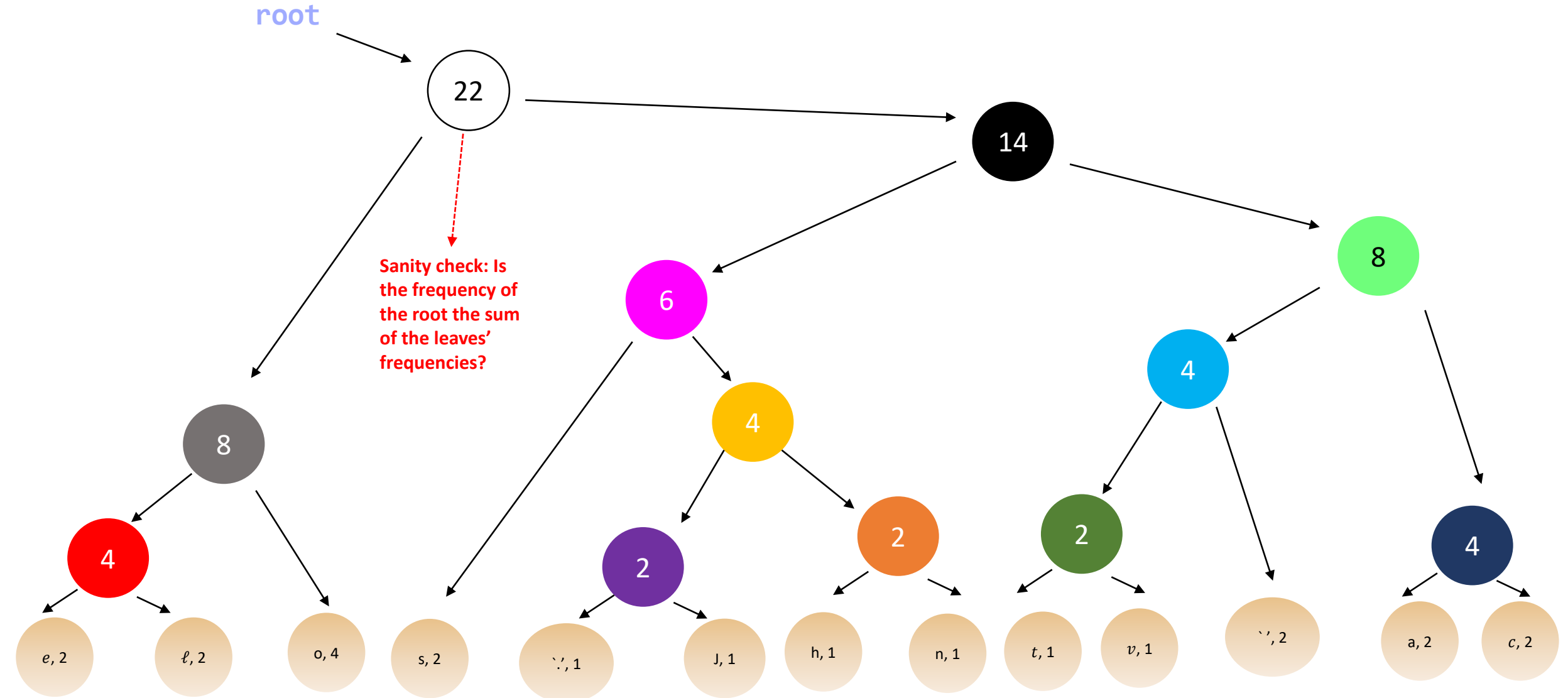
Building the binary tree....

4. Last node in PQ becomes **root of tree**

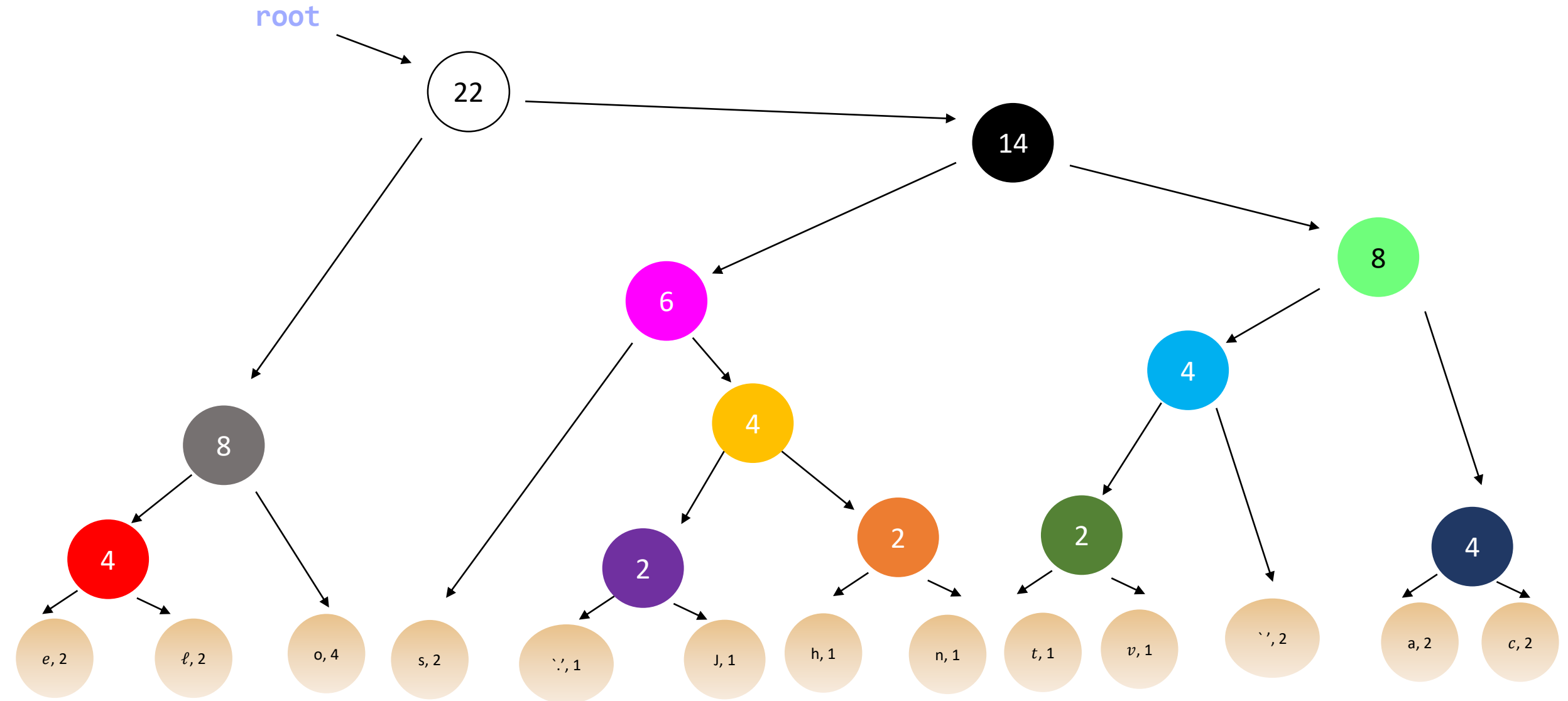


Building the binary tree....

4. Last node in PQ becomes **root of tree**

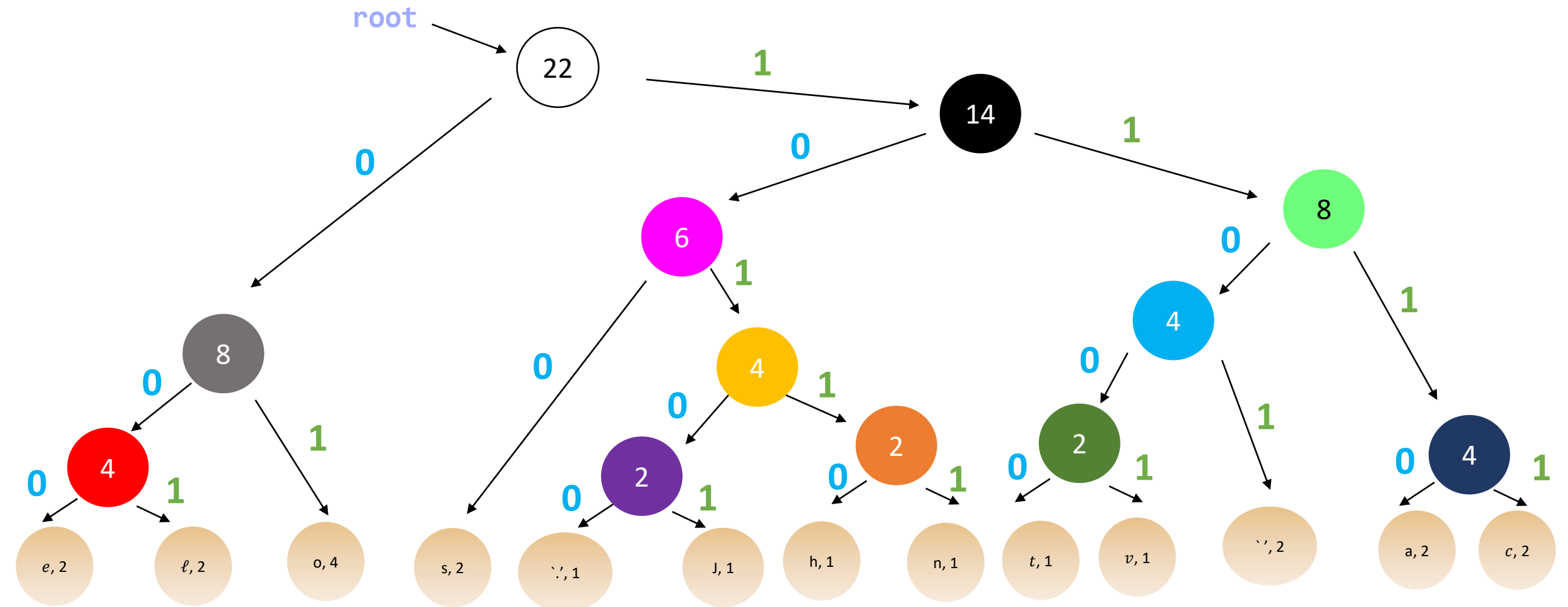


AND NOW FOR THE KILLER....



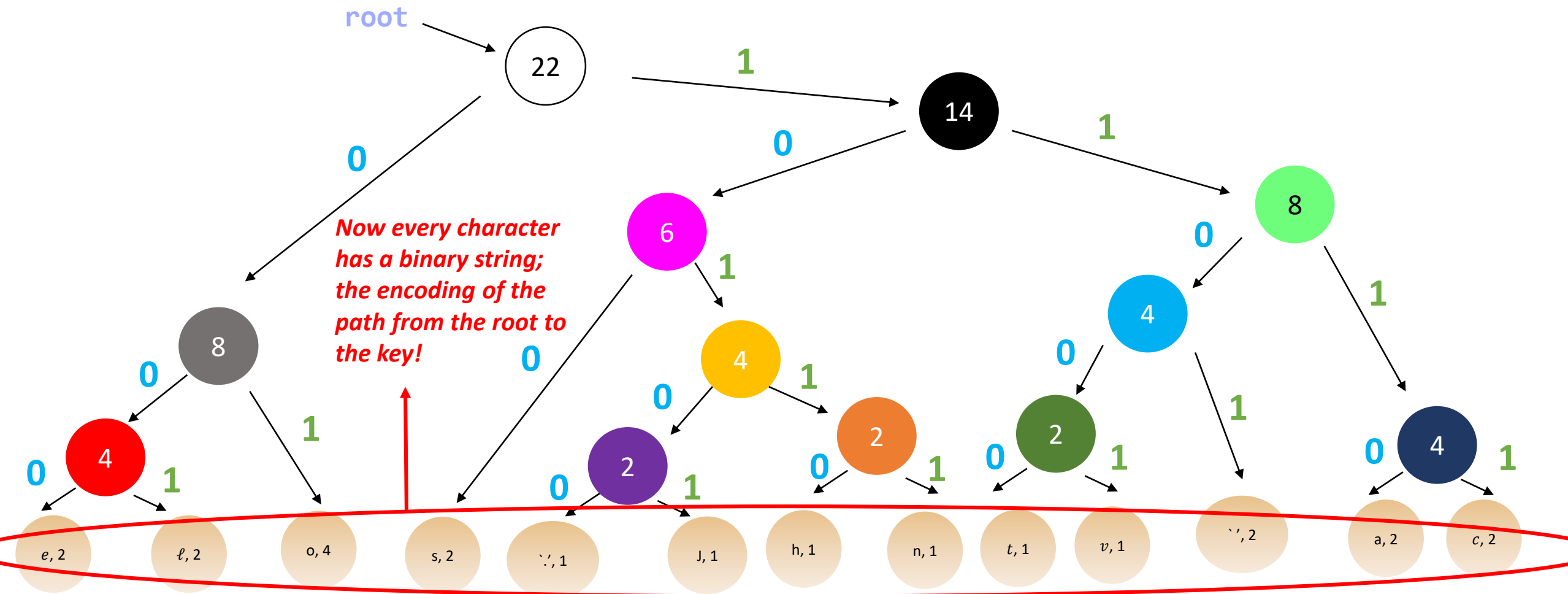
Building the **binary trie**!

5. Associate **left links with 0s** and **right links with 1s**, turning the binary tree into a **binary trie**!



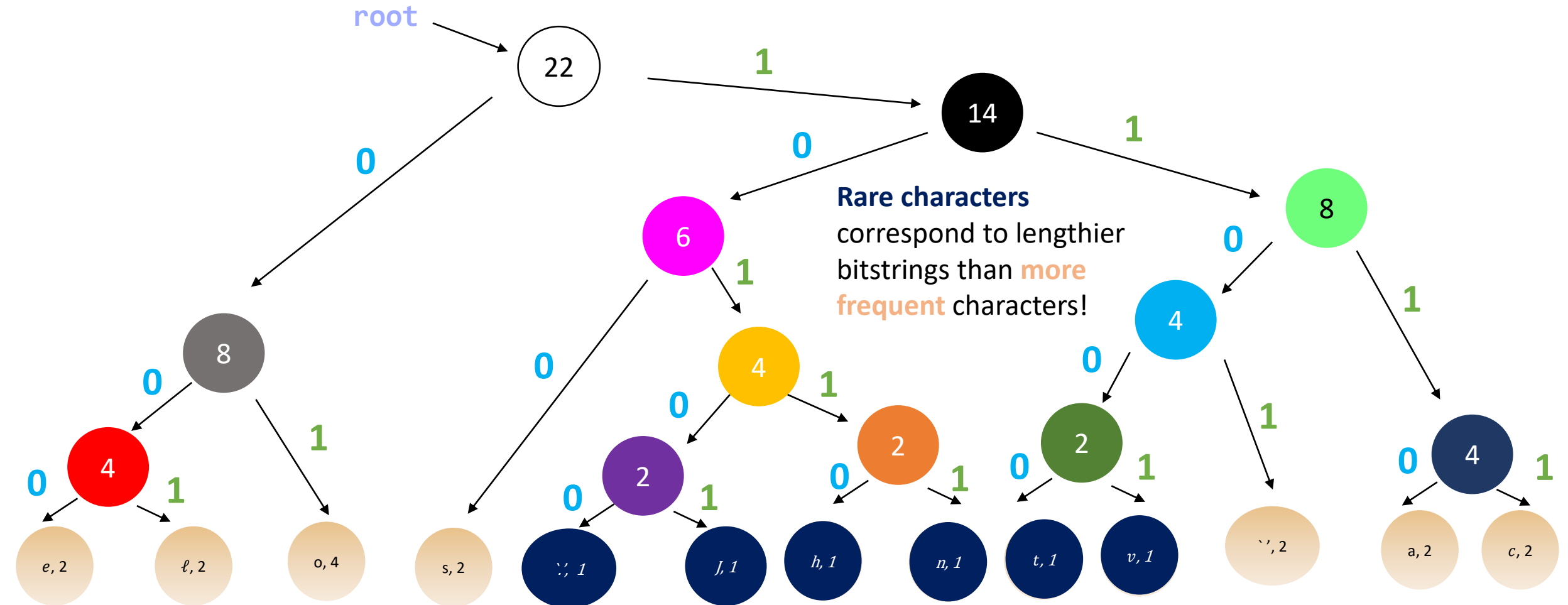
Building the **binary trie**!

5. Associate **left links with 0s** and **right links with 1s**, turning the binary tree into a **binary trie**!



Building the **binary trie**!

5. Associate **left links with 0s** and **right links with 1s**, turning the binary tree into a **binary trie**!



Final structure: lookup table

6. Perform a preorder traversal of the trie to build a **two-way lookup table** that **associates characters with Huffman binary encodings** and vice versa
 - Likeliest inner implementation of this lookup table: **two hash tables**.
 - In exams, **you won't have to do this**, since you can immediately “see” the encoding without the use of a lookup table.

Final structure: lookup table

Character	Binary encoding
e	000
<i>ℓ</i>	001
o	01
s	100
.	10100
J	10101
h	10110
n	10111
t	11000
v	11001
SPACE	1101
a	1110
c	1111

Encoding our message in two ways

Character	Huffman Binary encoding	ASCII (7-bit) encoding
e	000	1100101
ℓ	001	1101100
o	01	1101111
s	100	1110011
.	10100	101110
J	10101	1001010
h	10110	1101000
n	10111	1101110
t	11000	1110100
v	11001	1110110
SPACE	1101	0010000
a	1110	1100001
c	1111	1100101

Encoding our message in two ways

Character	Huffman Binary encoding	ASCII (7-bit) encoding
e	000	1100101
ℓ	001	1101100
o	01	1101111
s	100	1110011
.	10100	101110
J	10101	1001010
h	10110	1101000
n	10111	1101110
t	11000	1110100
v	11001	1110110
SPACE	1101	0010000
a	1110	1100001
c	1111	1100101

- Note that the most expensive Huffman encodings in this table **are still under 7 bits!**

Encoding our message in two ways

Character	Huffman Binary encoding	ASCII (7-bit) encoding
e	000	1100101
l	001	1101100
o	01	1101111
s	100	1110011
.	10100	101110
J	10101	1001010
h	10110	1101000
n	10111	1101110
t	11000	1110100
v	11001	1110110
SPACE	1101	0010000
a	1110	1100001
c	1111	1100101

- Let's compare how the message

“Jason loves chocolate”

decodes through both ASCII and Huffman:

Encoding our message in two ways

Character	Huffman Binary encoding	ASCII (7-bit) encoding
e	000	1100101
l	001	1101100
o	01	1101111
s	100	1110011
.	10100	101110
J	10101	1001010
h	10110	1101000
n	10111	1101110
t	11000	1110100
v	11001	1110110
SPACE	1101	0010000
a	1110	1100001
c	1111	1100101

- Let's compare how the message

"Jason loves chocolate"

decodes through both ASCII and Huffman:

- ASCII:** 1001010110000111100111101111
110111000100001101100110111111011011001011
11001100100001100011 1101000 1101111 1100011
1101111 11011001101111 1100001 1110100
1100101, 147 bits total.

Encoding our message in two ways

Character	Huffman Binary encoding	ASCII (7-bit) encoding
e	000	1100101
l	001	1101100
o	01	1101111
s	100	1110011
.	10100	101110
J	10101	1001010
h	10110	1101000
n	10111	1101110
t	11000	1110100
v	11001	1110110
SPACE	1101	0010000
a	1110	1100001
c	1111	1100101

- Let's compare how the message

"Jason loves chocolate"

decodes through both ASCII and Huffman:

- ASCII:** 1001010110000111100111101111
1101110001000011011001101111110110110010111
1001100100001100011 1101000 1101111 1100011
1101111 11011001101111 1100001 1110100 1100101,
147 bits total.
- Huffman:**
10101111010001101111101001011100100010011011
1111011001111101001111011000000, 75 bits total

Encoding our message in two ways

Character	Huffman Binary encoding	ASCII (7-bit) encoding
e	000	1100101
l	001	1101100
o	01	1101111
s	100	1110011
.	10100	101110
J	10101	1001010
h	10110	1101000
n	10111	1101110
t	11000	1110100
v	11001	1110110
SPACE	1101	0010000
a	1110	1100001
c	1111	1100101

- Let's compare how the message

"Jason loves chocolate"

decodes through both ASCII and Huffman:

- ASCII:** 1001010110000111100111101111
1101110001000011011001101111111011011001011
11001100100001100011 1101000 1101111 1100011
1101111 11011001101111 1100001 1110100
1100101, 147 bits total.



- Huffman:**
10101111010001101111110100101110010001001101
11111011001111101001111011000000, 75 bits total
≈ 50% of ASCII !

Encoding our message in two ways

Character	Huffman Binary encoding	ASCII (7-bit) encoding
e	000	1100101
l	001	1101100
o	01	1101111
s	100	1110011
.	10100	101110
J	10101	1001010
h	10110	1101000
n	10111	1101110
t	11000	1110100
v	11001	1110110
SPACE	1101	0010000
a	1110	1100001
c	1111	1100101

- Can I encode other strings using this table (remember, we created it using the string “Jason loves chocolate”)?

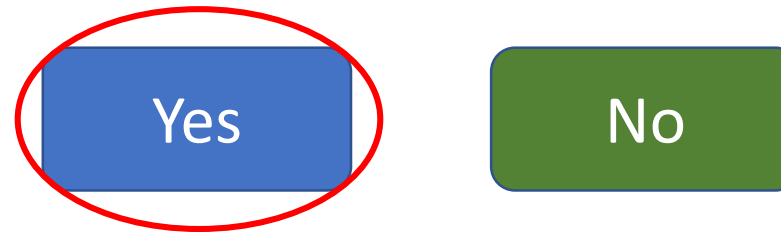
Yes

No

Encoding our message in two ways

Character	Huffman Binary encoding	ASCII (7-bit) encoding
e	000	1100101
ℓ	001	1101100
o	01	1101111
s	100	1110011
.	10100	101110
J	10101	1001010
h	10110	1101000
n	10111	1101110
t	11000	1110100
v	11001	1110110
SPACE	1101	0010000
a	1110	1100001
c	1111	1100101

- Can I encode other strings using this table (remember, we created it using the string “Jason loves chocolate”)?



- As long as the table covers the subset of the alphabet they need!
- As an example, let's encode the string “Jonas loves nachos” in both ASCII and Huffman, using this table 😊

Encoding our message in two ways

Don't worry, Jason has made sure that the alphabet is the same! :)

"Jonas loves nachos"

Character	Huffman Binary encoding	ASCII (7-bit) encoding
e	000	1100101
l	001	1101100
o	01	1101111
s	100	1110011
.	10100	101110
J	10101	1001010
h	10110	1101000
n	10111	1101110
t	11000	1110100
v	11001	1110110
SPACE	1101	0010000
a	1110	1100001
c	1111	1100101

- ASCII: 100101011011111101110110000
11110011001000011011001101111111
01101100101111001100100001101110
11000011100011110100011011111110
011 = 126 bits

Encoding our message in two ways

Alphabet is the same!
 $\Sigma = \{a, c, e, h, j, l, n, o, s, v\}$

“Jonas loves nachos”

Character	Huffman Binary encoding	ASCII (7-bit) encoding
e	000	1100101
l	001	1101100
o	01	1101111
s	100	1110011
.	10100	1011110
J	10101	1001010
h	10110	1101000
n	10111	1101110
t	11000	1110100
v	11001	1110110
SPACE	1101	0010000
a	1110	1100001
c	1111	1100101

- **ASCII:** 10010101101111110111011000011110011001000011011001101111110110110010111100110010000110111011000011100011110100011011111110011 = 126 bits
- **Huffman:** 101011110100011011111010010111001000100110110111111011111011001100 = 66 bits $\approx 52.39\%$ of ASCII

Encoding our message in two ways

Alphabet is the same!
 $\Sigma = \{a, c, e, h, j, l, n, o, s, v\}$

“Jonas loves nachos”

Character	Huffman Binary encoding	ASCII (7-bit) encoding
e	000	1100101
l	001	1101100
o	01	1101111
s	100	1110011
.	10100	1011110
J	10101	1001010
h	10110	1101000
n	10111	1101110
t	11000	1110100
v	11001	1110110
SPACE	1101	0010000
a	1110	1100001
c	1111	1100101

- ASCII: 10010101101111110111011000011110011001000011011001101111110110110010111100110010000110111011000011100011110100011011111110011 = 126 bits

- Huffman: 101011110100011011111010010111001000100110110111111011111011001100 = 66 bits $\approx 52.39\%$ of ASCII
 - Slightly worse than before!
 - Why?



Your turn!

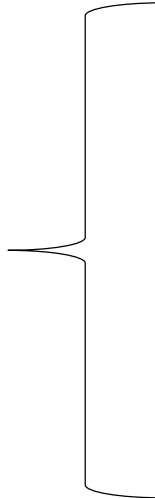
Jason: Give students an ASCII table on a different projector....

- Build the Huffman trie for the string

“Data structures”

Solution

“Data structures”

- 
- SPC: 1
 - 'a': 2
 - 'c': 1
 - 'D': 1
 - 'e': 1
 - 'r': 2
 - 's': 2
 - 't': 3
 - 'u': 2

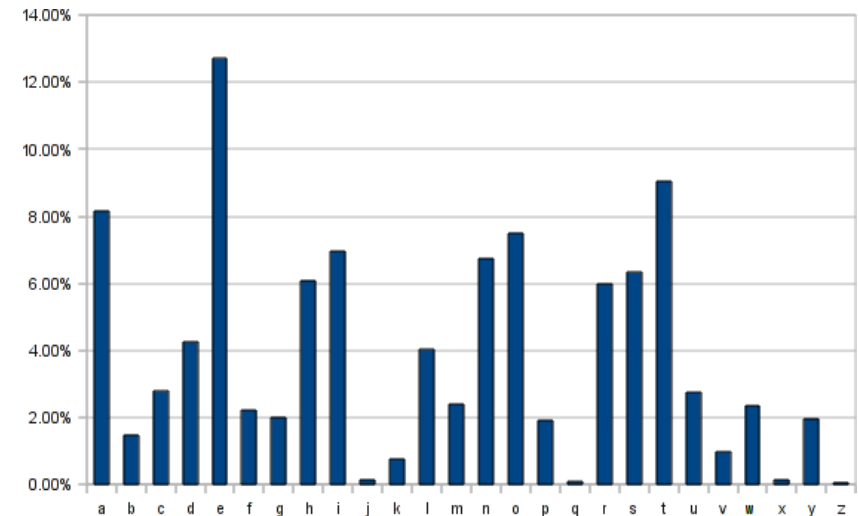
Solution

“Data structures”

- SPC: 1
- ‘a’: 2
- ‘c’: 1
- ‘d’: 1
- ‘e’: 1
- ‘r’: 2
- ‘s’: 2
- ‘t’: 3
- ‘u’: 2

In this text, there are more ‘t’s and ‘u’s than ‘e’s. This is unrealistic in the full spectrum of English.

As our text grows, your histogram will look more and more like the one we’ve seen:



Solution

"Data structures"

- SPC: 1
- 'a': 2
- 'c': 1
- 'd': 1
- 'e': 1
- 'r': 2
- 's': 2
- 't': 3
- 'u': 2



SPC, 1

D, 1

c, 1

e, 1

a, 2

r, 2

s, 2

u, 2

t, 3

Solution

“Data structures”

- SPC: 1
- 'a': 2
- 'c': 1
- 'D': 1
- 'e': 1
- 'r': 2
- 's': 2
- 't': 3
- 'u': 2

Note: In ASCII, Space precedes caps, which precede lowercase!

SPC, 1

D, 1

c, 1

e, 1

a, 2

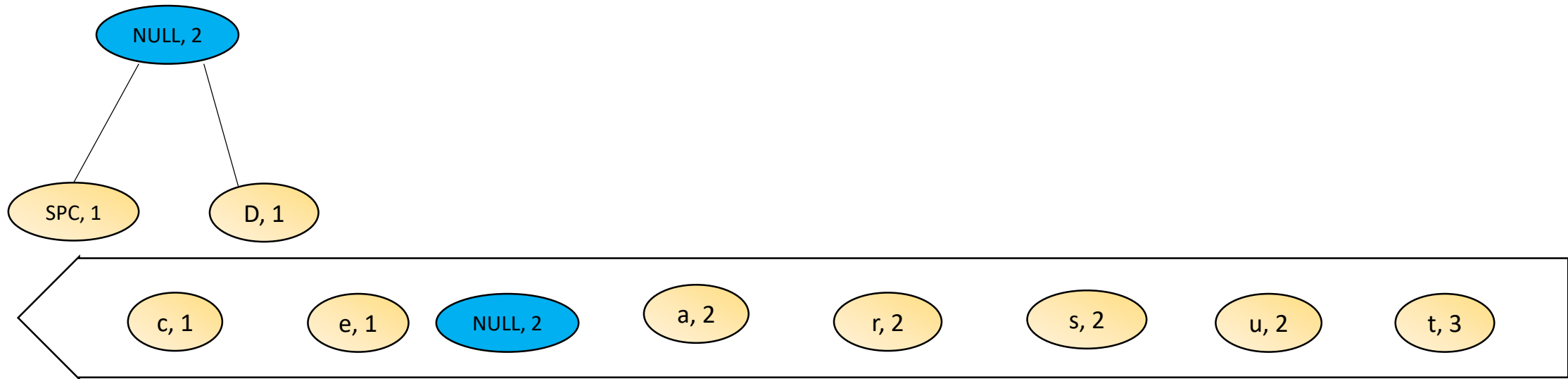
r, 2

s, 2

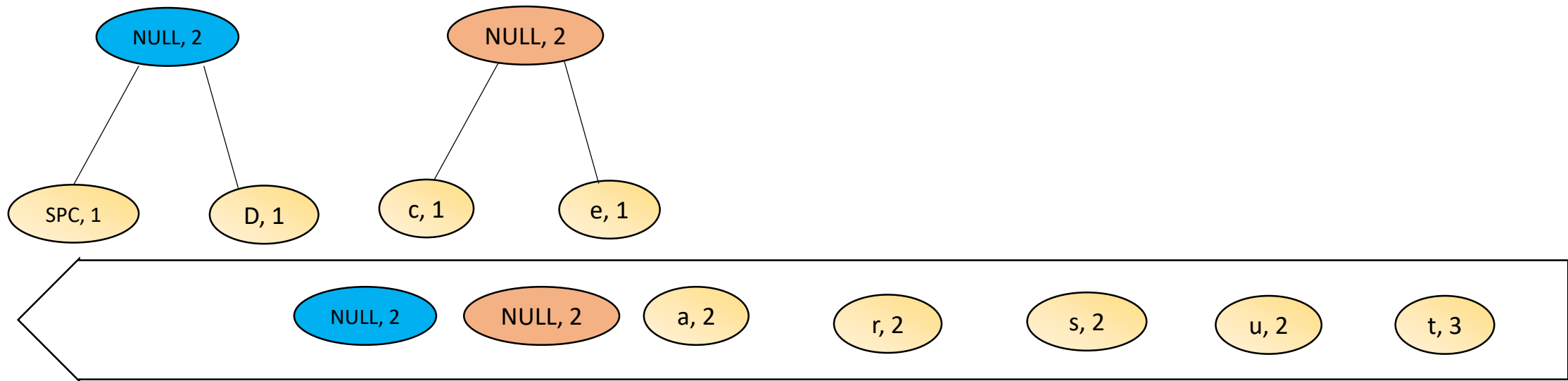
u, 2

t, 3

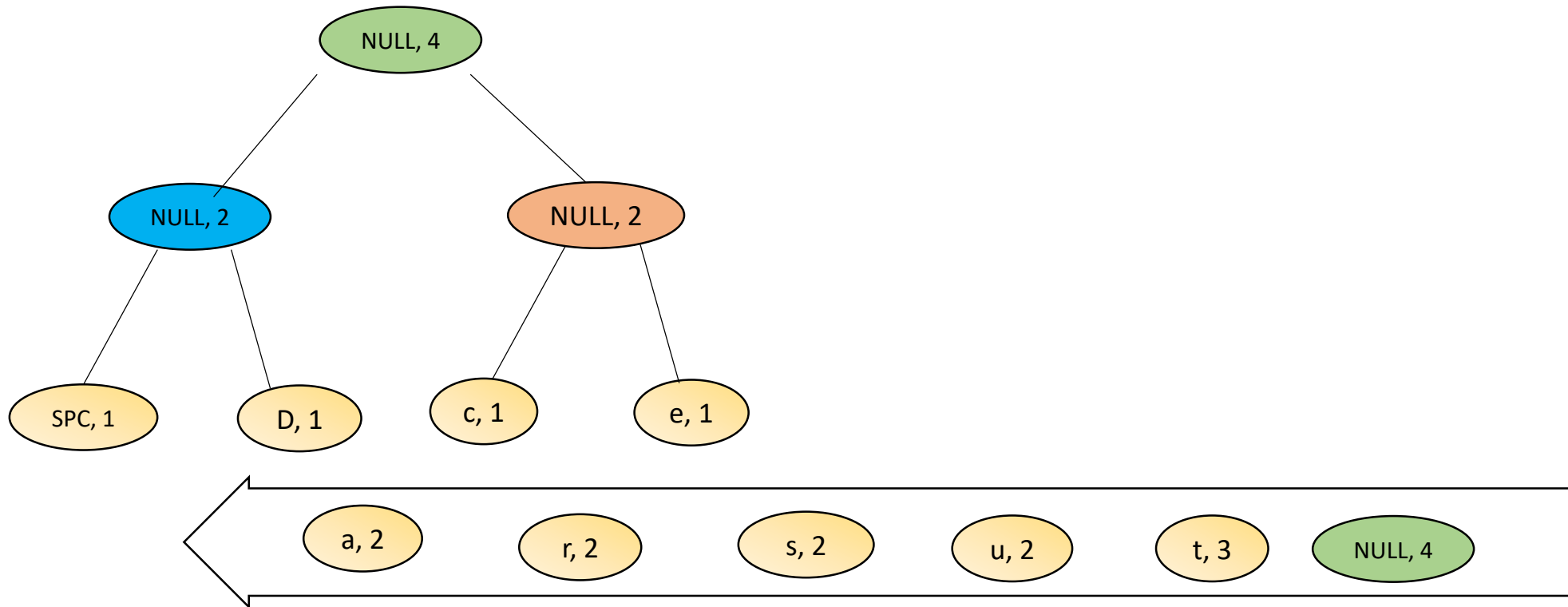
Solution



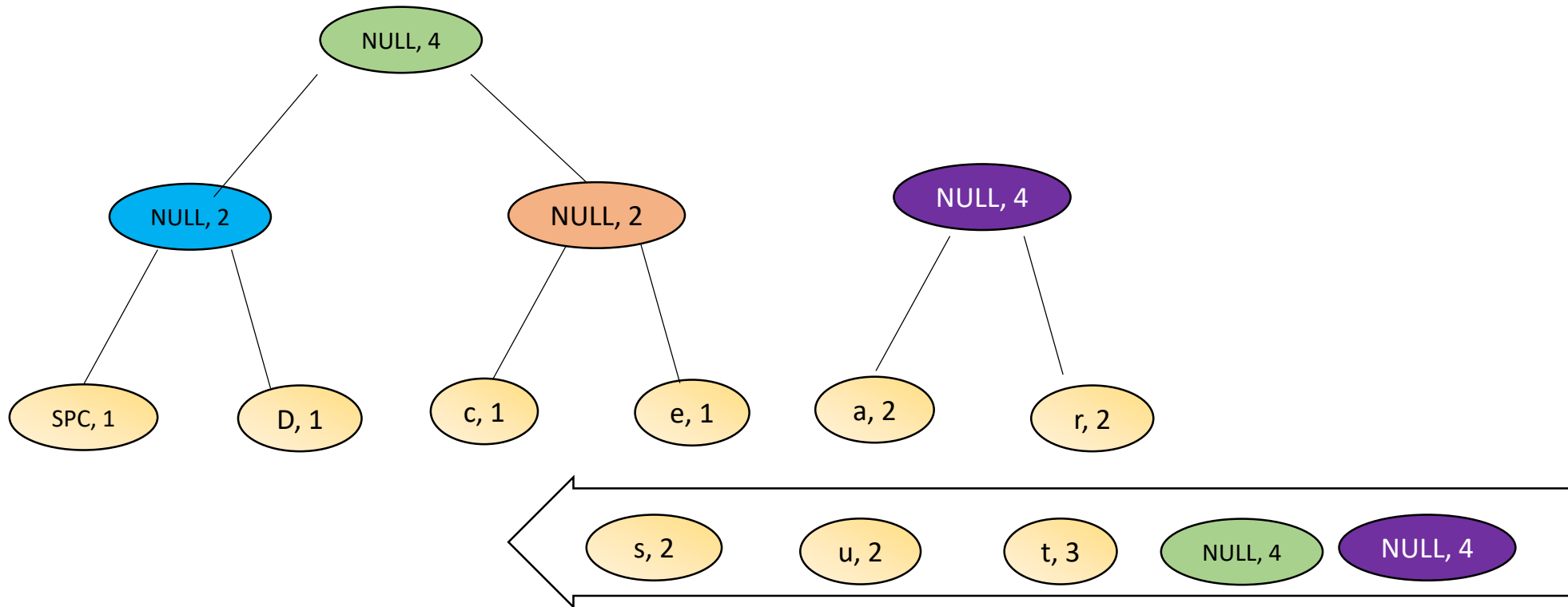
Solution



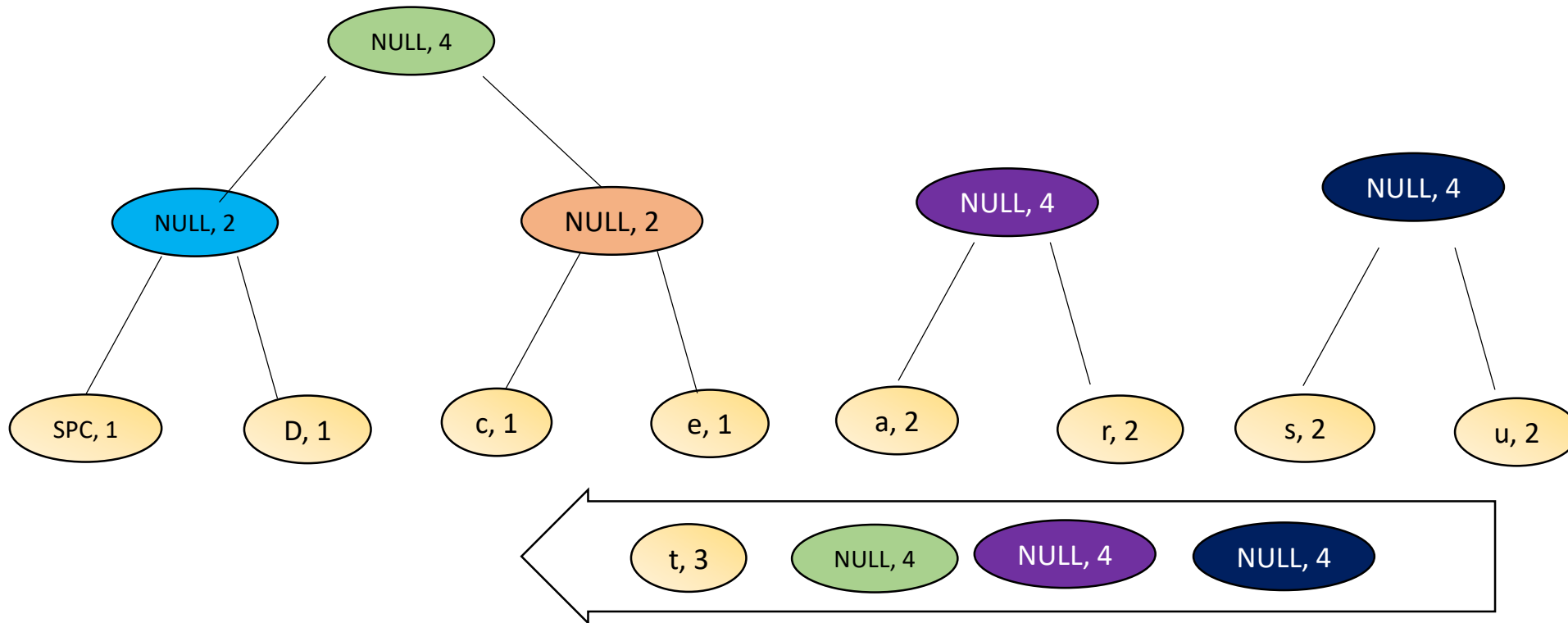
Solution



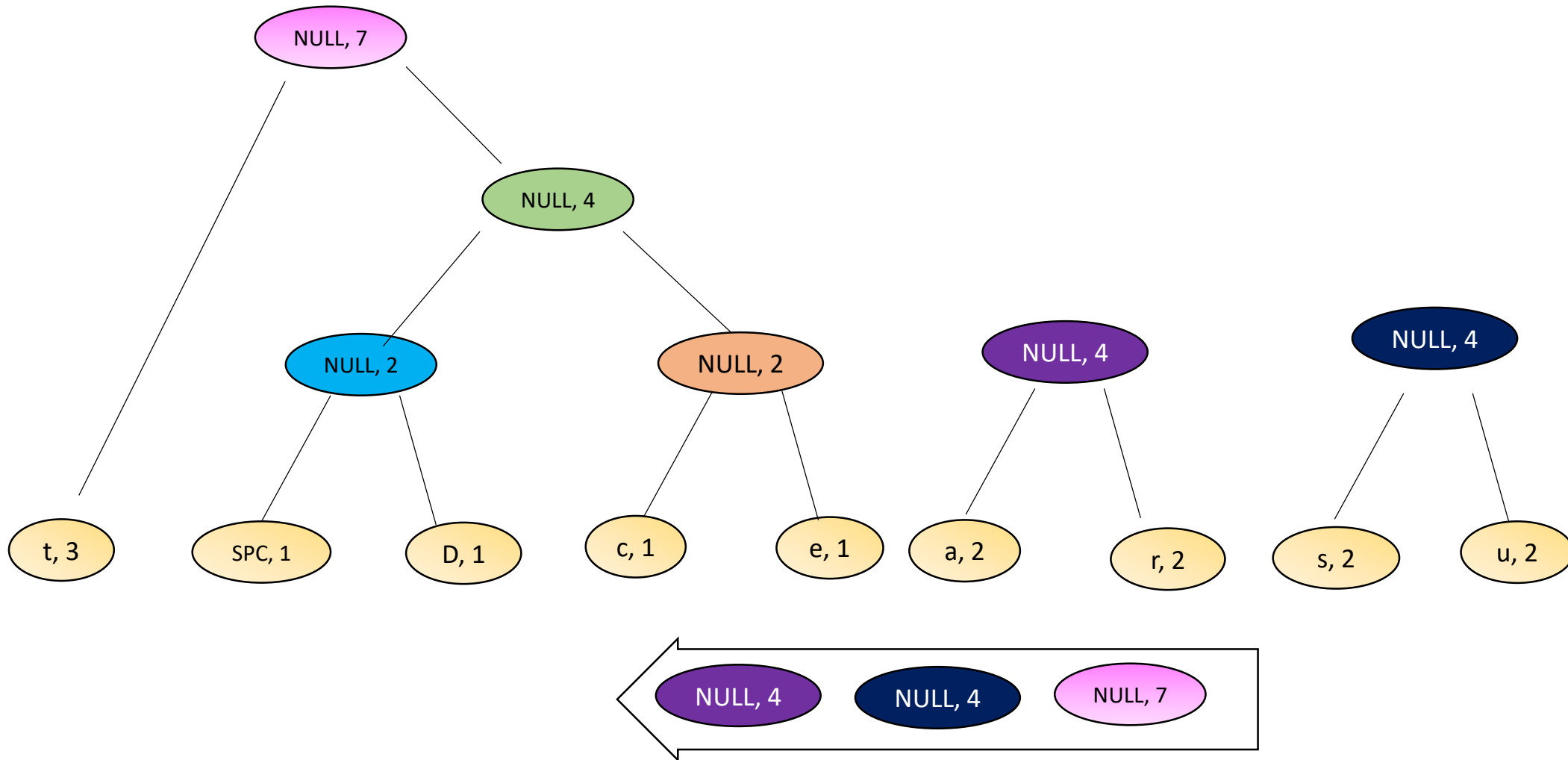
Solution



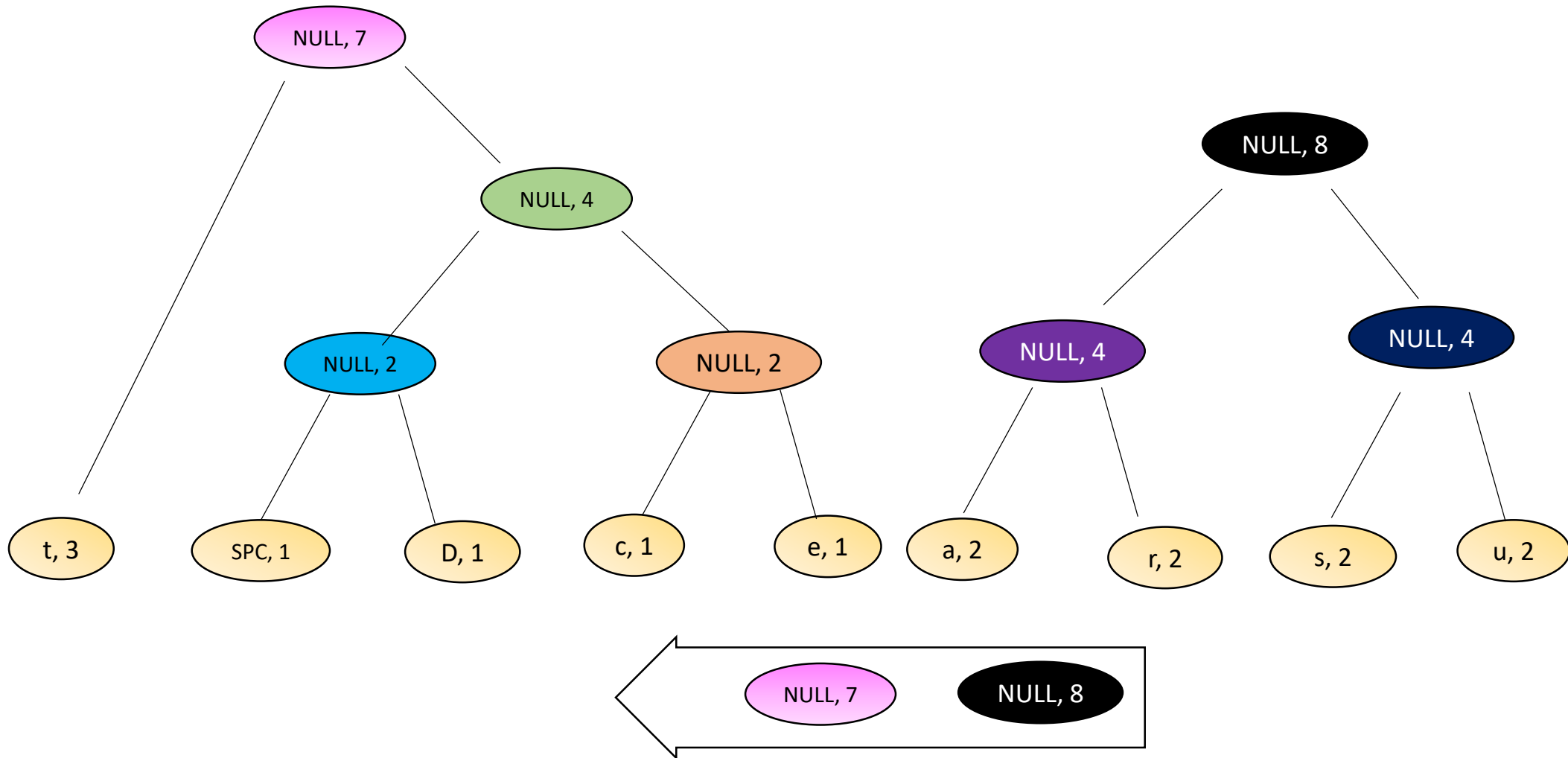
Solution



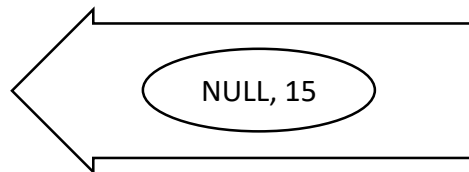
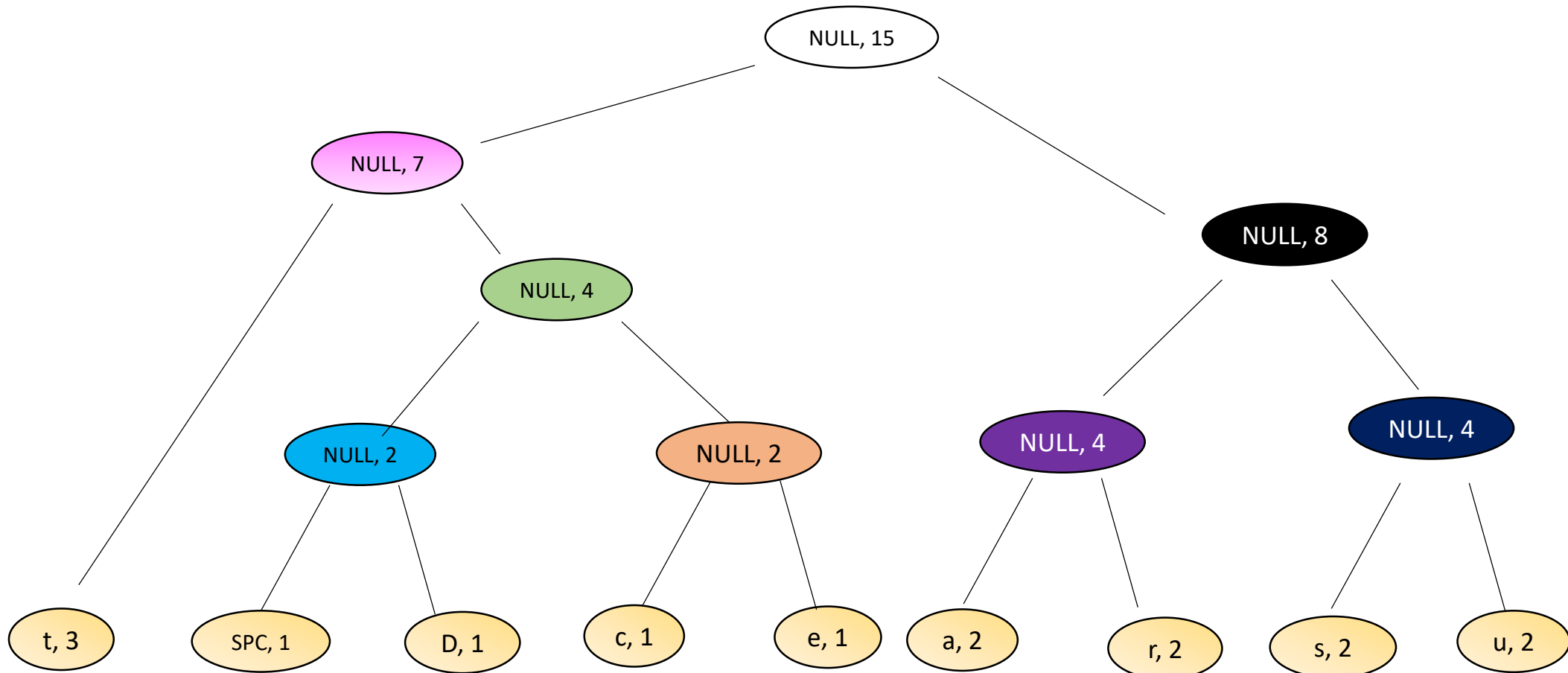
Solution



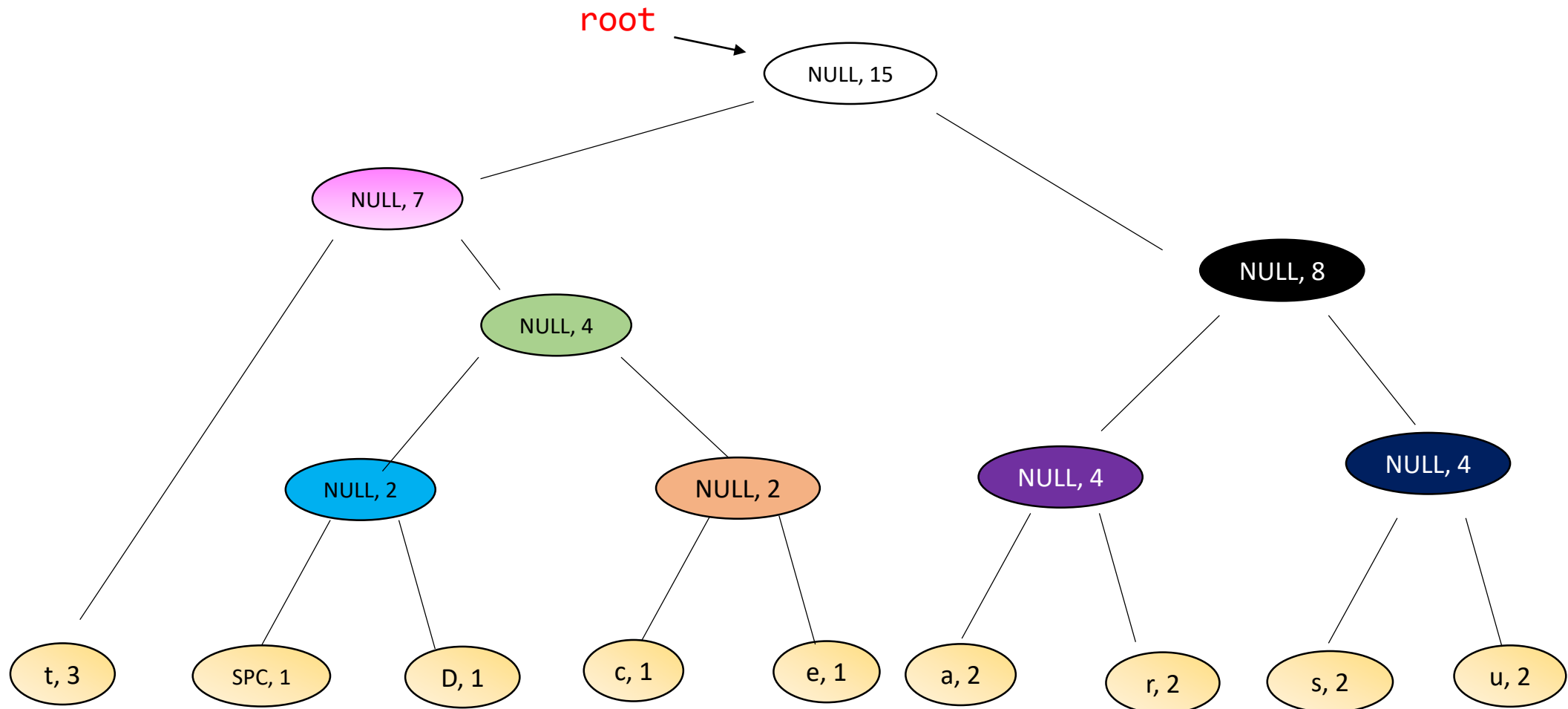
Solution



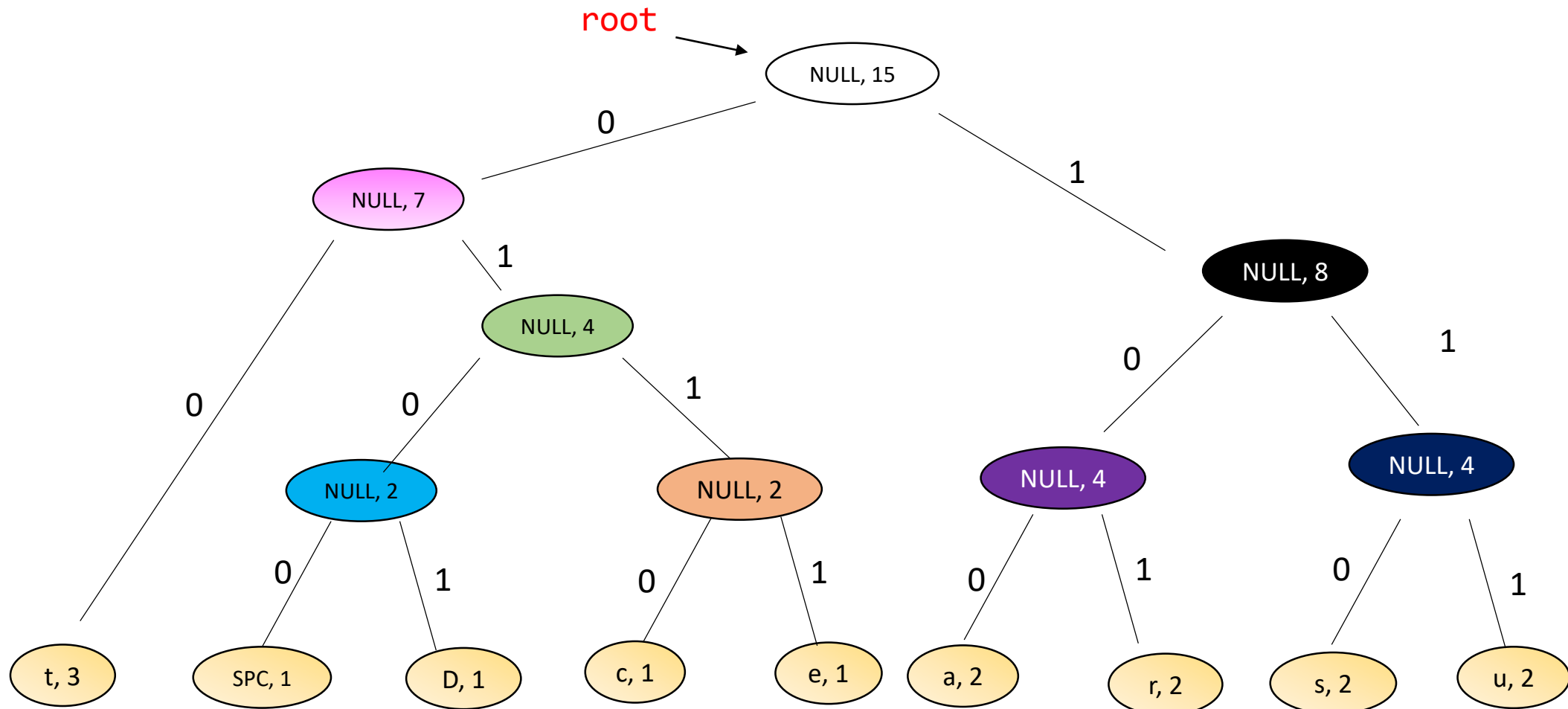
Solution



Solution

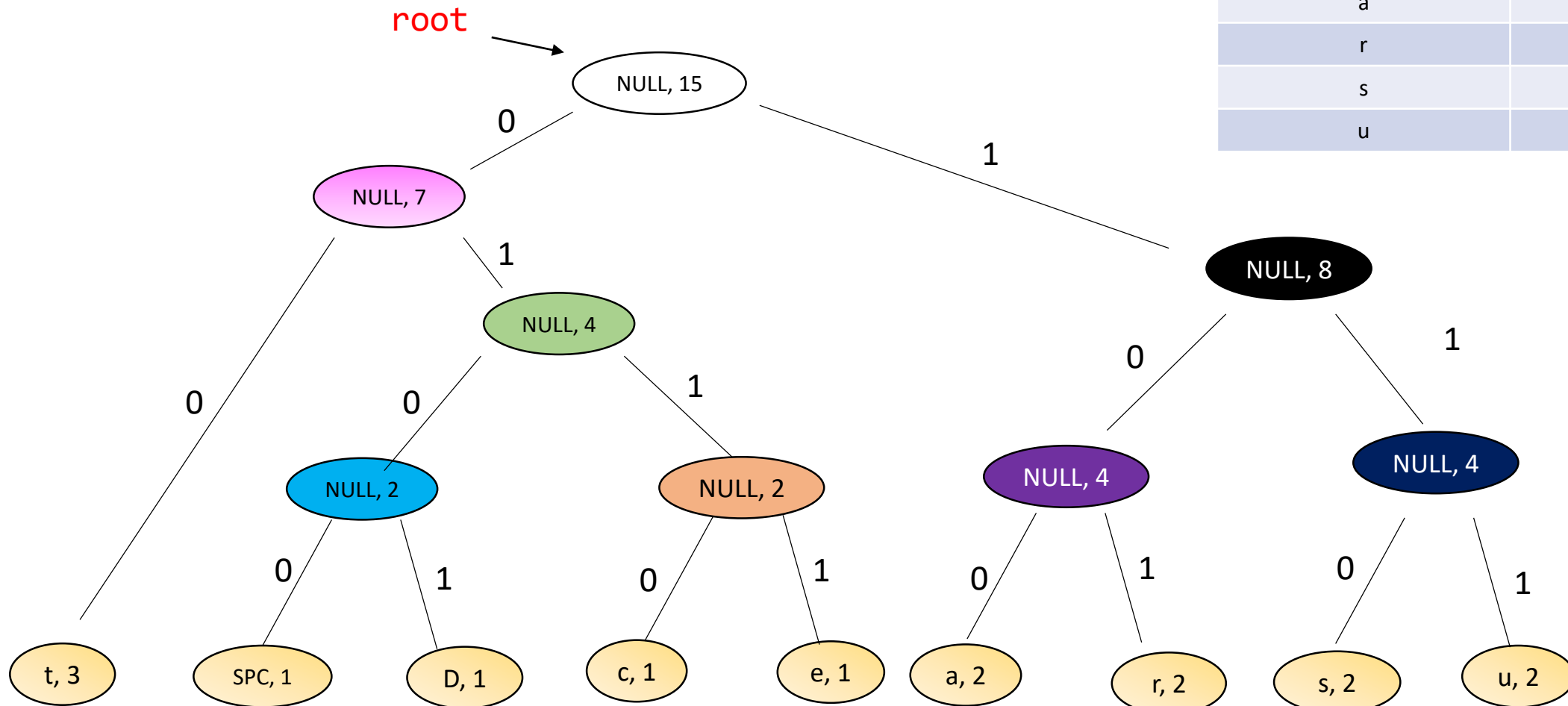


Solution



Solution

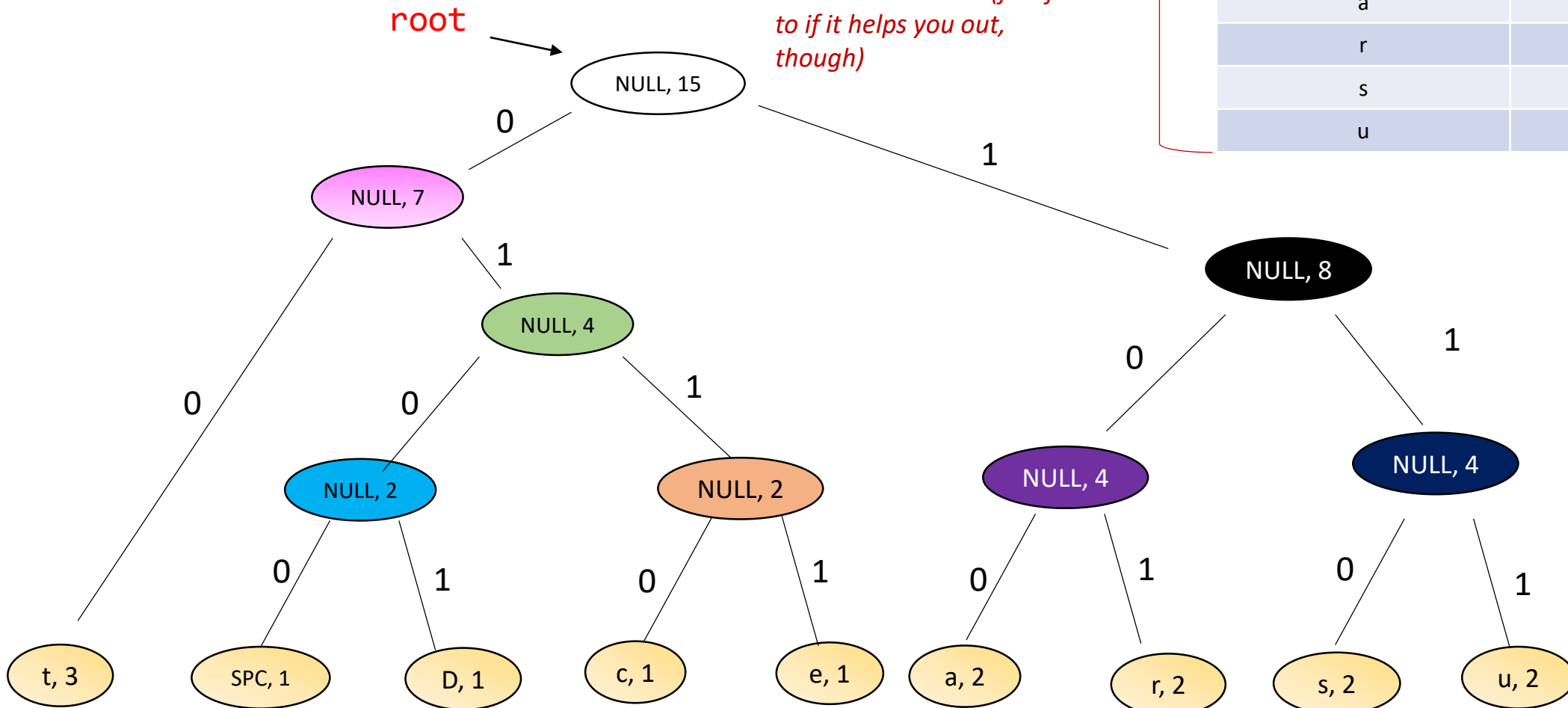
Character	Huffman Encoding
t	00
SPC	0100
D	0101
c	0110
e	0111
a	100
r	101
s	110
u	111



Solution

Remember: This step is necessary to encode the string, but you won't have to do it in an exam (feel free to if it helps you out, though)

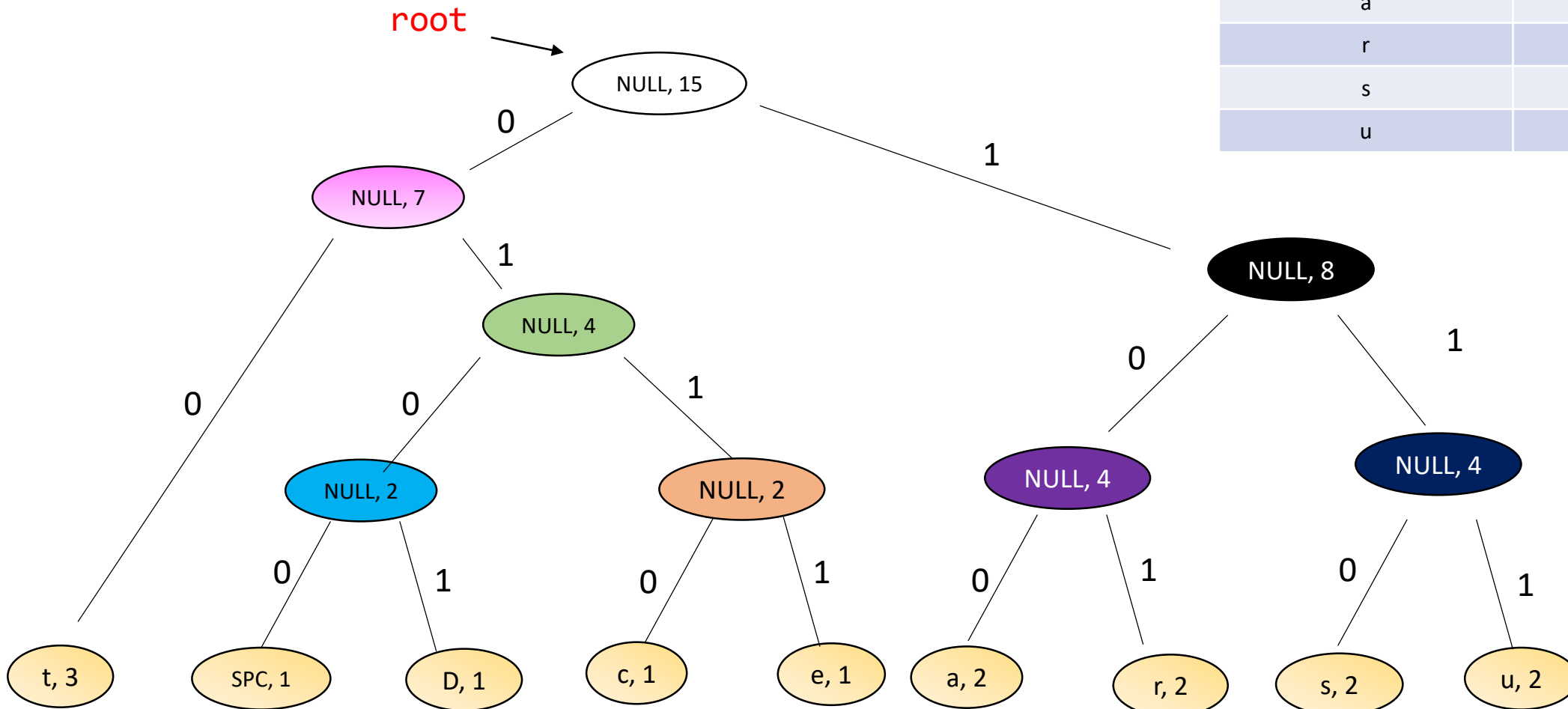
Character	Huffman Encoding
t	00
SPC	0100
D	0101
c	0110
e	0111
a	100
r	101
s	110
u	111



Transmission: "Data structures" =
0101100001000100110001011
110110001111010111110

Solution

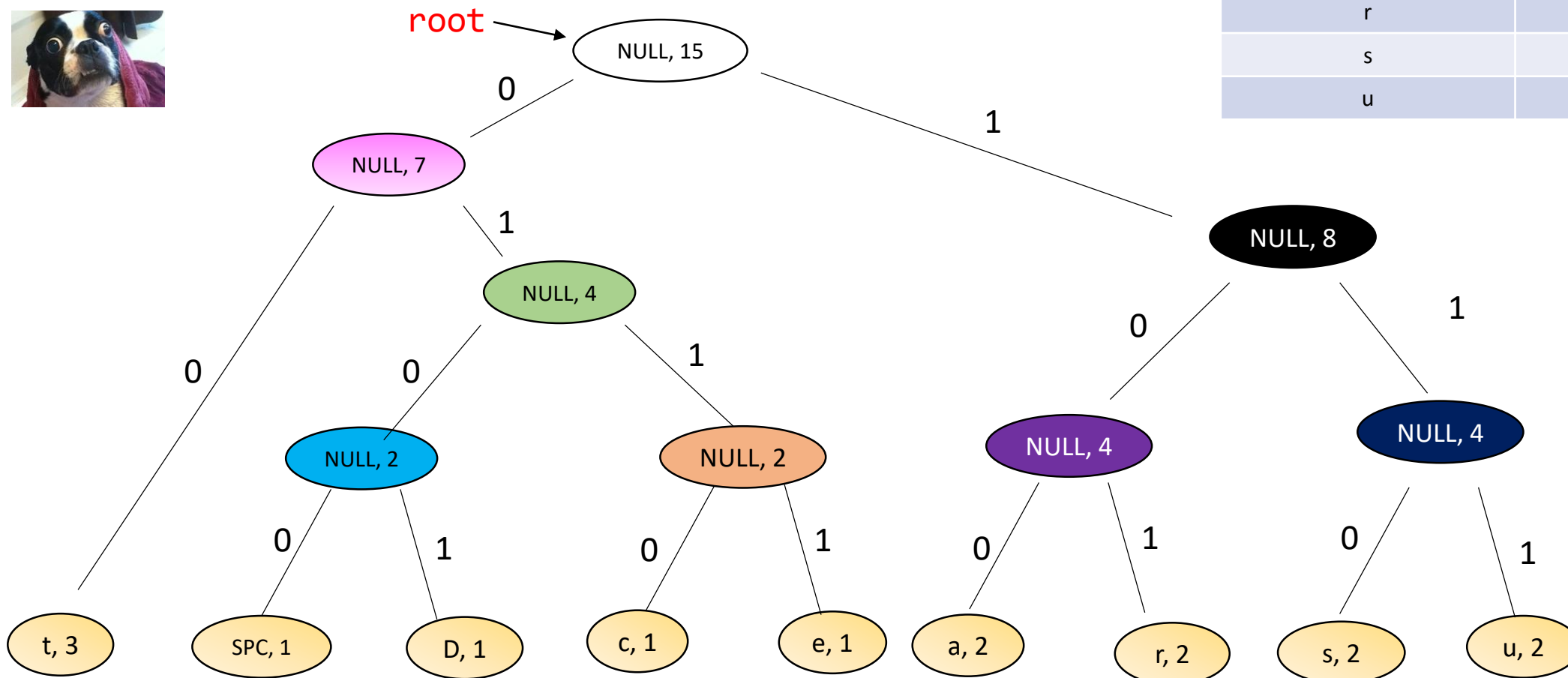
Character	Huffman Encoding
t	00
SPC	0100
D	0101
c	0110
e	0111
a	100
r	101
s	110
u	111



Solution

Transmission: "Data structures" =
0101100001000100110001011
110110001111010111110

46 bits. Compare with $15 * 7 = 105$
for ASCII... (43.8%)



Character	Huffman Encoding
t	00
SPC	0100
D	0101
c	0110
e	0111
a	100
r	101
s	110
u	111

Encoding phase

- Formally, what we did is known as **the encoding phase of Huffman**.
- We encoded characters from some alphabet Σ **based on their frequencies** in a (**hopefully representative!**) text.

Encoding phase

- Formally, what we did is known as **the encoding phase of Huffman**.
- We encoded characters from some alphabet Σ **based on their frequencies** in a (**hopefully representative!**) text.
- **But what if some characters from Σ never appear in the text?**

Encoding phase

- Formally, what we did is known as **the encoding phase of Huffman**.
- We encoded characters from some alphabet Σ **based on their frequencies** in a (**hopefully representative!**) text.
- **But what if some characters from Σ never appear in the text?**
 - Well, then maybe there's something to be said about the "**representativeness**" of the text...

Encoding phase

- Formally, what we did is known as **the encoding phase of Huffman**.
- We encoded characters from some alphabet Σ **based on their frequencies** in a (**hopefully representative!**) text.
- **But what if some characters from Σ never appear in the text?**
 - Well, then maybe there's something to be said about the "**representativeness**" of the text...
 - **Easy solution:** Add 1 to every other node's values and set the values for all the non-appearing characters be the previous minimum value. This we do while we compute the frequencies, by another linear scan over a $|\Sigma|$ –large histogram, not the text T! **So we still only have to scan the text once to build the trie.**

Encoding phase

- **Invariant satisfaction:** Recall; one of our goals was for **no code to be the prefix of another.**

Encoding phase

- **Invariant satisfaction:** Recall; one of our goals was that **no code to be the prefix of another.**
- This is **trivially satisfied by Huffman Encoding**, since **all codes are generated by following paths from the root to the leaves!**
 - And there's nothing "below" the leaves!



Criticism of Huffman's encoding phase

- Huffman has to traverse the string **twice**, **once for frequency computation** and **once for encoding**. 😞
- This is one of the **main criticisms** of Huffman encoding.

Criticism of Huffman's encoding phase

- Huffman has to traverse the string **twice**, **once for frequency computation** and **once for encoding**. 😞
- This is one of the **main criticisms** of Huffman encoding.
- Huffman encoding can be seen as a method for **string compression**, since we save memory space per character.

Criticism of Huffman's encoding phase

- Huffman has to traverse the string **twice**, once for frequency computation and once for encoding. 😞
- This is one of the **main criticisms** of Huffman encoding.
- Huffman encoding can be seen as a method for **string compression**, since we save memory space per character.
 - However, the LZW compression algorithm **does better** 😊