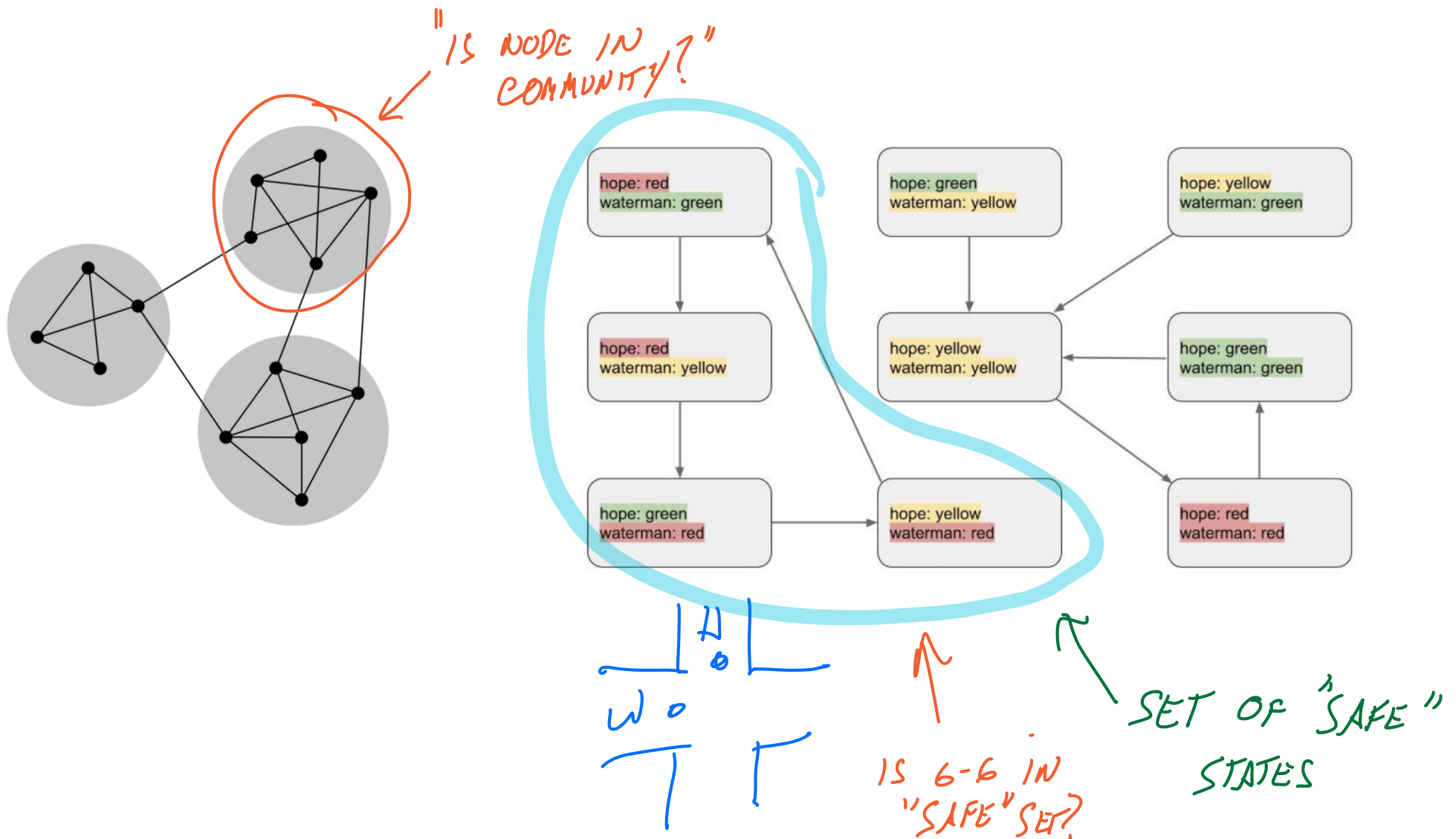


## How would we represent a collection of nodes within, eg., a community?

=> Problem: what if we have a LOT of nodes?



How would represent membership in a collection?

Good initial choice is a hashset...

**Problem: how do we name the nodes in a set?**

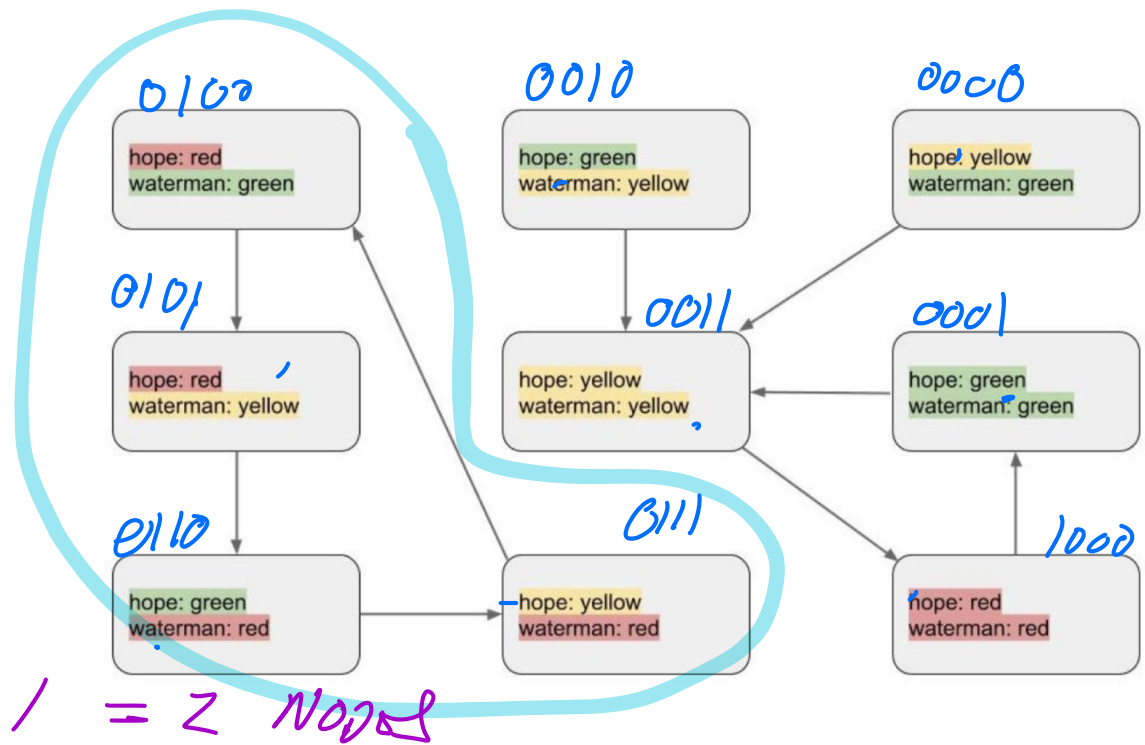
**=> Can be especially problematic when there are a lot of nodes (eg. really large graph)**

- Idea: Could assign a name to each Node object: {"node1", "node2", ...}
- Another idea: Could just use the memory addresses: {@1127, @4452, ...}  
=> Storing a set of numbers (addresses) would use less space than strings)

What else could we do?

Claim: if we get creative about how we name the objects, can represent membership in a space-efficient way.

Idea: smallest piece of information in a computer is a bit (0, 1)  
 - Assign a name as a sequence of bits (bit vector, or a bit string)



1 BIT  $\Rightarrow$  0 1 = 2 Nodes

2 BITS  $\Rightarrow$  00 01 10 11  $\Rightarrow$  4 Nodes.

3 BITS  $\Rightarrow$  000 001 010 011 100 101 110 111  $\Rightarrow$  8 Nodes.

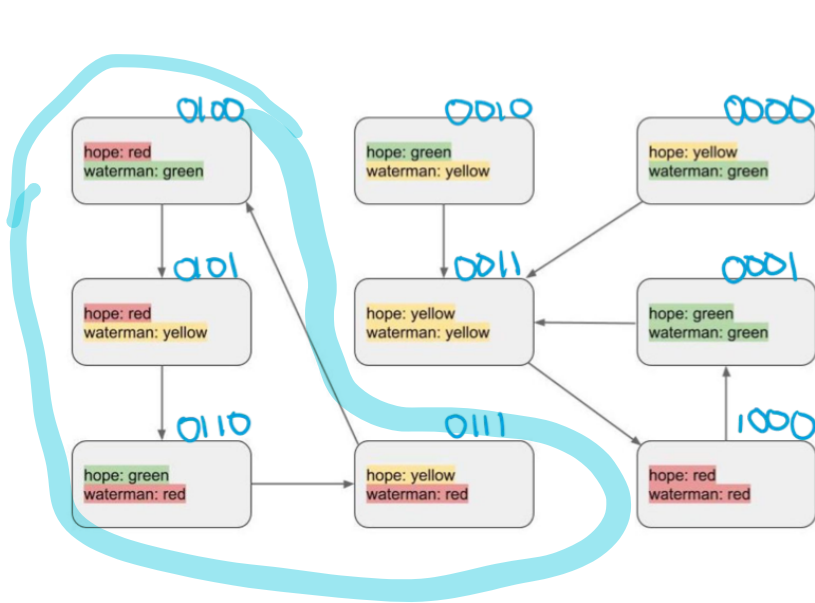
4 BITS

11

16 Nodes

Set[0100, 0101, 0110, 0111]

Claim: if we get creative about how we name the objects, can represent membership in a space-efficient way



①

Idea: give each state a name

- Smallest piece of information in a computer is a bit (0 or 1)
  - Assign name to each state "bit string" (sequence of bits)
- For n states, how many bits?  $\log_2(N)$

2 STATES  $\Rightarrow$  REPRESENT AS BITS: 0, 1  $\Rightarrow$  1 BIT

4 STATES  $\Rightarrow$  4 POSSIBLE BITSTRINGS: 00, 01, 10, 11  $\Rightarrow$  2 BITS

8 STATES  $\Rightarrow$  8 POSSIBLE BITSTRINGS: 000, 001, 010, 011, 100, 101, 110, 111  $\Rightarrow$  3 BITS

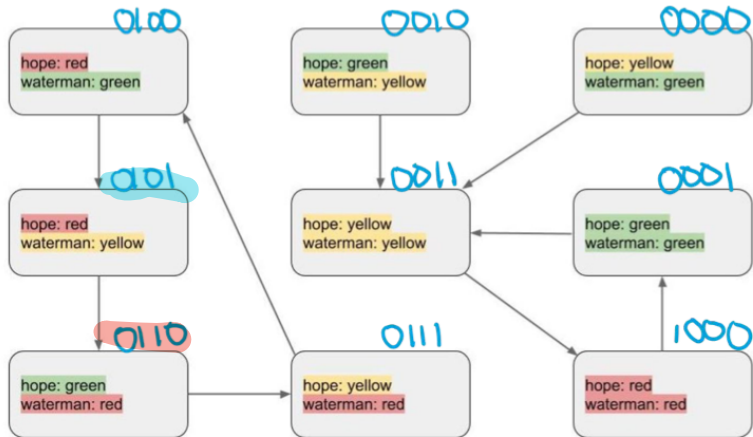
In general:  $N \text{ bits} \Rightarrow 2^N \text{ possible "things" (in this case states)}$   
OR

$M \text{ states} \Rightarrow \log_2(M) \text{ bits}$   
 $\log_2(9) = 3.1 \Rightarrow 4$  (need all names to be same length)

Using bitstrings would let us make the set:  
Set[0100, 0101, 0110, 0111]

$\Rightarrow$  This is much smaller than using objects or addresses, but still grows linearly!

Can we do better?



=> Using this idea, we can represent each state as a bitstring of 0's and 1's

0 1 0 0  
                     
 b<sub>1</sub> b<sub>2</sub> b<sub>3</sub> b<sub>4</sub>

> Notation: we write state of individual bit as b<sub>0</sub>, b<sub>1</sub>, ..., where b<sub>0</sub> could be 0 or 1 (like a Boolean variable)

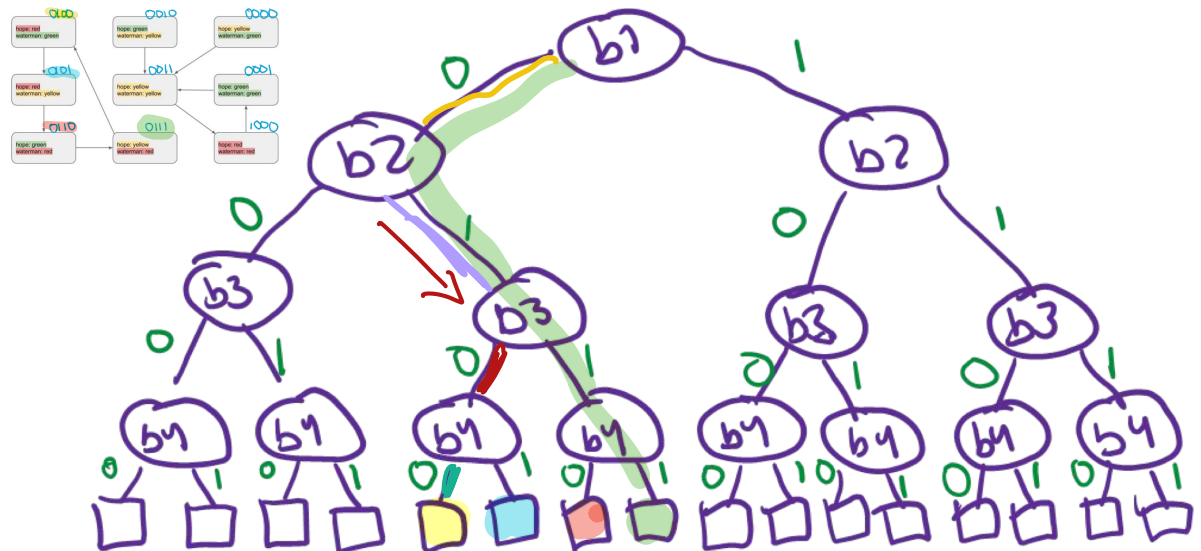
What can we do with this?

- Represent states VERY succinctly (lower space complexity)
- Can represent set of reachable states as a decision tree => Binary Decision Diagram (BDD)

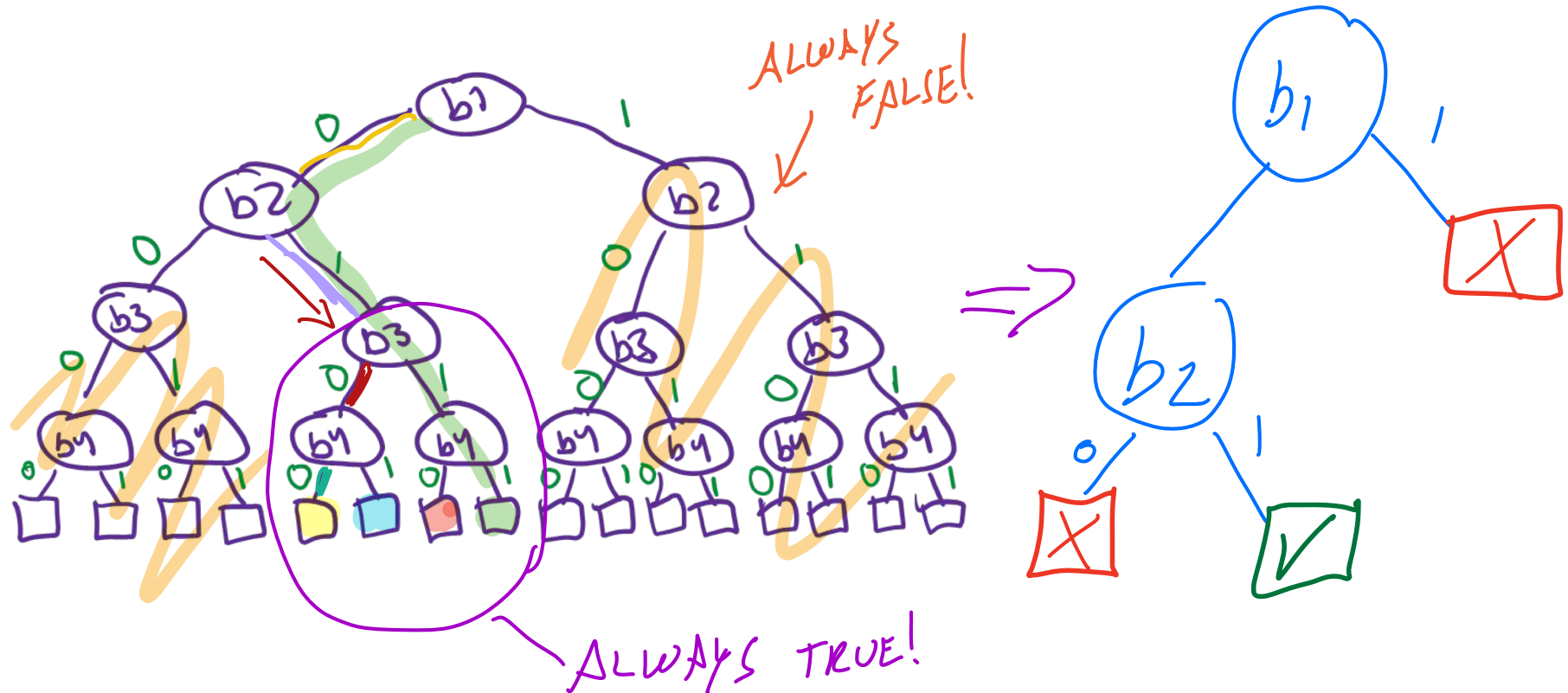
BDD: How it works (high level)

- Each bit is like a step in a decision tree: value is either 0 or 1. After that, consider the next bit
- Each "leaf" is a decision on whether or not the node is in the set

- Can make this very compact...



If we are clever about how to select the names, can reduce a large decision tree into a small one (similar to the decision tree project!). For example, if  $b1 == 1 \Rightarrow$  always false



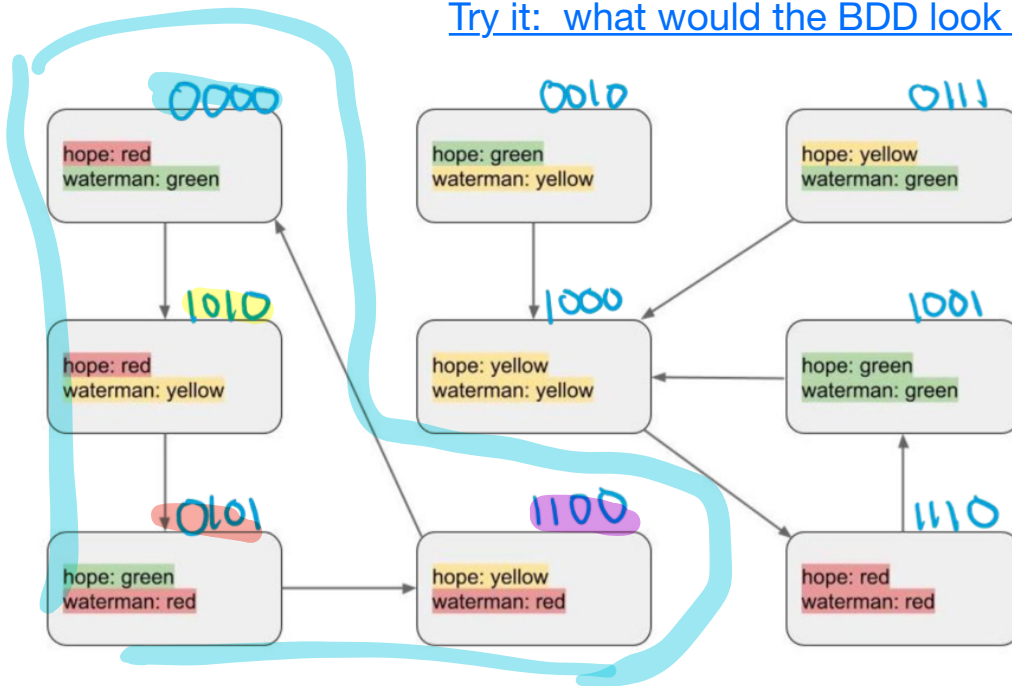
Can also write this as a formula:

Reachable:  $(b1 == 0) \text{ AND } (b2 == 1)$

Could also write as:  $(\text{not } b1) \text{ and } (b2)$

**$\Rightarrow$  To check if the state is safe, Don't even need to store a list of reachable states! Just need to check if the bitstring matches the formula!**

Try it: what would the BDD look like if we picked these names instead?



0 0 0 0  
B<sub>1</sub> B<sub>2</sub> B<sub>3</sub> B<sub>4</sub>

Formula would be a lot more complicated!

((not b<sub>1</sub>) and (not b<sub>2</sub>) and (not b<sub>3</sub>) and (not b<sub>4</sub>)) OR  
((b<sub>1</sub>) and (not b<sub>2</sub>) and (b<sub>3</sub>) and (not b<sub>4</sub>)) OR

...

### Takeaways:

- Coming up with a good naming of states is nontrivial

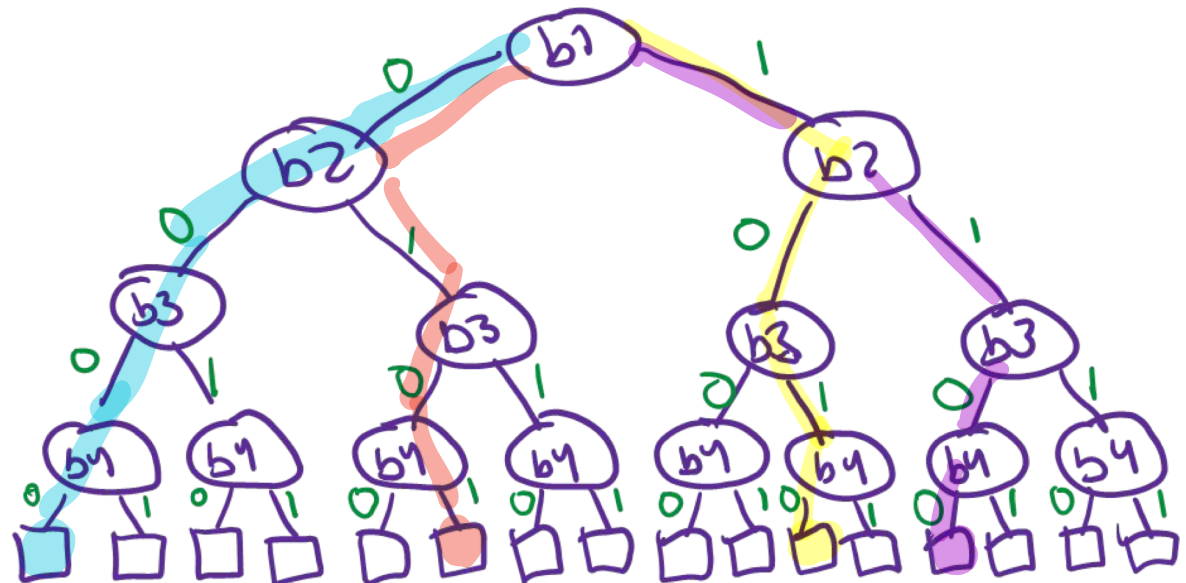
=> Problems are provably computationally hard, people work on this for specific settings

- Depending on what the BDD looks like, logical formula may get complex

Want to learn more?

Consider taking: Logic for Systems

=> Learn about algorithms people have made to represent these in more concise ways, make tools to help, software to use





## Compression: how to store text efficiently?

"NATS"  
↙ 8 BITS/CHAR (ASCII)  
32 BITS/CHAR (UNICODE)

ASCII: Use 8 bit bitstring for each character  
=>  $4 * 8 \text{ bits} = 32 \text{ bits total}$

Storage: increases linearly with number of characters

BigWiki: 141 million characters \* 8 bits  $\approx$  136MB

Idea: language (and most data) isn't random

Compression: use fewer bits to represent data that appears more frequently

"e" appears more often than "x"—so maybe we can use fewer bits?

### Huffman coding:

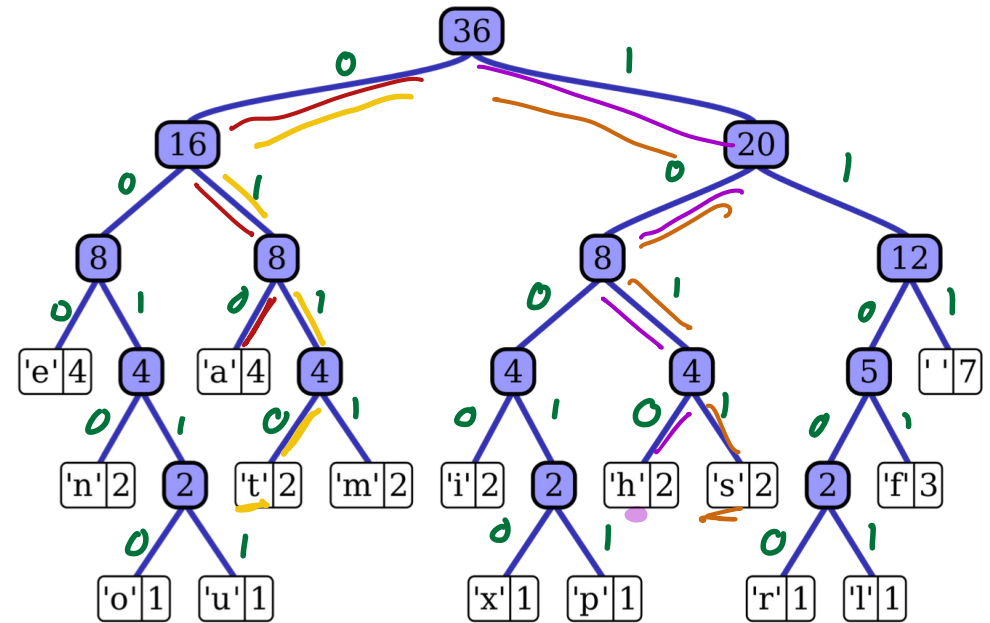
- for some input data, find unique bitstring to represent characters, use fewer bits for more common letters

### Takeaways:

=> Leveraging patterns in data for efficiency

=> Encoding pattern in BDD structure (use to convert)

"this is an example of a huffman tree"



[https://en.wikipedia.org/wiki/Huffman\\_coding](https://en.wikipedia.org/wiki/Huffman_coding)

"NATS"

h: 1010 4

a: 010 3

t: 0110 4

s: 1011 4

=> 15 BITS

de ↕		
a	4	010
e	4	000
f	3	1101
h	2	1010
i	2	1000
m	2	0111
n	2	0010
s	2	1011
t	2	0110
l	1	11001
o	1	00110
p	1	10011
r	1	11000
u	1	00111
x	1	10010