



University of
Pittsburgh

Algorithms and Data Structures 2

CS 1501



Fall 2022

Sherif Khattab

ksm73@pitt.edu

(Slides are adapted from Dr. Ramirez's and Dr. Farnan's CS1501 slides.)

Announcements

- Upcoming Deadlines
 - Homework 7: next Friday @ 11:59 pm
 - Lab 6: Monday 10/31 @ 11:59 pm
 - Nothing due this week
- Midterm Exam
 - Wednesday 10/19 (MW Section) and Thursday 10/20 (TuTh Section)
 - in-person, closed-book
- Weekly Live QA Session on Piazza
 - Friday 4:30-5:30 pm

Previous lecture

- ADT Priority Queue (PQ)
 - Heap implementation
- Heap Sort
- Indexable PQ

This Lecture

- Muddiest Points
- Introduction to ADT Graph

Muddiest Points

- **Q: what is going to be on the midterm?**
- Up to and including material covered on Monday 10/10 and Tuesday 10/11
- Please check the study guide on Canvas
 - practice test on GradeScope
 - old exam
 - will try to post the answer as soon as possible

Muddiest Points

- **Q: I'm confused about entropy and the equations for that. How is it useful?**
- It helps us determine the “information content” in a file.
- In lossless compression, a file cannot be compressed to less than its entropy

Muddiest Points

- **Q: How do we determine the entropy of a given file we're trying to compress?**
- Given that
 - the file has n total characters and K unique characters,
 - $f(c)$ is the frequency of character c in the file
- Shannon's entropy: $H(\text{file}) = -1 * \sum_{c=0}^K \frac{f(c)}{n} \log_2\left(\frac{f(c)}{n}\right)$ bits/character
- Underlying assumption:
 - the file has been generated by a source that produces **independent characters**
 - **Huffman Compression is optimal under that assumption**
 - This assumption may be wrong though!
 - repeated long strings \rightarrow use LZW
 - long sequences of identical characters \rightarrow use RLE (Run Length Encoding)
 - We may need to try different compression algorithms on the file

Muddiest Points

- **Q: What might be some examples in where we might pick one compression type over the other?**
- Depends on the structure of the input file
- long strings of identical values
 - → RLE
- repeated long strings
 - → LZW
- frequently occurring values
 - → Huffman
- few different characters
 - → fixed-length codewords with < 8 bits

Muddiest Points

- **Q: Does entropy play a role in the implementation of the compression algorithm itself? Or is it only used for selecting the best algorithm?**
- Some compression algorithms attempt to reach Shannon's Entropy lower bound on the file size
 - e.g., Huffman Encoding and Arithmetic Encoding

Muddiest Points

- **Q: Not exactly clear on move to front encoding**

Input:

e a e d e e

0

a

1

b

2

c

3

d

4

e

Output:

Muddiest Points

- **Q: Not exactly clear on move to front encoding**

Input:

e a e d e e

0

a

1

b

2

c

3

d

4

e

Output:

4

Muddiest Points

- **Q: Not exactly clear on move to front encoding**

Input:

e a e d e e

0

a

e

1

b

a

2

c

b

3

d

c

4

e

d

Output:

4

Muddiest Points

- **Q: Not exactly clear on move to front encoding**

Input:

e a e d e e

0

a

e

a

1

b

a

e

2

c

b

b

3

d

c

c

4

e

d

d

Output:

4

1

Muddiest Points

- **Q: Not exactly clear on move to front encoding**

Input:

e a e d e e

0	a	e	a	e
1	b	a	e	a
2	c	b	b	b
3	d	c	c	c
4	e	d	d	d

Output:

4 1 1

Muddiest Points

- **Q: Not exactly clear on move to front encoding**

Input:

e a e d e e

0	a	e	a	e	d
1	b	a	e	a	e
2	c	b	b	b	a
3	d	c	c	c	b
4	e	d	d	d	c

Output:

4 1 1 4

Muddiest Points

- **Q: Not exactly clear on move to front encoding**

Input:

e a e d e e

0	a	e	a	e	d	e
1	b	a	e	a	e	d
2	c	b	b	b	a	a
3	d	c	c	c	b	b
4	e	d	d	d	c	c

Output:

4 1 1 4 1

Muddiest Points

- **Q: Not exactly clear on move to front encoding**

Input:

e a e d e e

0	a	e	a	e	d	e	e
1	b	a	e	a	e	d	d
2	c	b	b	b	a	a	a
3	d	c	c	c	b	b	b
4	e	d	d	d	c	c	c

Output:

4 1 1 4 1 0

Muddiest Points

- **Q: Not exactly clear on move to front encoding**
- **Decoding**

Input:

4 1 1 4 1 0

0

a

1

b

2

c

3

d

4

e

Output:

e

Muddiest Points

- **Q: Not exactly clear on move to front encoding**
- **Decoding**

Input:

4 1 1 4 1 0

0	a	e
1	b	a
2	c	b
3	d	c
4	e	d

Output:

e

Muddiest Points

- **Q: Not exactly clear on move to front encoding**
- **Decoding**

Input:

4 1 1 4 1 0

0	a	e	a
1	b	a	e
2	c	b	b
3	d	c	c
4	e	d	d

Output:

e a

Muddiest Points

- **Q: Not exactly clear on move to front encoding**
- Decoding

Input:

4 1 1 4 1 0

0	a	e	a	e
1	b	a	e	a
2	c	b	b	b
3	d	c	c	c
4	e	d	d	d

Output:

e a e

Muddiest Points

- **Q: Not exactly clear on move to front encoding**
- Decoding

Input:

4 1 1 4 1 0

0	a	e	a	e	d
1	b	a	e	a	e
2	c	b	b	b	a
3	d	c	c	c	b
4	e	d	d	d	c

Output:

e a e d

Muddiest Points

- **Q: Not exactly clear on move to front encoding**
- Decoding

Input:

4 1 1 4 1 0

0	a	e	a	e	d	e
1	b	a	e	a	e	d
2	c	b	b	b	a	a
3	d	c	c	c	b	b
4	e	d	d	d	c	c

Output:

e a e d e

Muddiest Points

- **Q: Not exactly clear on move to front encoding**
- Decoding

Input:

4 1 1 4 1 0

0	a	e	a	e	d	e	e
1	b	a	e	a	e	d	d
2	c	b	b	b	a	a	a
3	d	c	c	c	b	b	b
4	e	d	d	d	c	c	c

Output:

e a e d e e

Muddiest Points

- **Q: LZW doesn't seem like prefix-free compression. Do we use delimiters between the codes or is there some way to tell the numbers apart. Even just 255 contains multiple interpretations of ASCII codes (2, 55 or 25, 5)**
- Great Question!
- Each integer is the same number of bits (e.g., 12 bits)
- So, the expansion program reads 12 bits at a time
- No need for delimiters nor prefix-free encoding of the integers

Muddiest Points

- **Q: calculating the bits for the lzw compression (last question on the tophat), we didn't go over this during the example in class**
- The second column represents the output of LZW compression, i.e., the compressed file
- Each integer is 12-bit codeword (per the question)
- Total compressed file size = # codewords * 12

Muddiest Points

- **Q: Corner case of lzw expansion**
- The tricky (corner) case happens when the longest match in compressions happens to be the string that was just added in the previous step
- Expansion sees that as a codeword that is not (yet) in its codebook
- Remember that expansion builds the same codebook as compression but is one step behind
- Handling the tricky case:
 - output: previous output + first character of previous output
 - add the same string to the codebook

LZW corner case example

- Compress, using 12 bit codewords: AAAAAA

Cur	Output	Add
A	65	AA:256
AA	256	AAA:257
AAA	257	--

- Expansion:

Cur	Output	Add
65	A	--
256	AA	256:AA
257	AAA	257:AAA

Muddiest Points

- **Q: How can the uncompressed file have more entropy than compressed if the entropy is the average number of bits to represent a word?**
- In lossless compression,
 - entropy of compressed file \geq entropy of uncompressed file
- Since compressed file has fewer characters than uncompressed
 - entropy/char of compressed file is $>$ entropy/char of uncompressed file

Muddiest Points

- **Q: I was wondering why LZW choose 12bits instead of any other number**
- 12 bits was used just as an example
- Actual implementation use an adaptive codeword size
 - start with 9 bits
 - when codebook full, change to 10 bits
 - when codebook full, change to 11 bits
 - ...
 - when codebook full and codeword is 16 bits
 - either stop adding to codebook or reset it
 - depends on the compression ratio

Muddiest Points

- **Q: please make assignments easier**
- Yep!

Muddiest Points

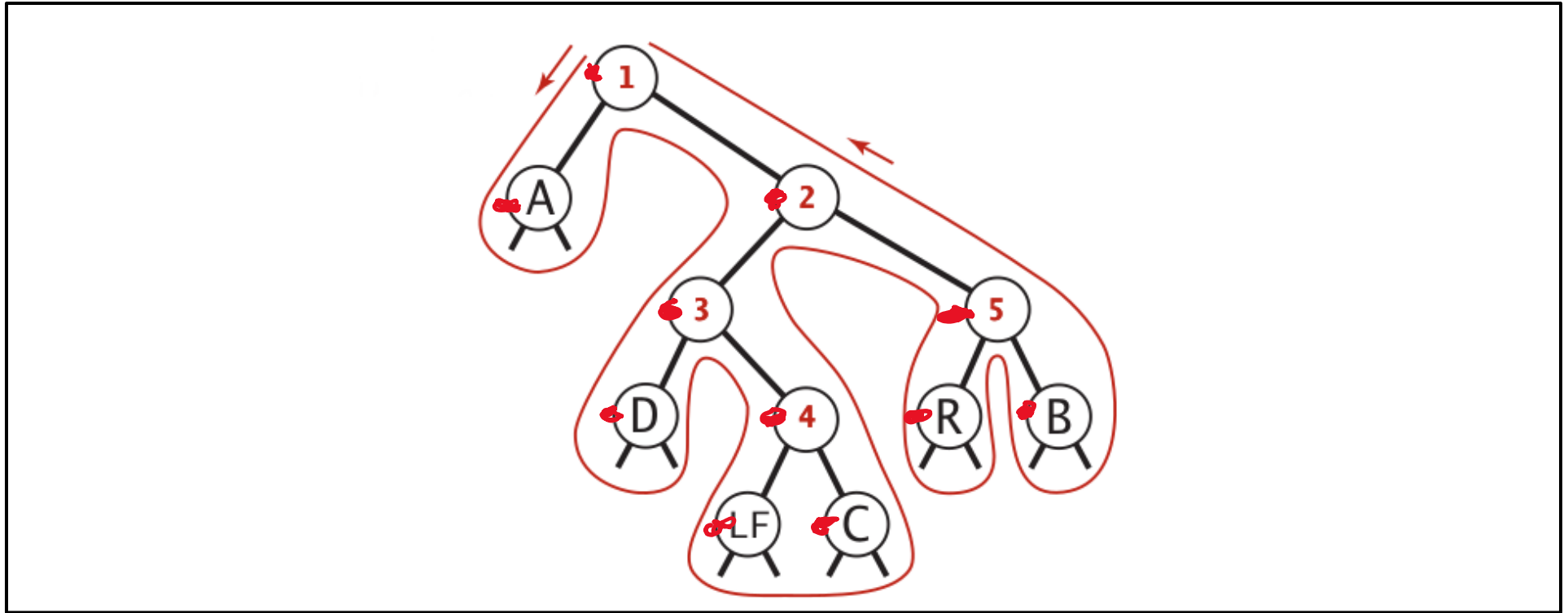
- **Q: Is there a best way to handle ties when doing Huffman compression or is it a purely arbitrary choice?**
- **Q: what are the rules for drawing out the tree in the huffman approach?**
- Ties are arbitrarily handled
- The compressed file size is the same no matter how ties are handled

Muddiest Points

- **Q: the bits needed for the compression post huffman encoding**
- Huffman compression may store the trie in the compressed file
- The trie is encoded using preorder traversal of the nodes
 - encoding internal nodes with 0
 - leaf nodes with 1 followed by the ASCII code of the character inside the leaf

Representing tries as bitstrings

Preorder traversal

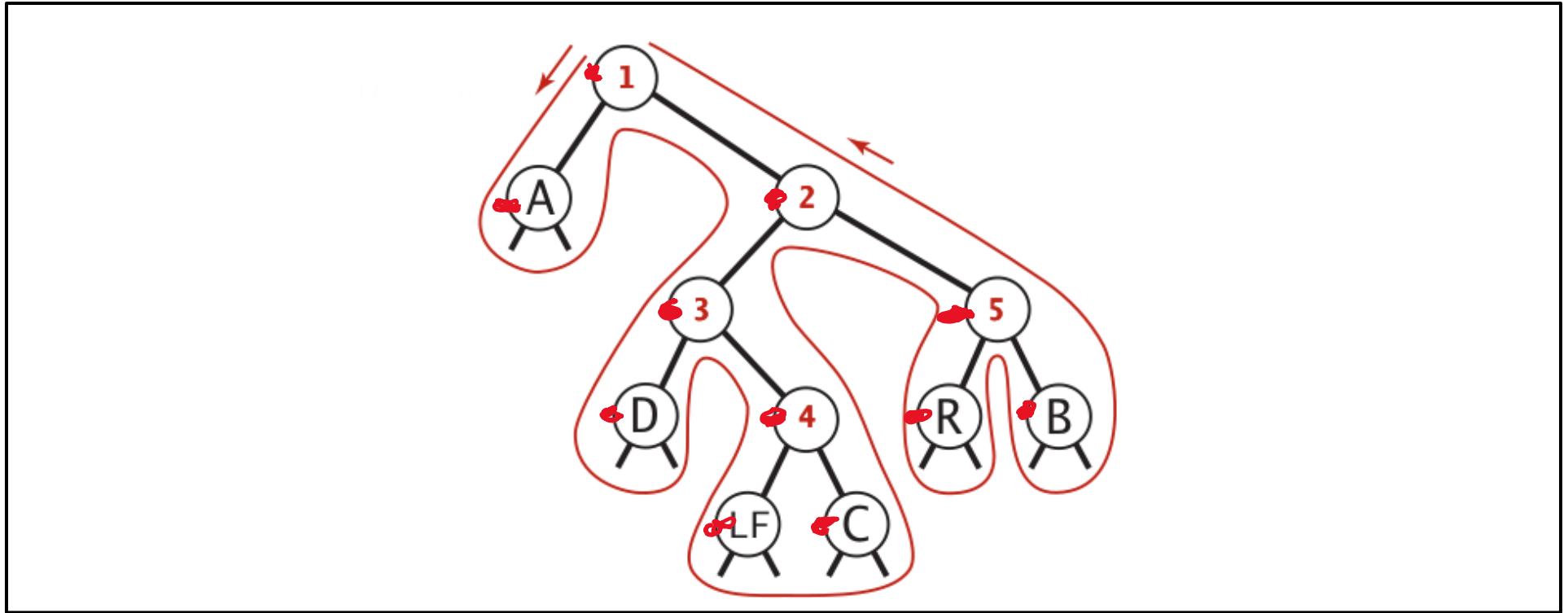


internal node $\rightarrow 0$

leaf node $\rightarrow 1$ followed by ASCII code of char inside

Representing tries as bitstrings

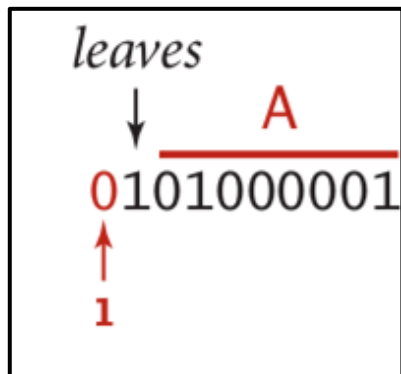
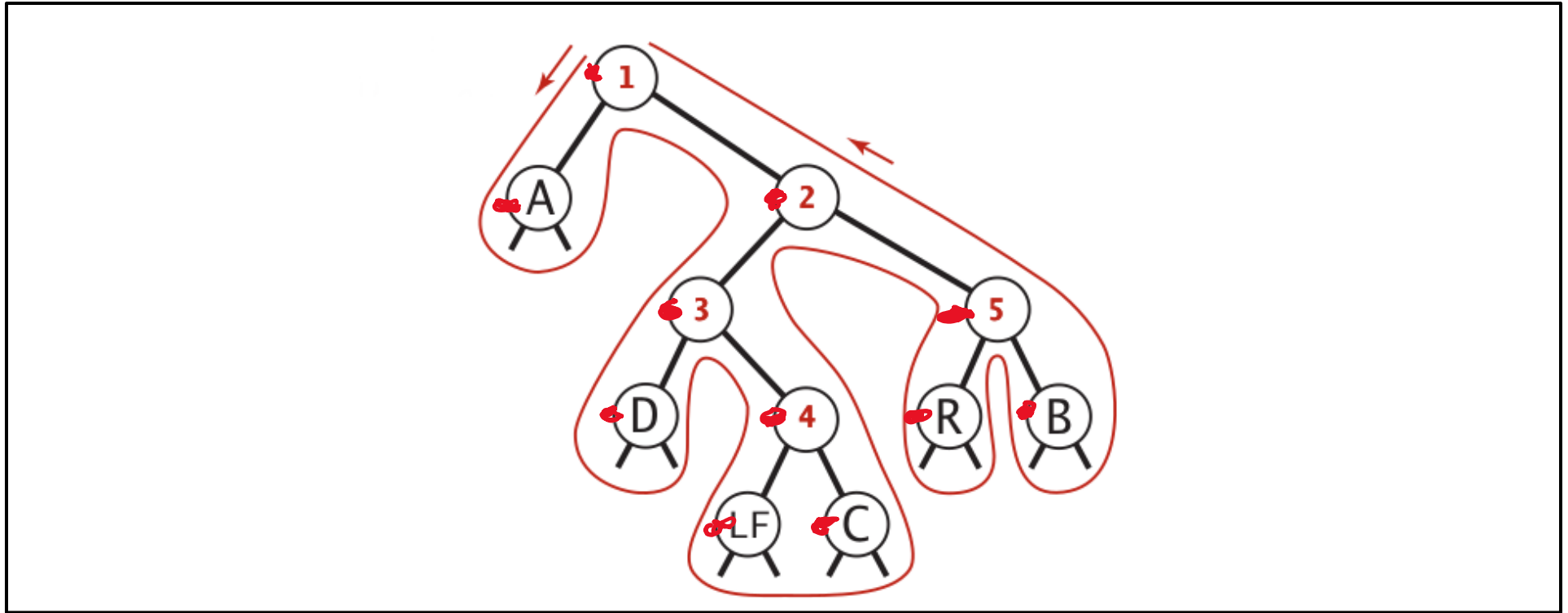
Preorder traversal



lea
0
↑
1

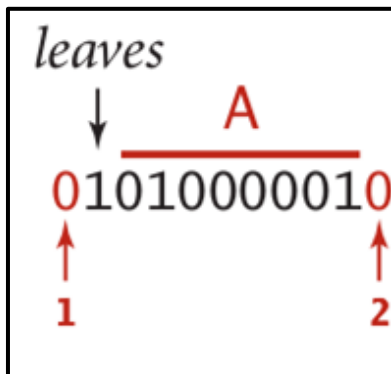
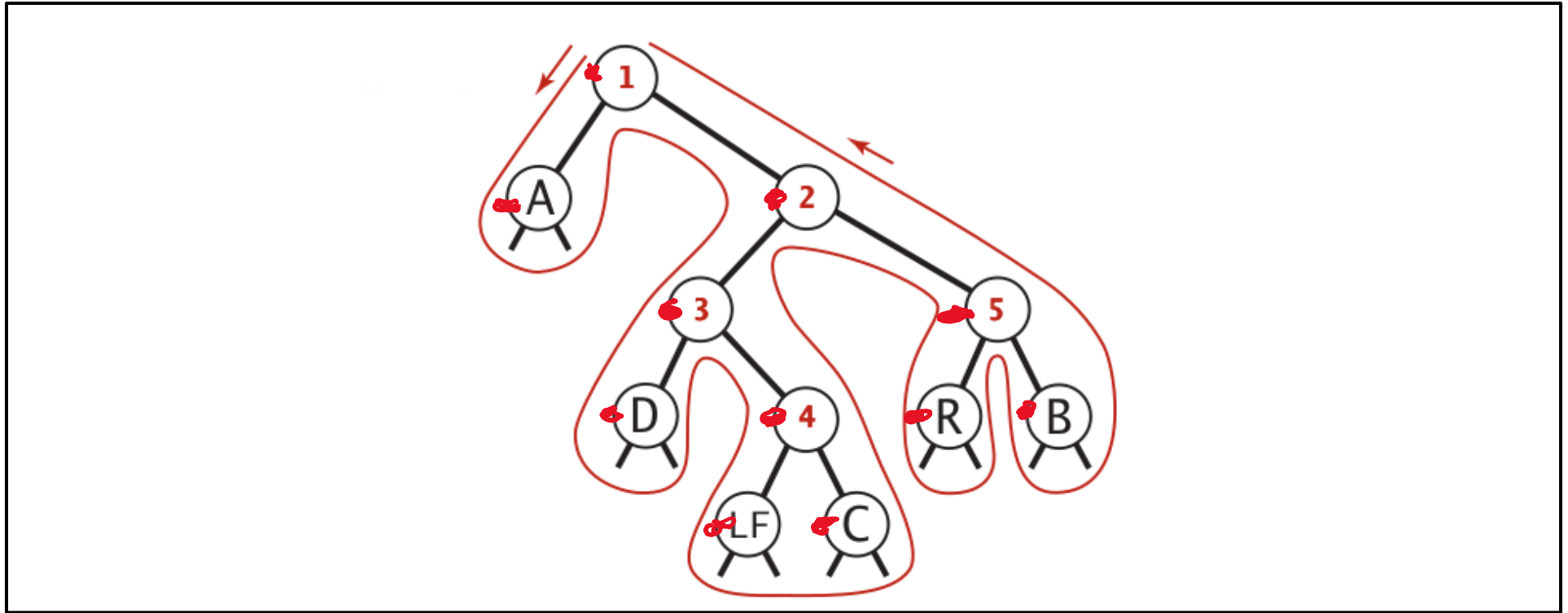
Representing tries as bitstrings

Preorder traversal



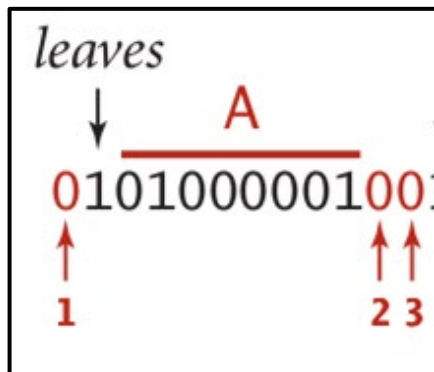
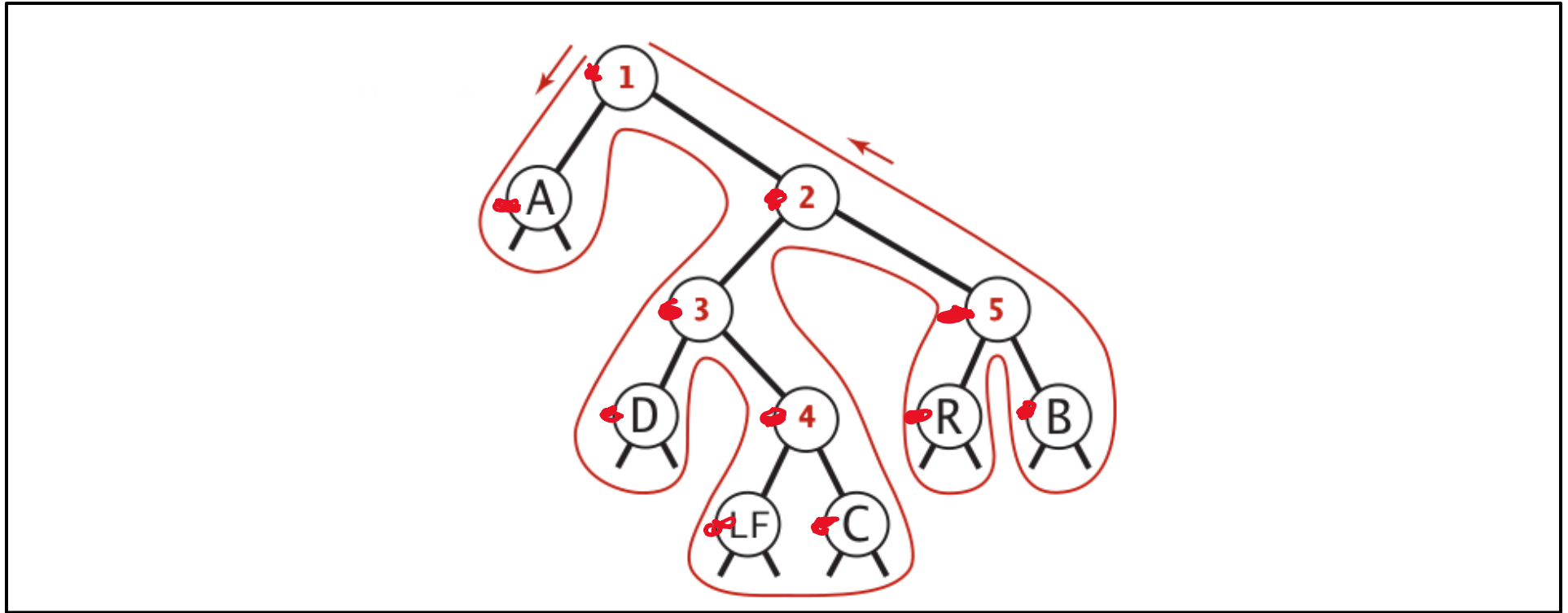
Representing tries as bitstrings

Preorder traversal



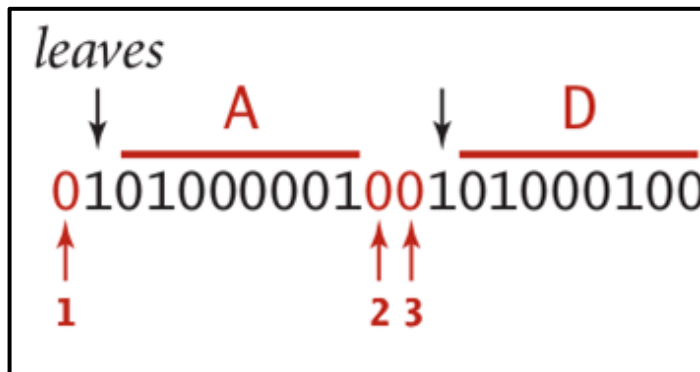
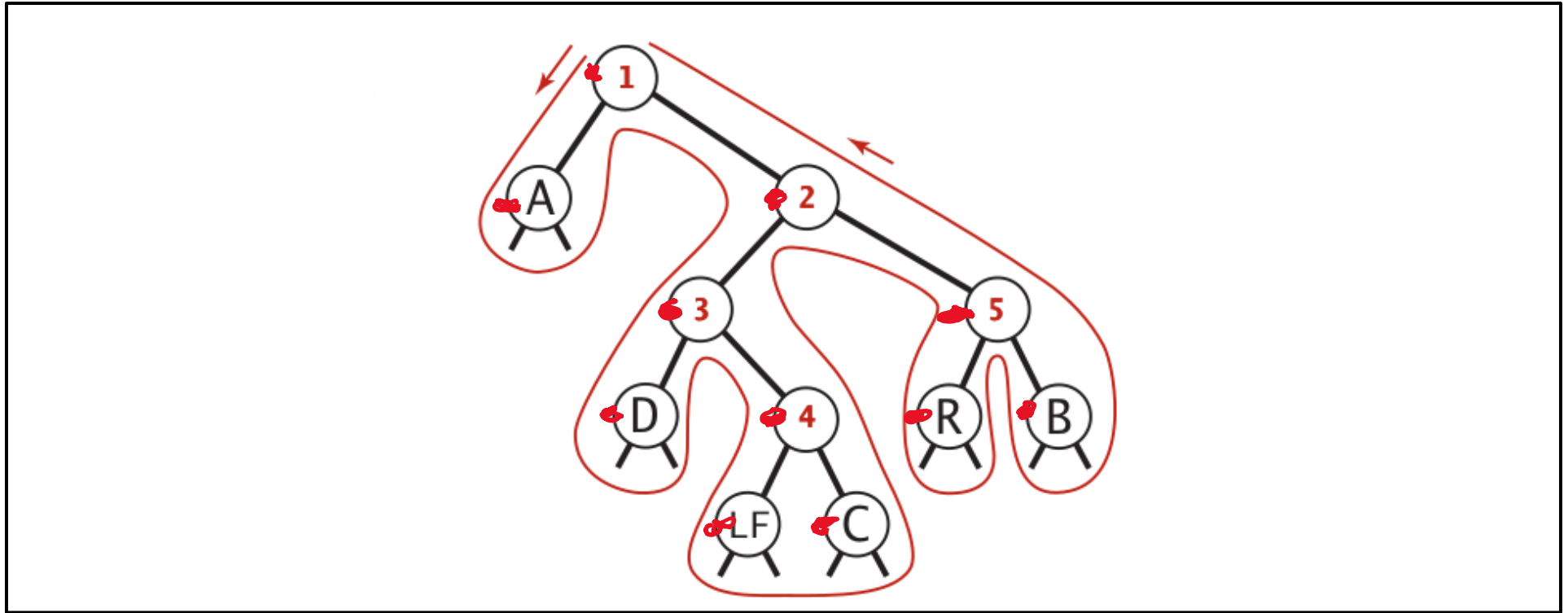
Representing tries as bitstrings

Preorder traversal



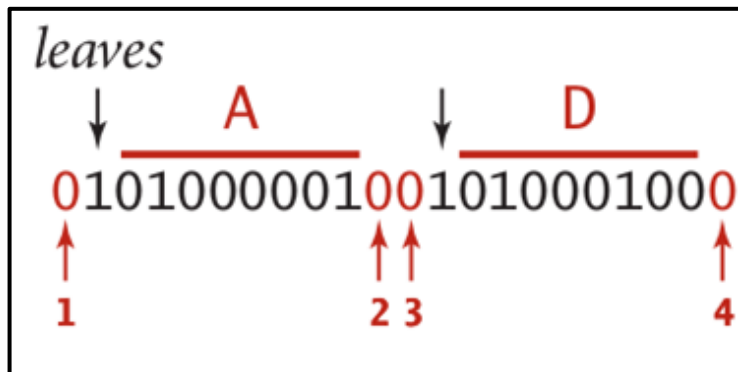
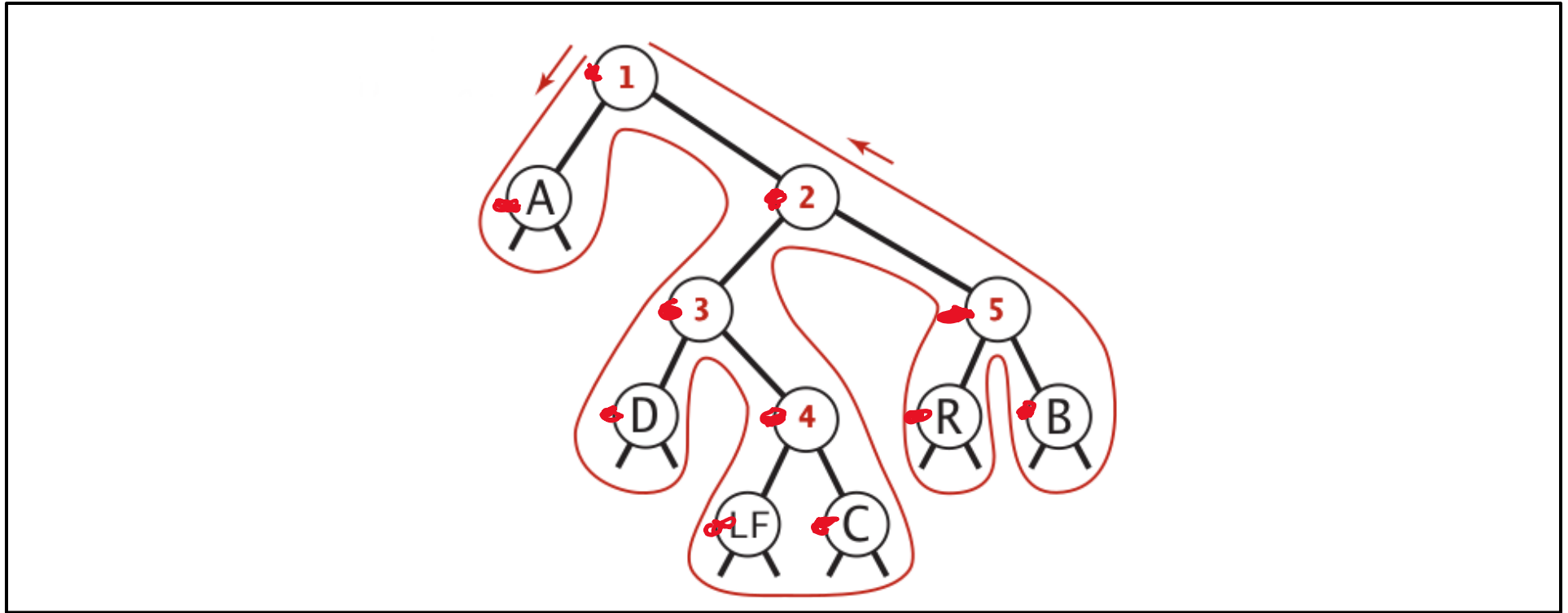
Representing tries as bitstrings

Preorder traversal



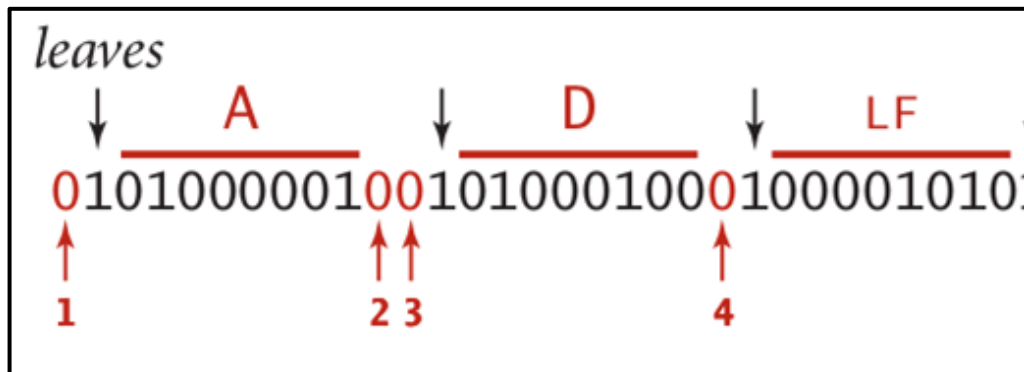
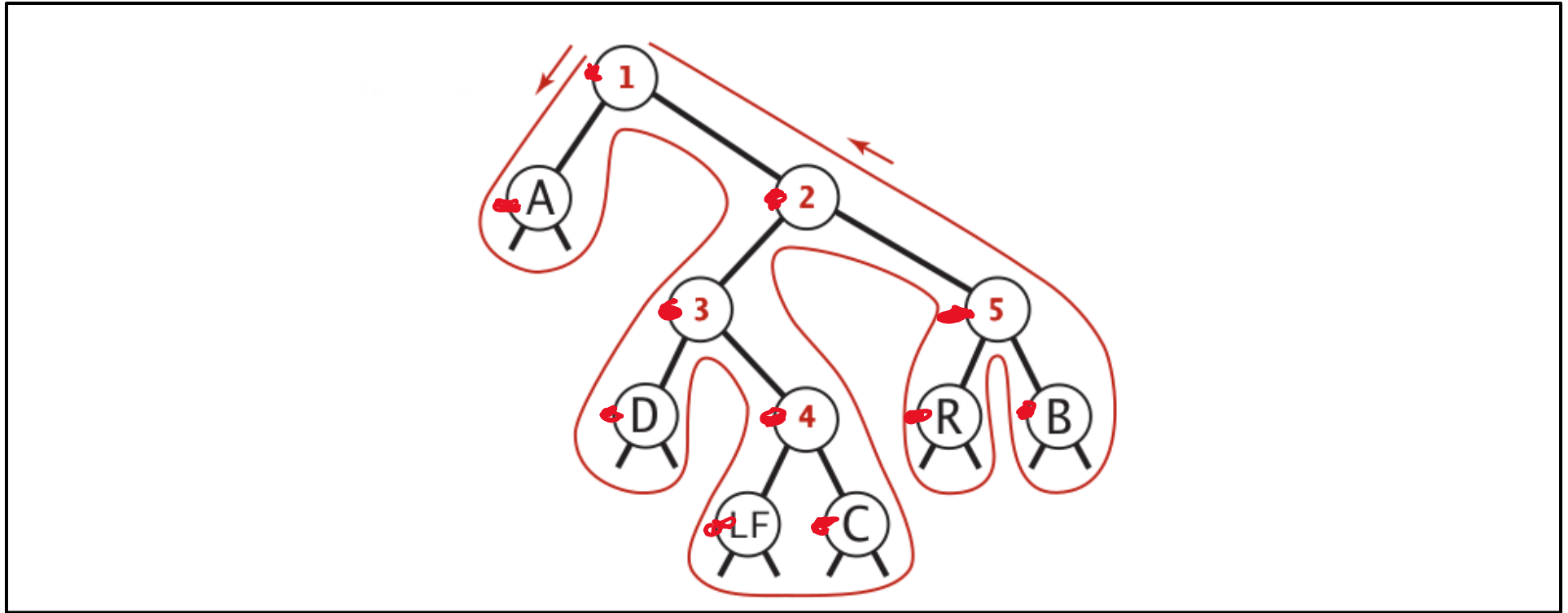
Representing tries as bitstrings

Preorder traversal



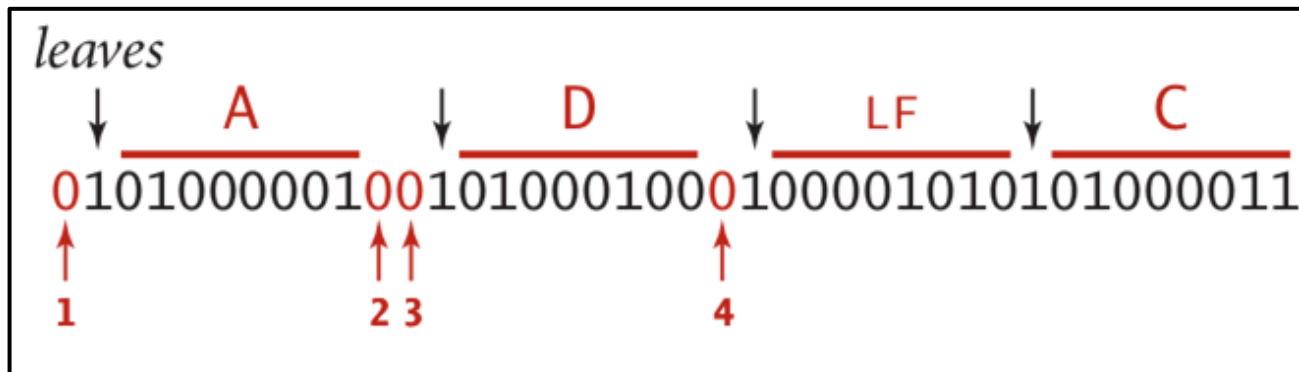
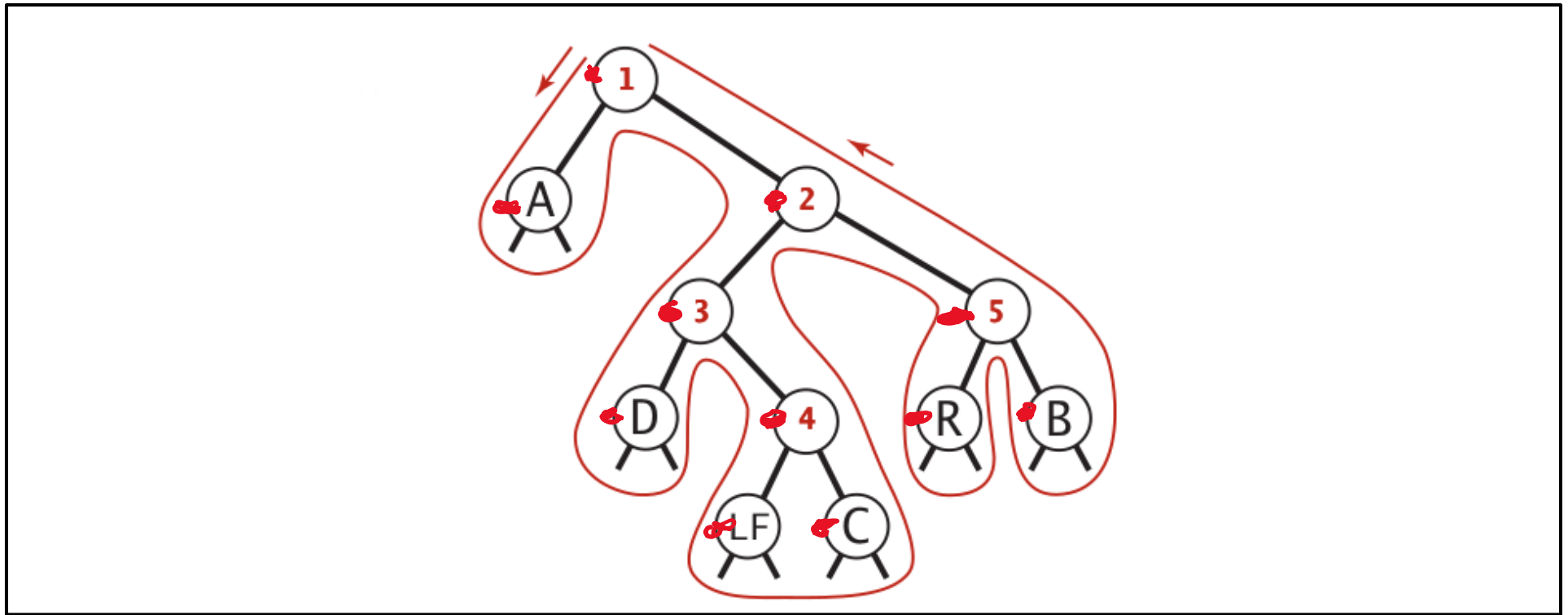
Representing tries as bitstrings

Preorder traversal



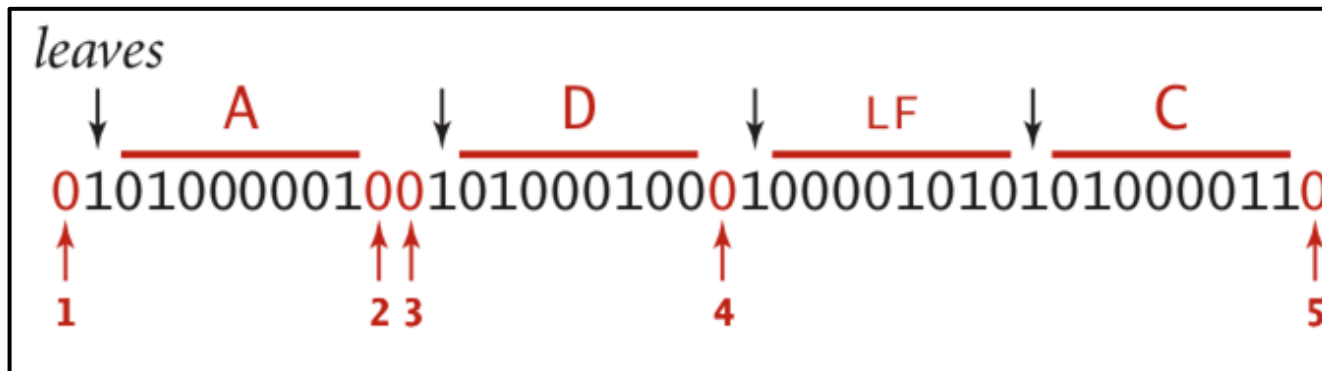
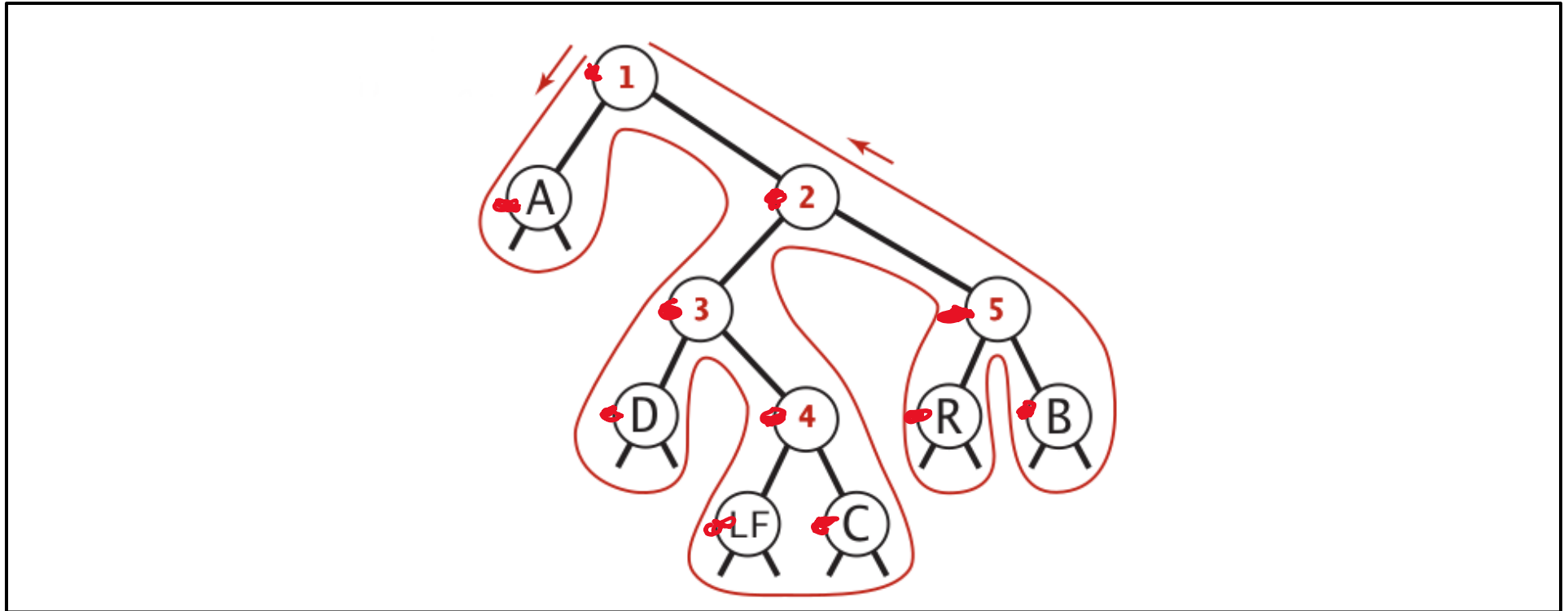
Representing tries as bitstrings

Preorder traversal



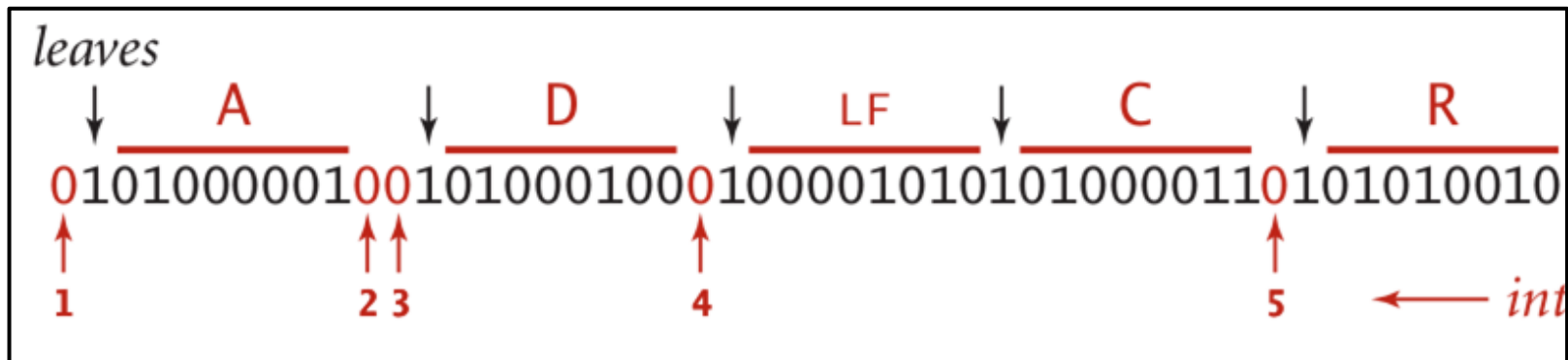
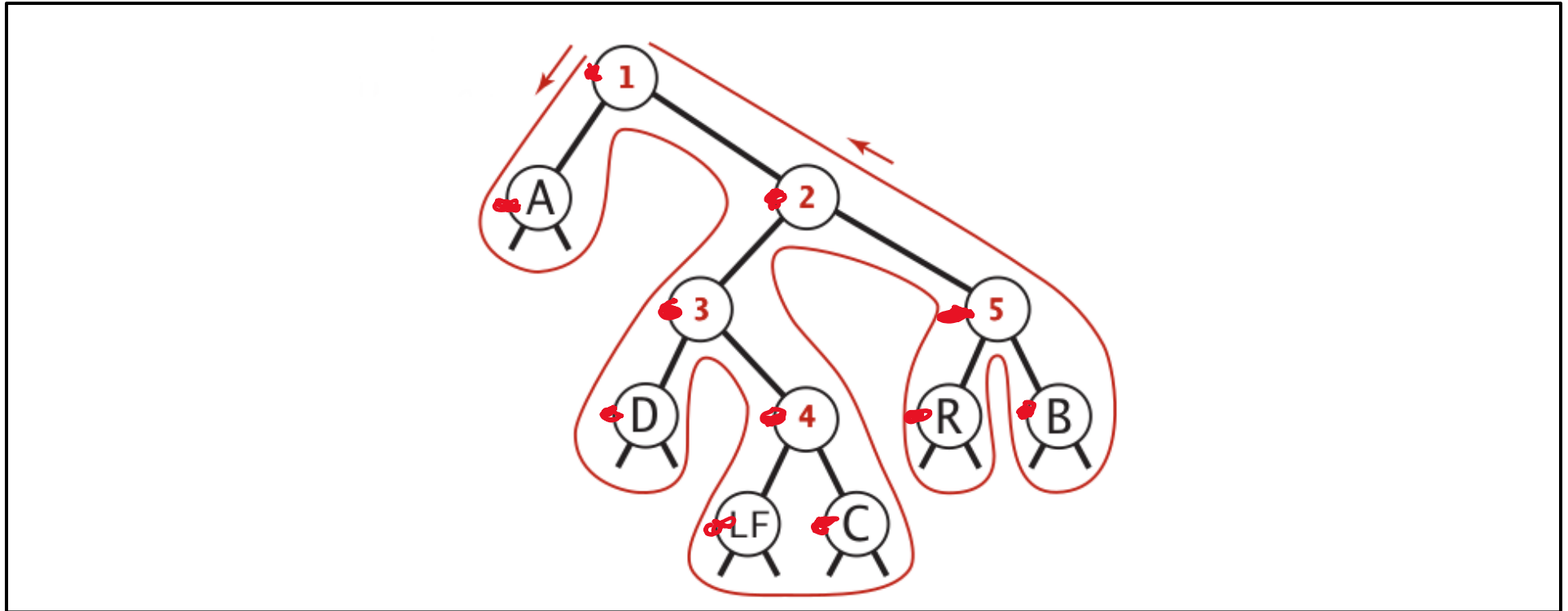
Representing tries as bitstrings

Preorder traversal



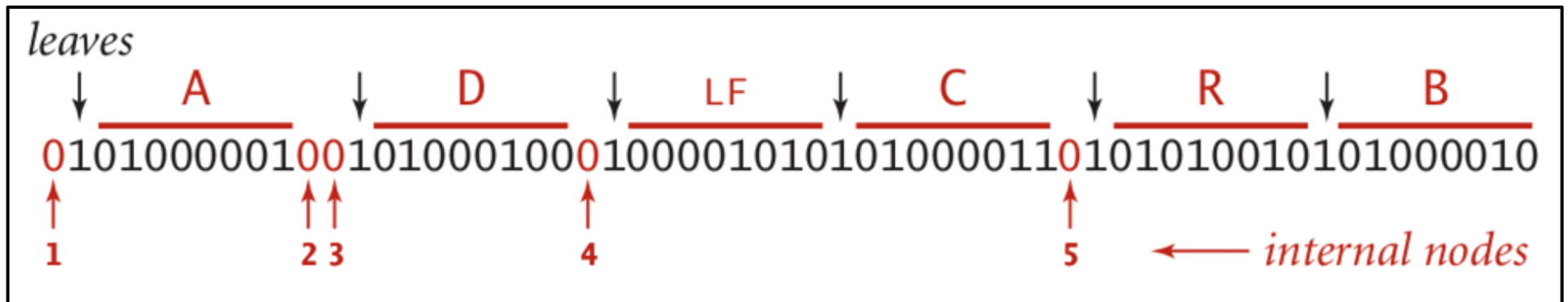
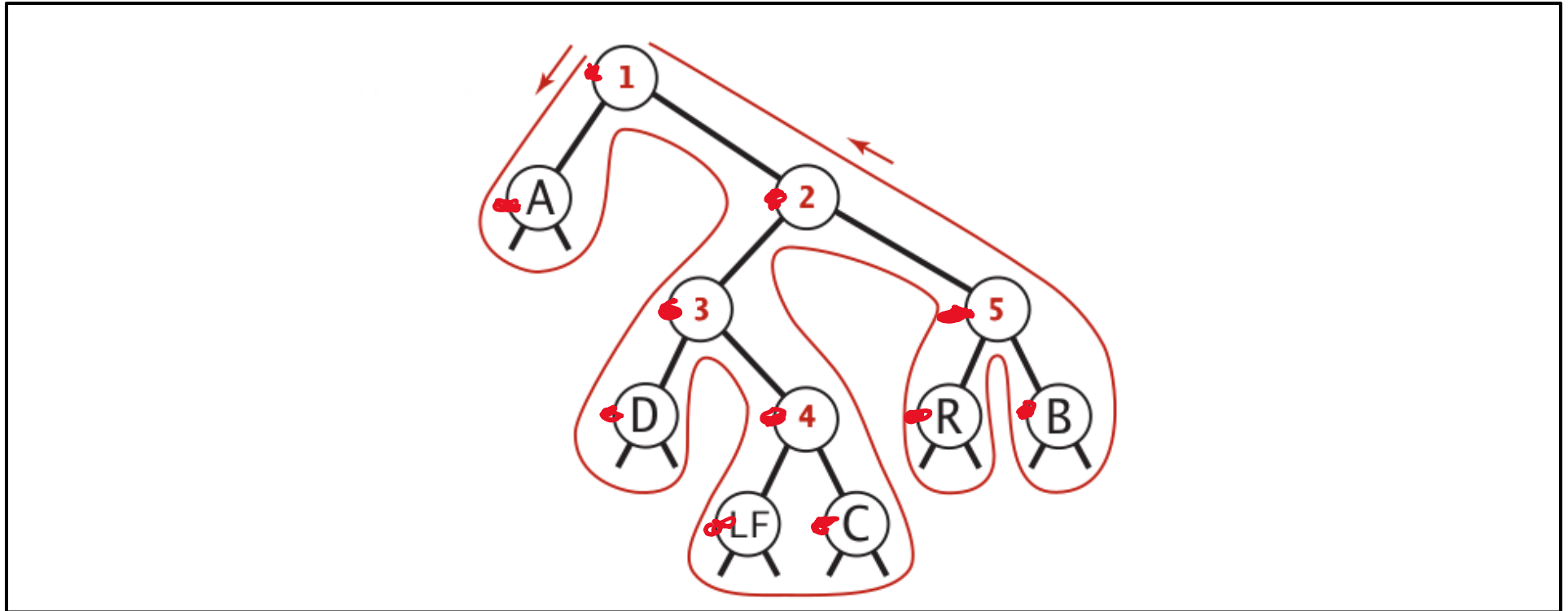
Representing tries as bitstrings

Preorder traversal



Representing tries as bitstrings

Preorder traversal



Muddiest Points

- **Q: After writing out the Trie as a bitstring, how can you tell where each character starts?**
- First, we are writing bits not characters!
- Each bit represents a node
 - except for leaves, their 1 bit is followed by 8 bits (ASCII code of char inside)
 - but that's also fixed length

Muddiest Points

- **Q: how compression and tries combine. Are you supposed to recreate the original trie while decompressing**
- The compressed file contains:
 - the trie representation in bits
 - the number of characters in the original file
 - the Huffman encoding of the file characters
- The original trie is reconstructed from the trie representation in the compressed file

Muddiest Points

- **Q: Questions are not getting responded to on Piazza in a timely manner. Students had asked several questions pertaining to items on Homework 4 last week but the questions were not addressed until AFTER Homework 4 was already graded. The questions I got marked wrong this week were the exact questions I had asked for clarification on but didn't get a response to (until after the homework was graded and returned...by that point it was too late).**
- I will hold a Live QA Session on Piazza every Friday 4:30-5:30 pm

Muddiest Points

- **Q: the overall concepts of code blocks/code words, like how do they fit into everything?**
- The input file is divided into code blocks
- Each code block is replaced by a codeword
- For Huffman:
 - code blocks are single characters
 - codewords are variable-length bit strings
- For LZW:
 - code blocks are the longest-match strings (variable length)
 - codewords are fixed-length integers (e.g., 12 bits)
- For RLE:
 - code block are long strings with identical characters (variable length)
 - codewords are fixed-length integer followed by fixed-length ASCII of the character

Muddiest Points

- **Q: Just making sure I got the order of the compression framework: We go from file to code block to code word via compression then back to code block via expansion**
- Yep!

Muddiest Points

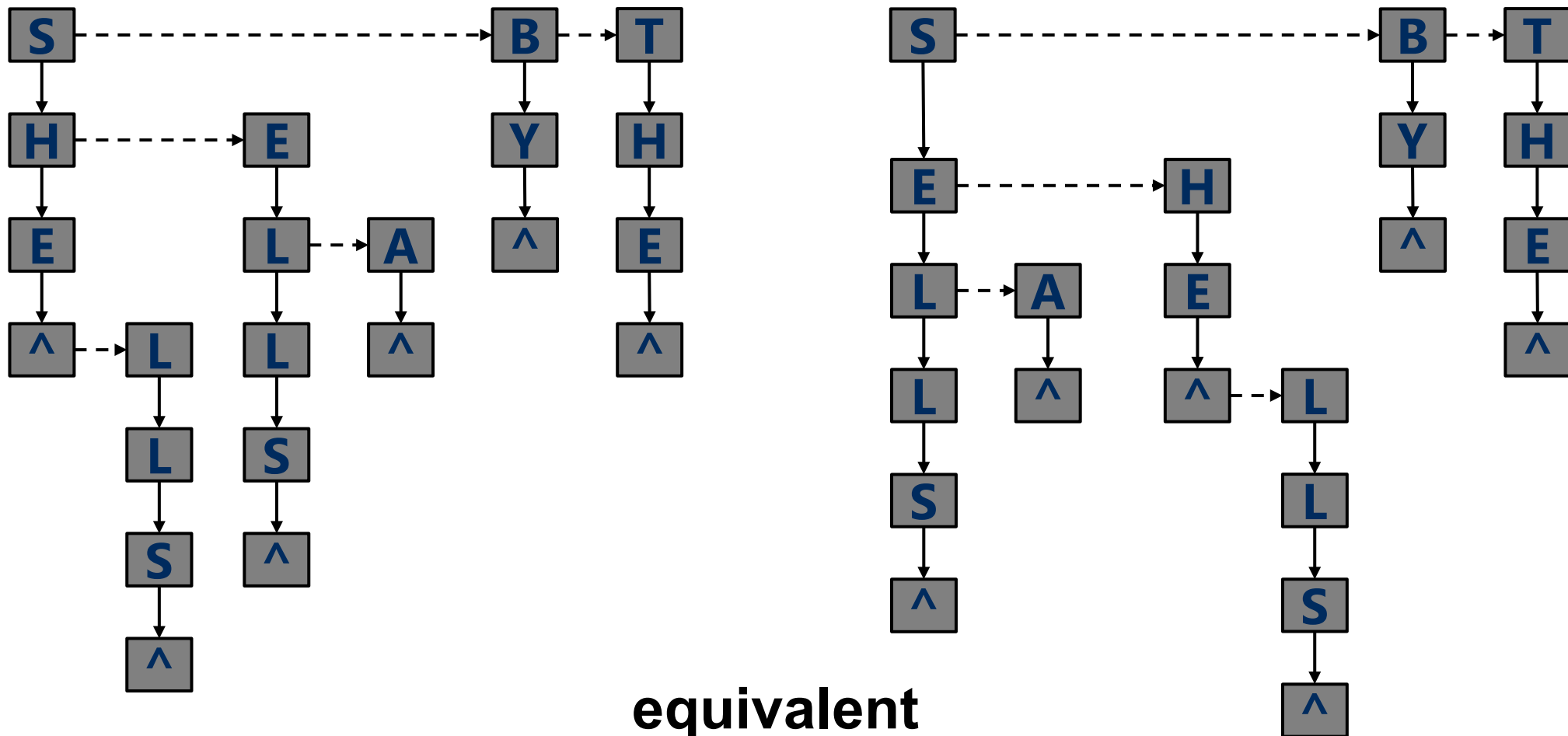
- **Q: Can we get partial credit on labs and projects?**
- We assign partial credit on projects only and only when autograder score $\leq 40\%$
 - partial credit has a ceiling of 60% of the autograder score
- Unless you think you lost points in the autograder because of an autograder error

Muddiest Points

- **Q:How are you able to create unique codes without leading to another prefix**
- Characters are leaves in the Huffman tree
- No two characters share the same path from root
- Codewords encode the root-to-leaf path
- No two codewords share a prefix

Muddiest Points

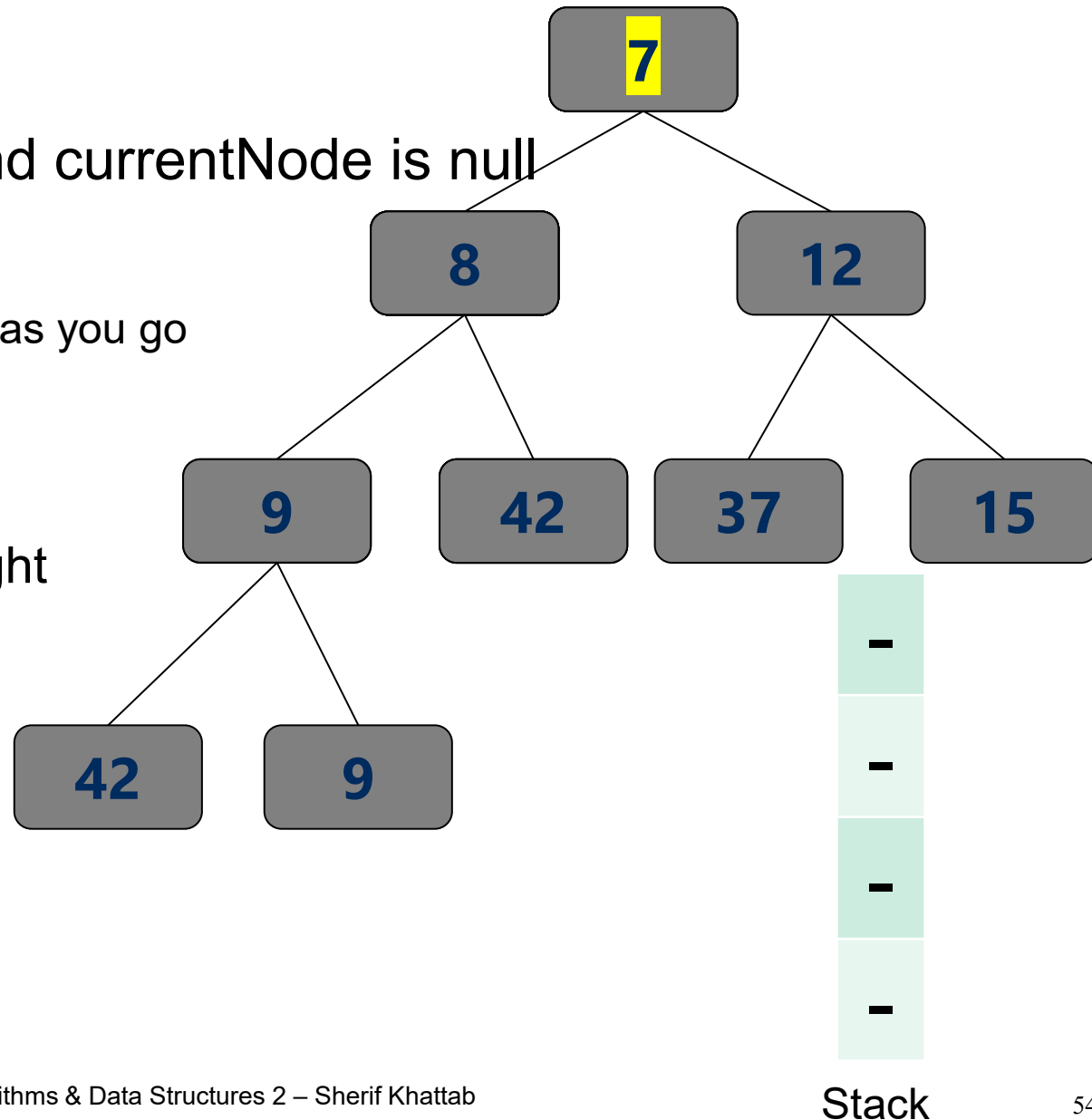
- **Q: In DLB tries, can you interchange a parent node's child with any one of the child's siblings? I think you can but you would have to change the sibling links, is that correct?**
- Correct!



Muddiest Points

- **Q: Review iterative inorder traversal**

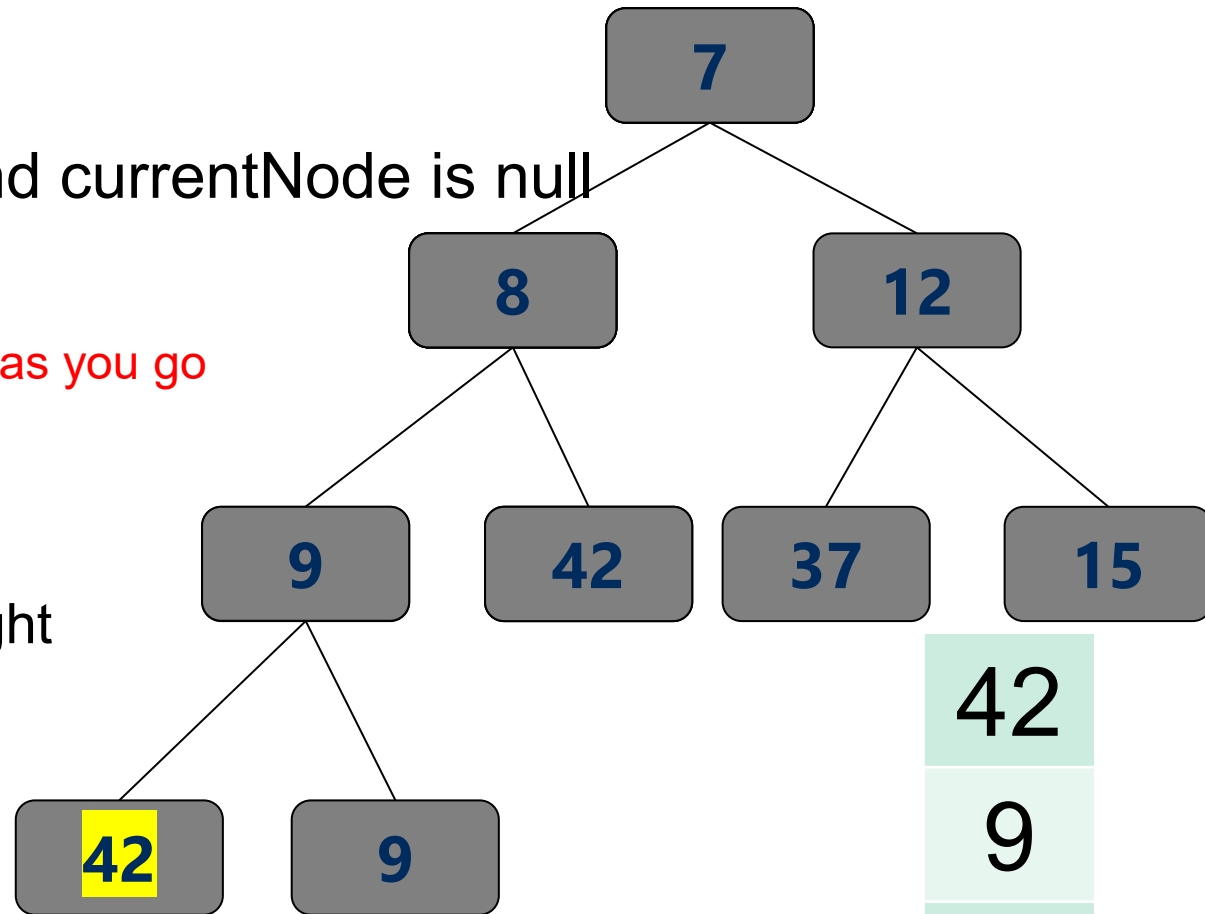
- `currentNode = root`
- repeat until stack is empty and `currentNode` is null
 - if `currentNode != null`
 - Move to leftmost node and push as you go
 - `currentNode = pop()`
 - visit `currentNode`
 - `currentNode = currentNode.right`



Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42

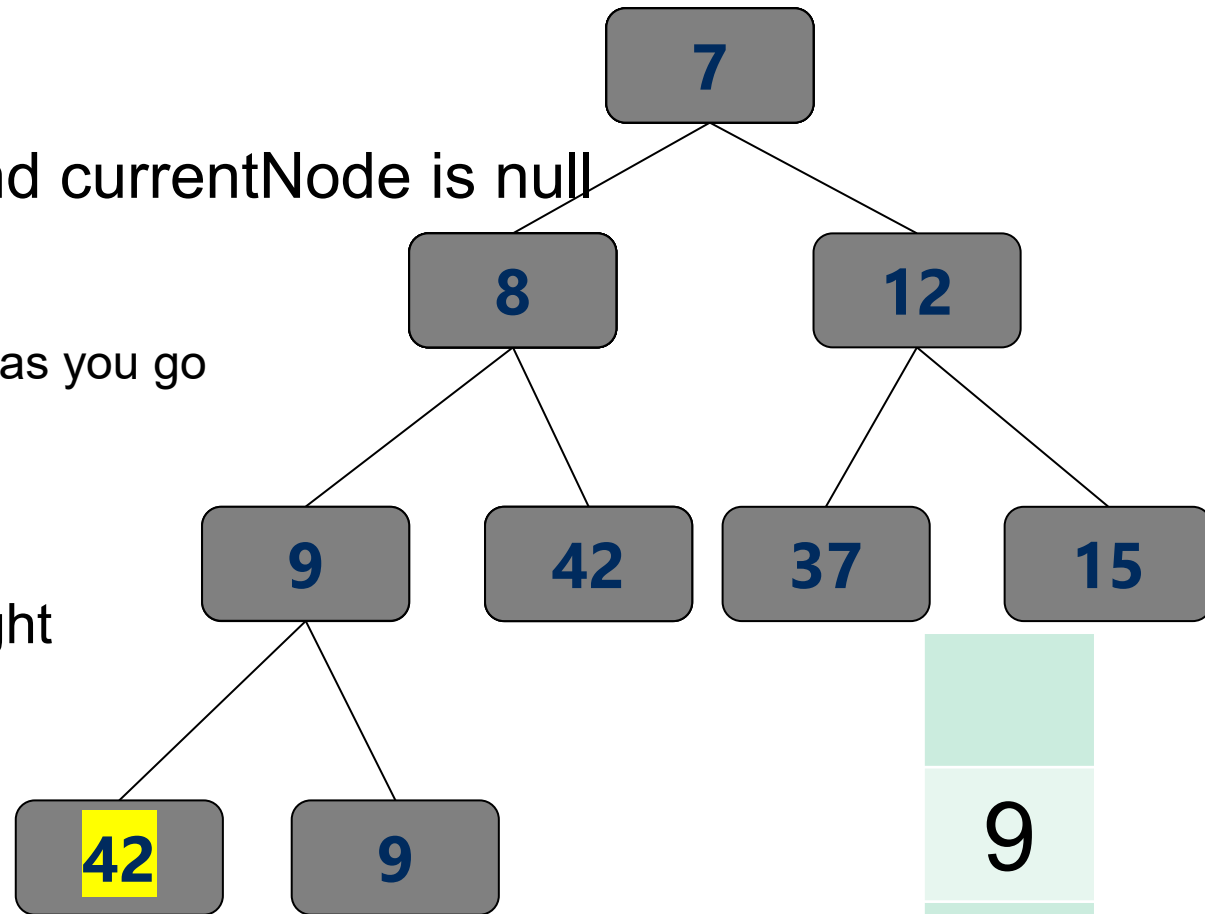
42
9
8
7

Stack

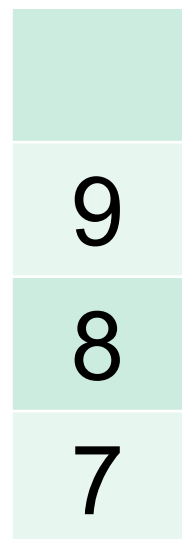
Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42

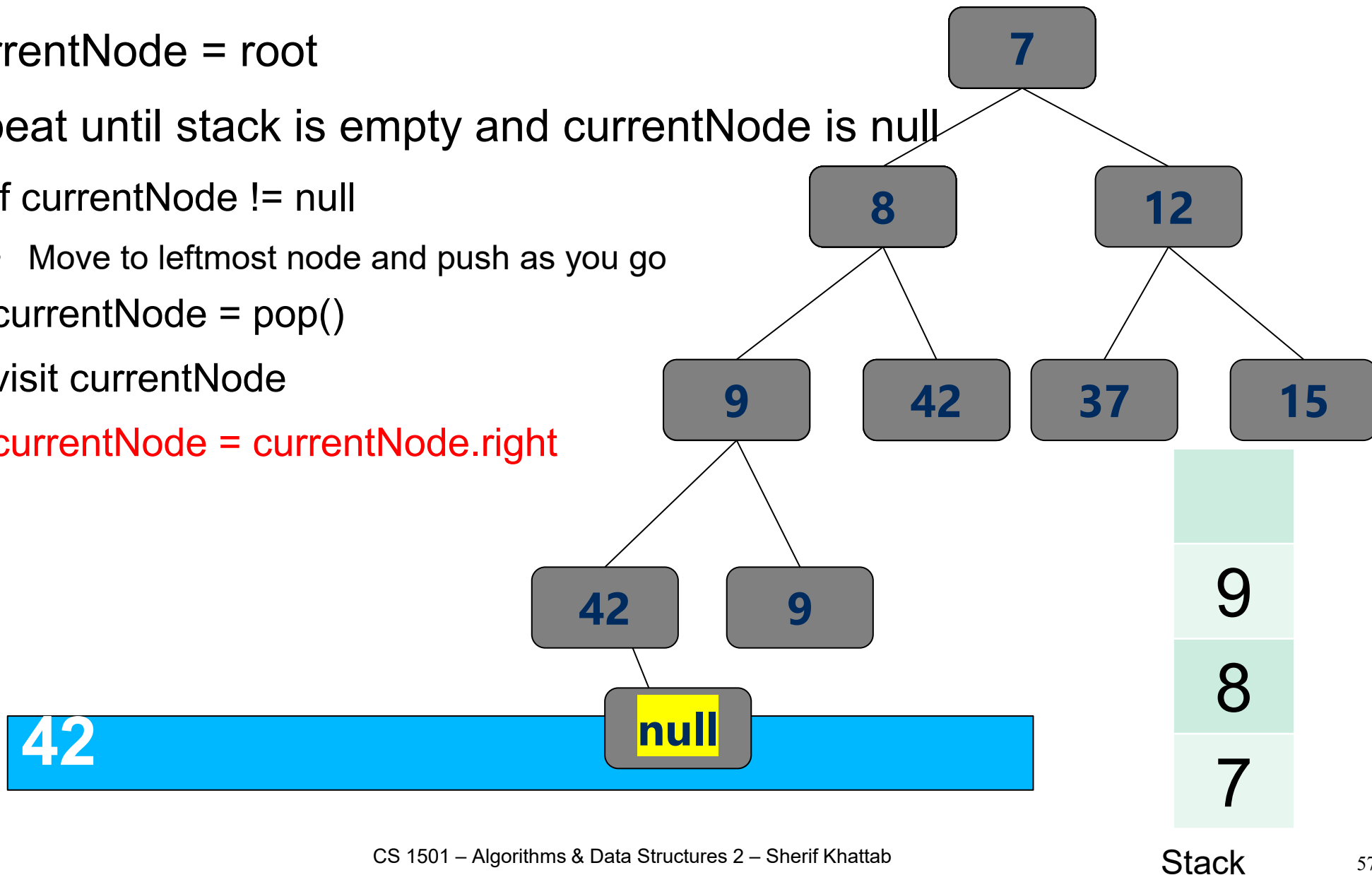


Stack

Muddiest Points

- Q: Review iterative inorder traversal

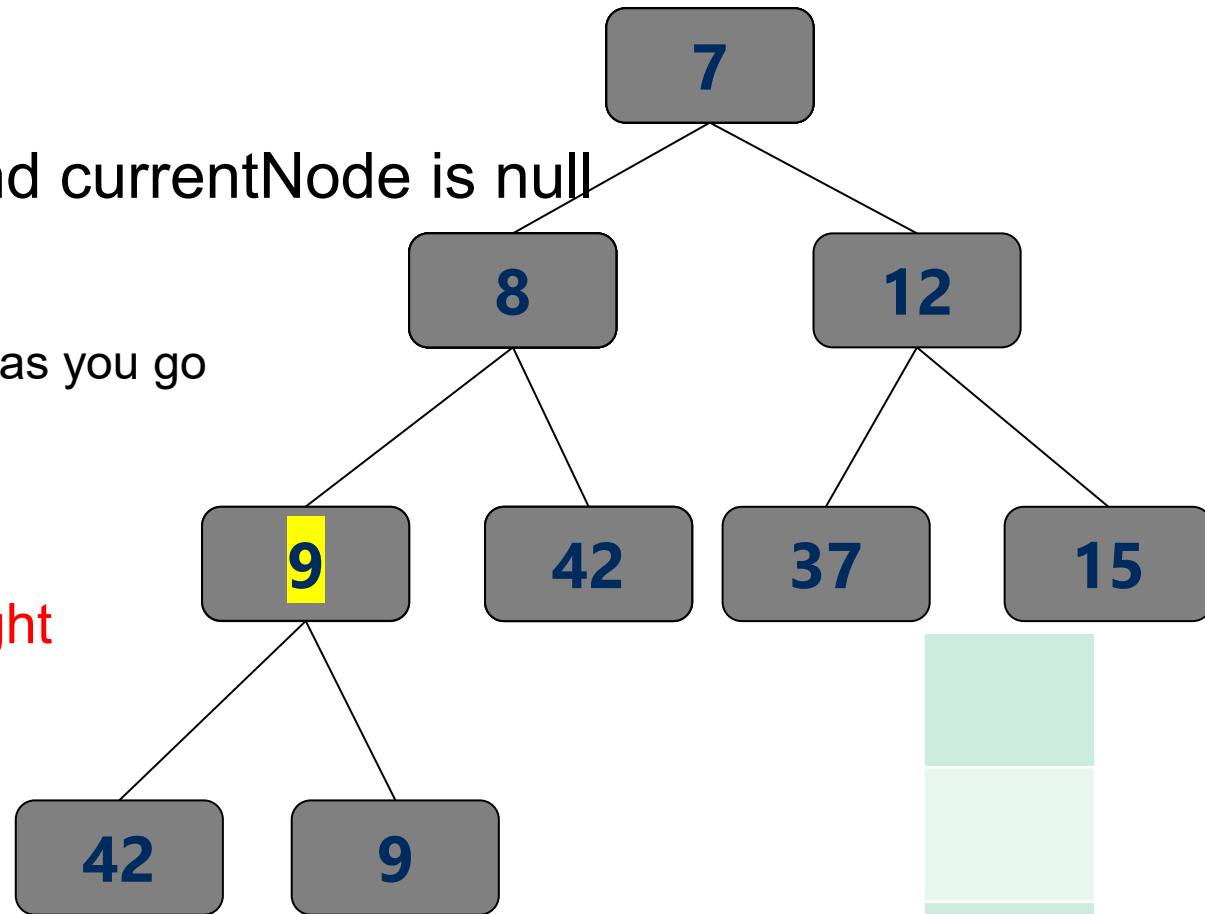
- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



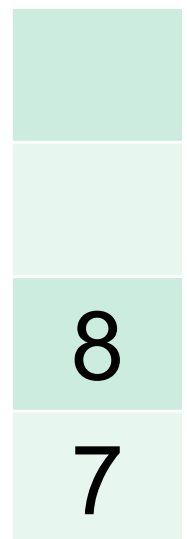
Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42, 9

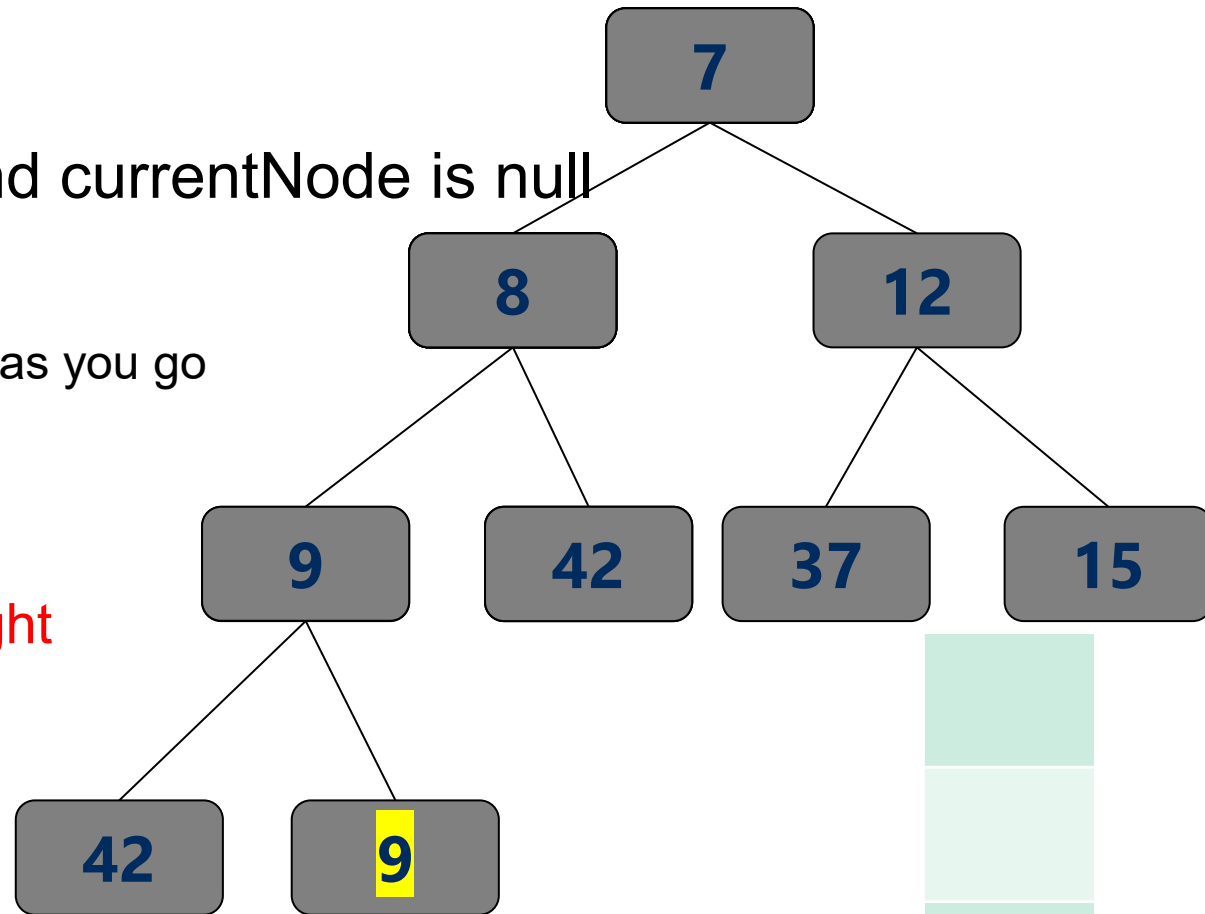


Stack

Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right

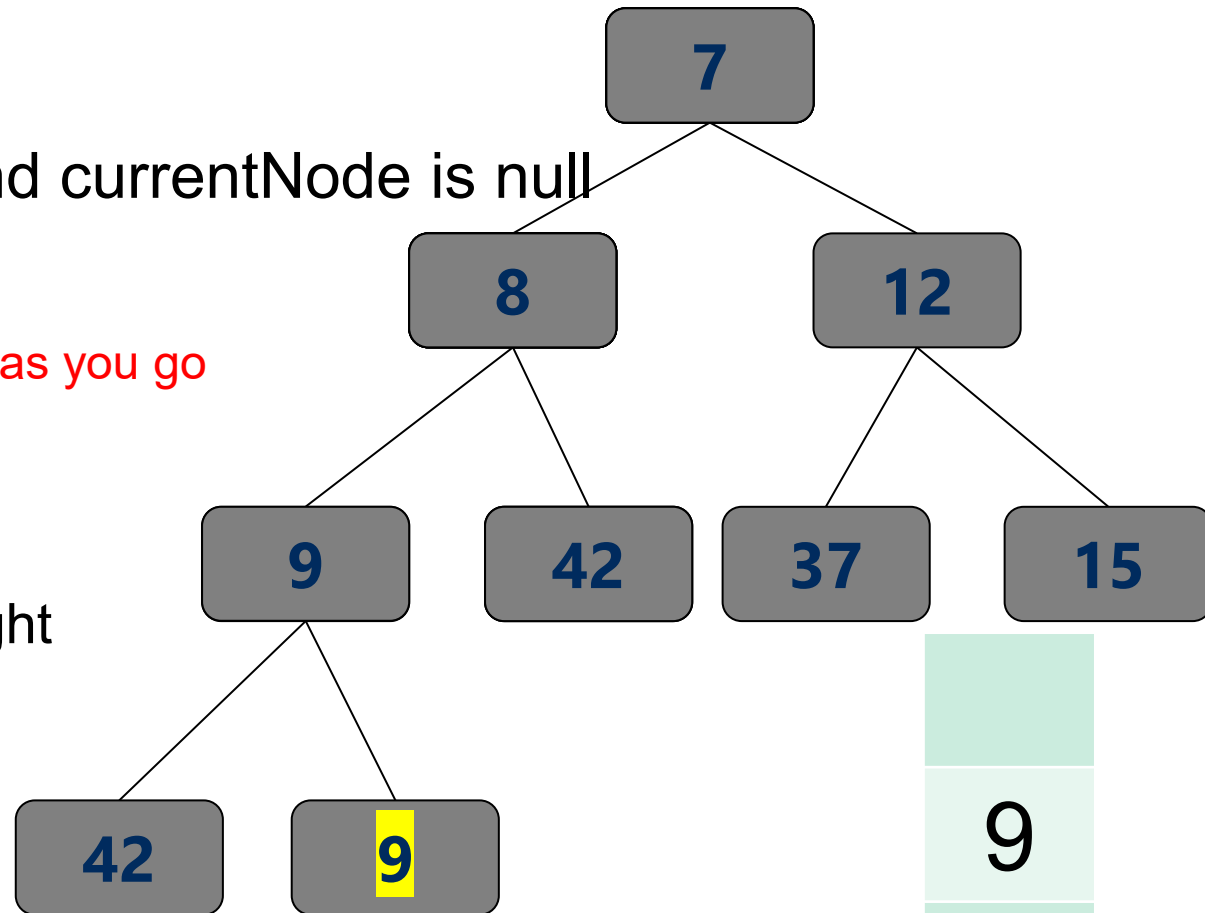


42, 9

Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right

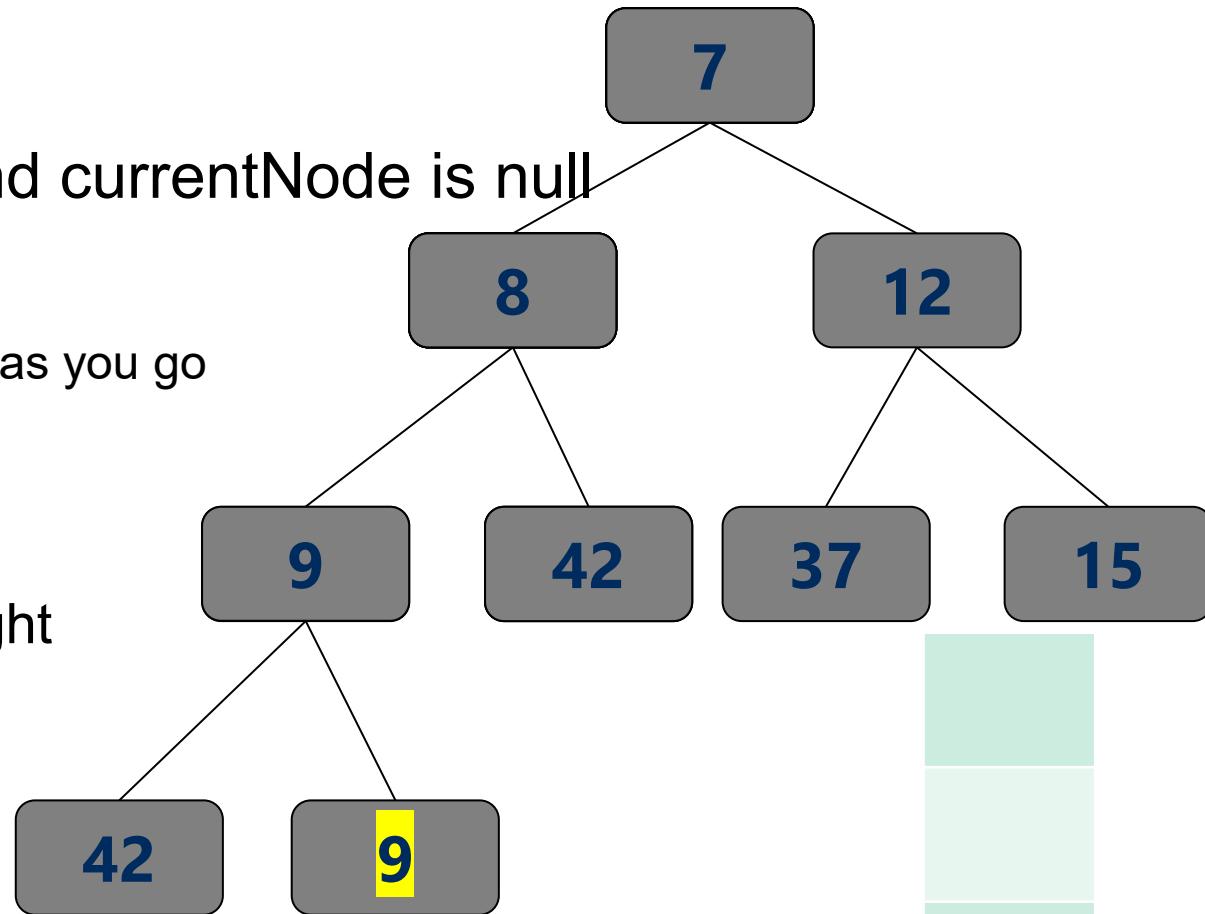


42, 9

Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right

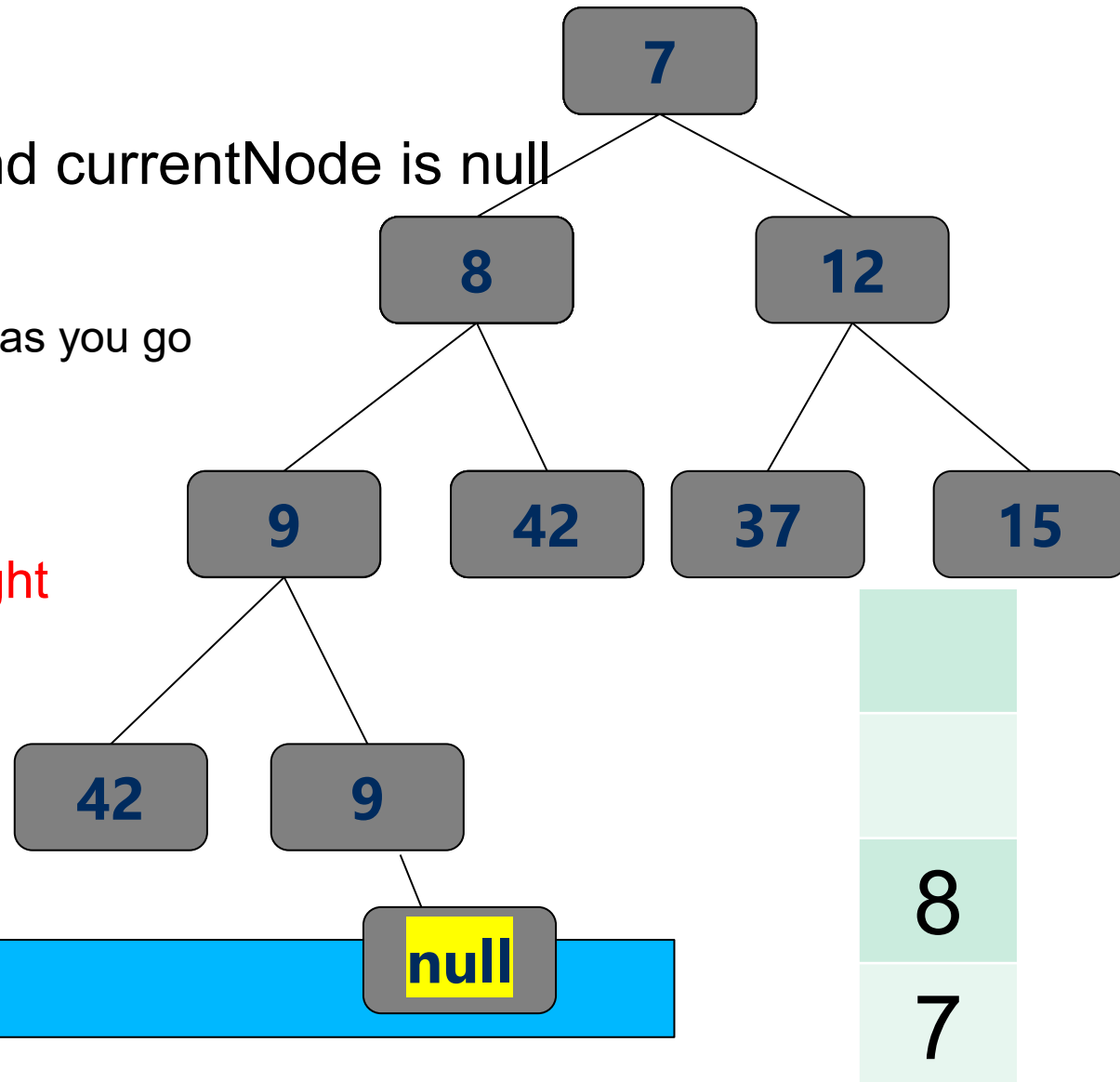


42, 9, 9

Muddiest Points

- Q: Review iterative inorder traversal

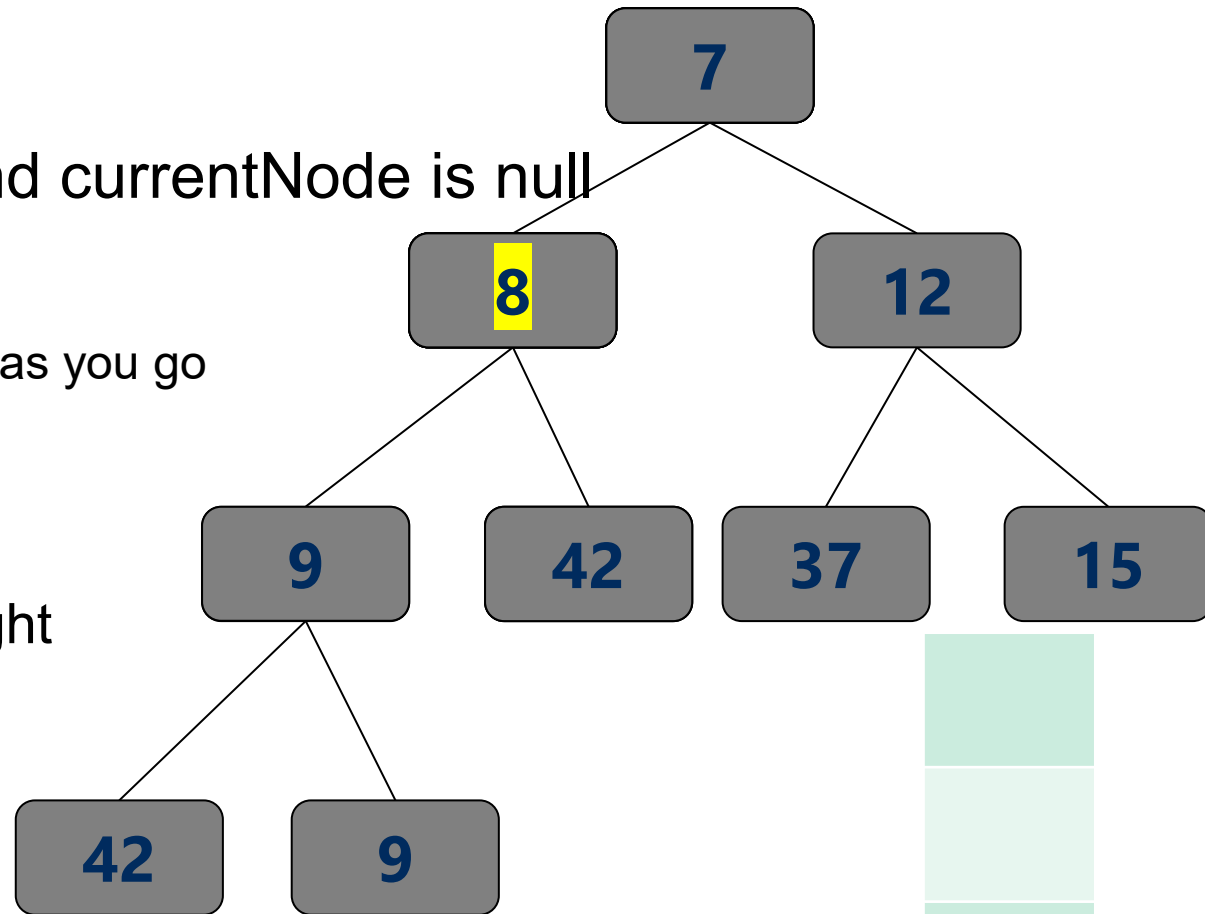
- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



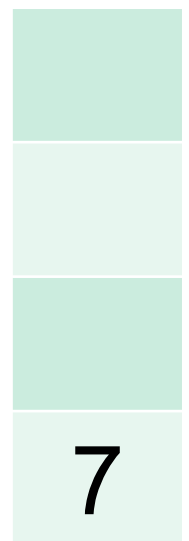
Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42, 9, 9, 8

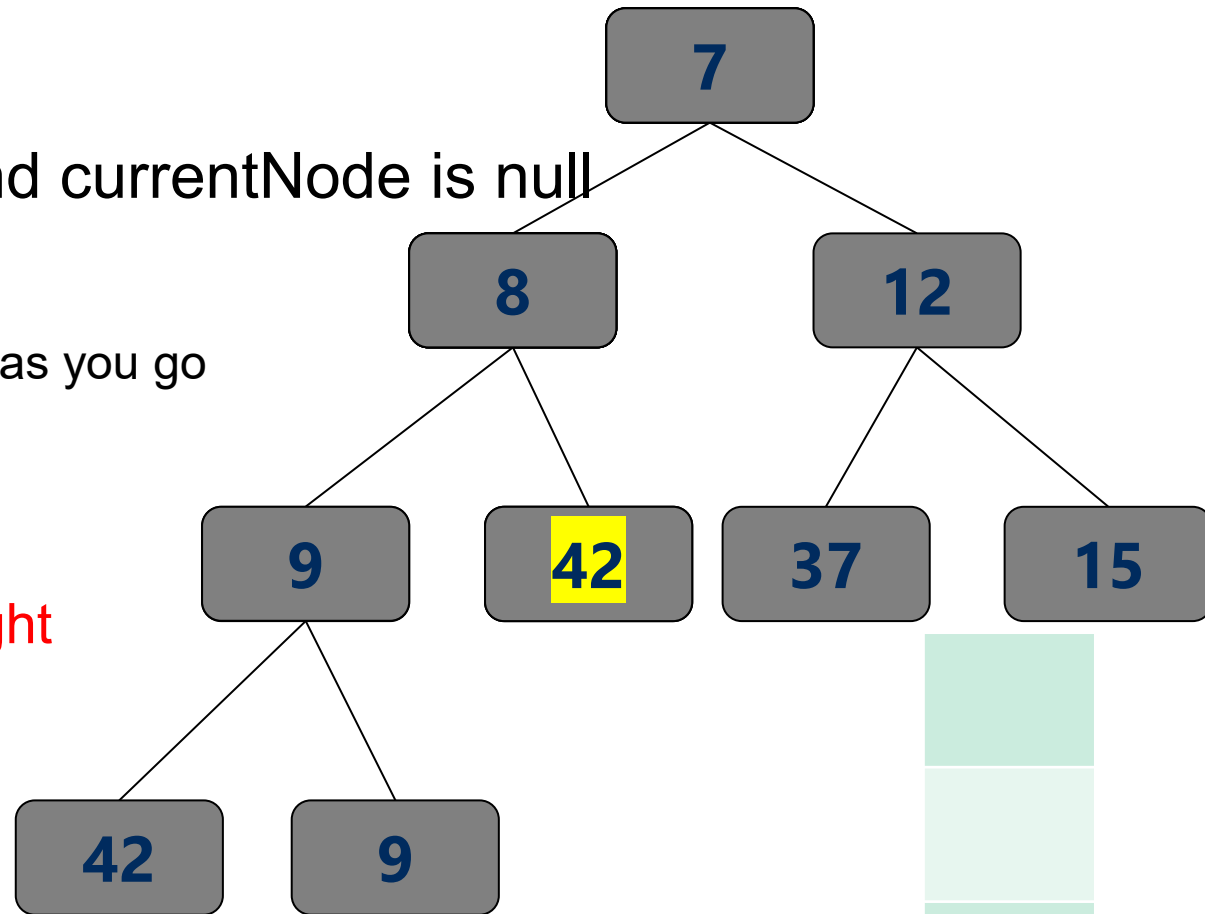


Stack

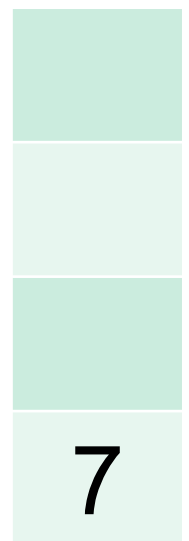
Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42, 9, 9, 8

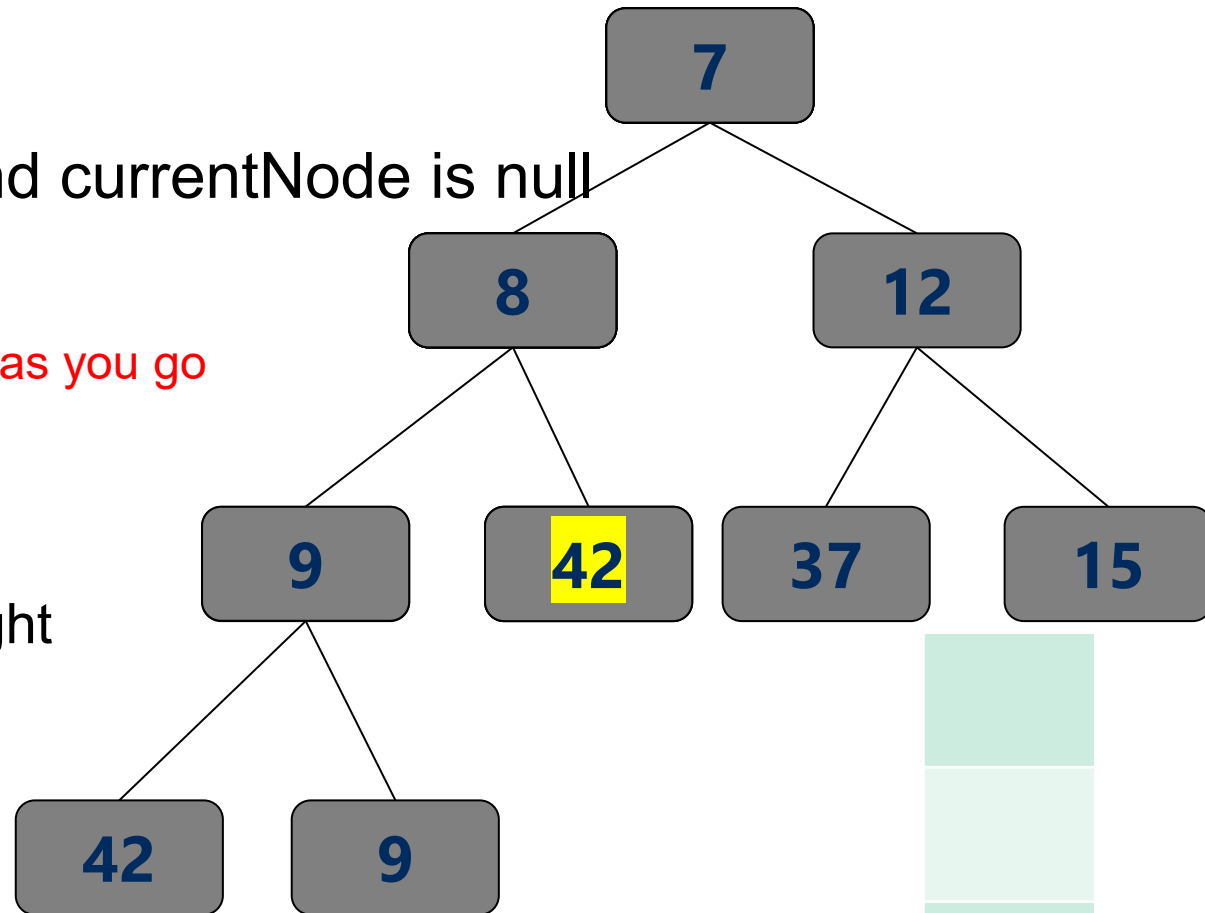


Stack

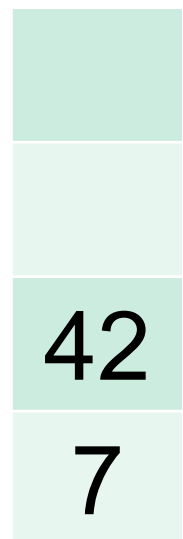
Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42, 9, 9, 8

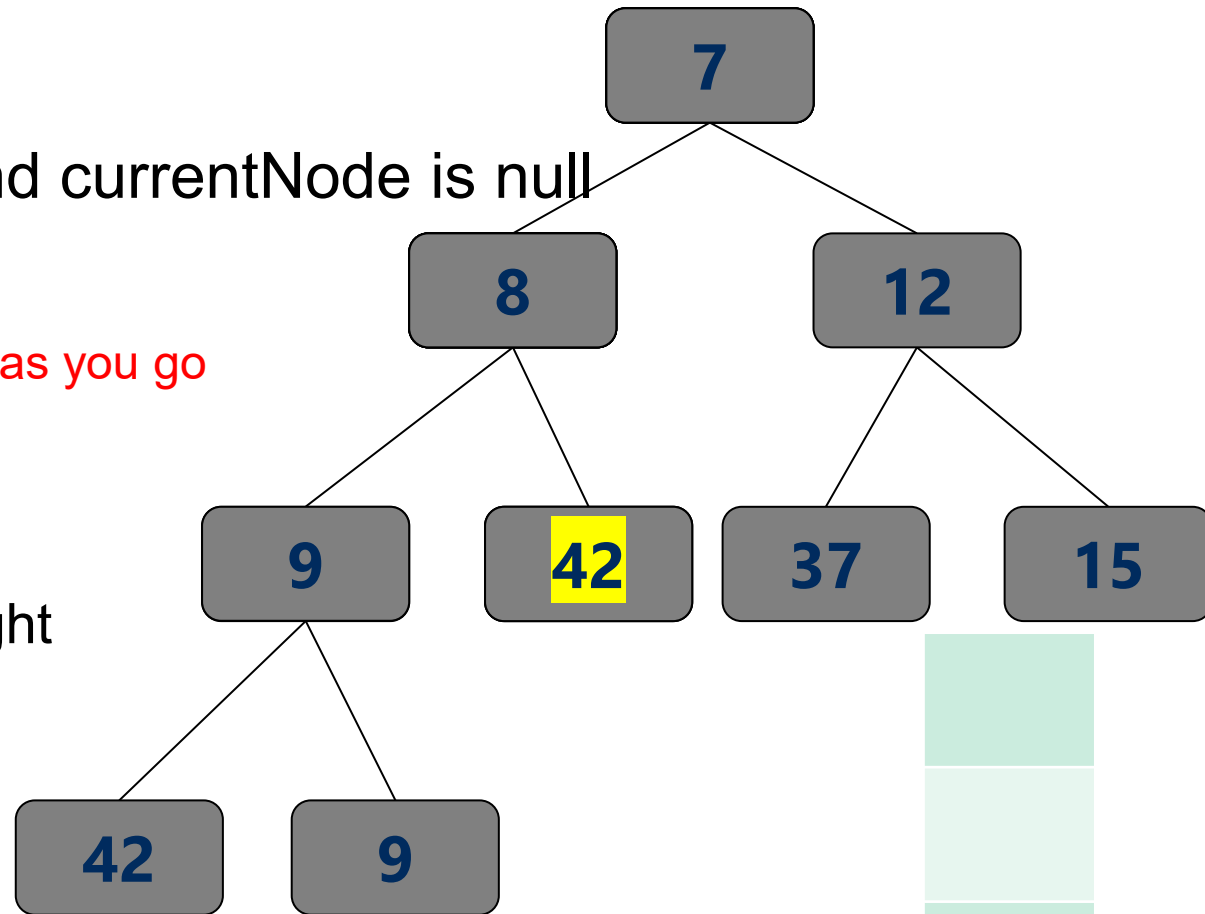


Stack

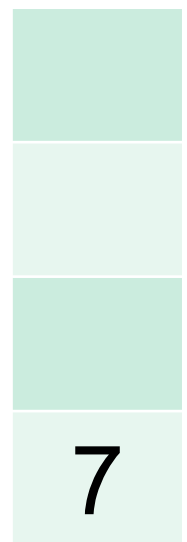
Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42, 9, 9, 8, 42

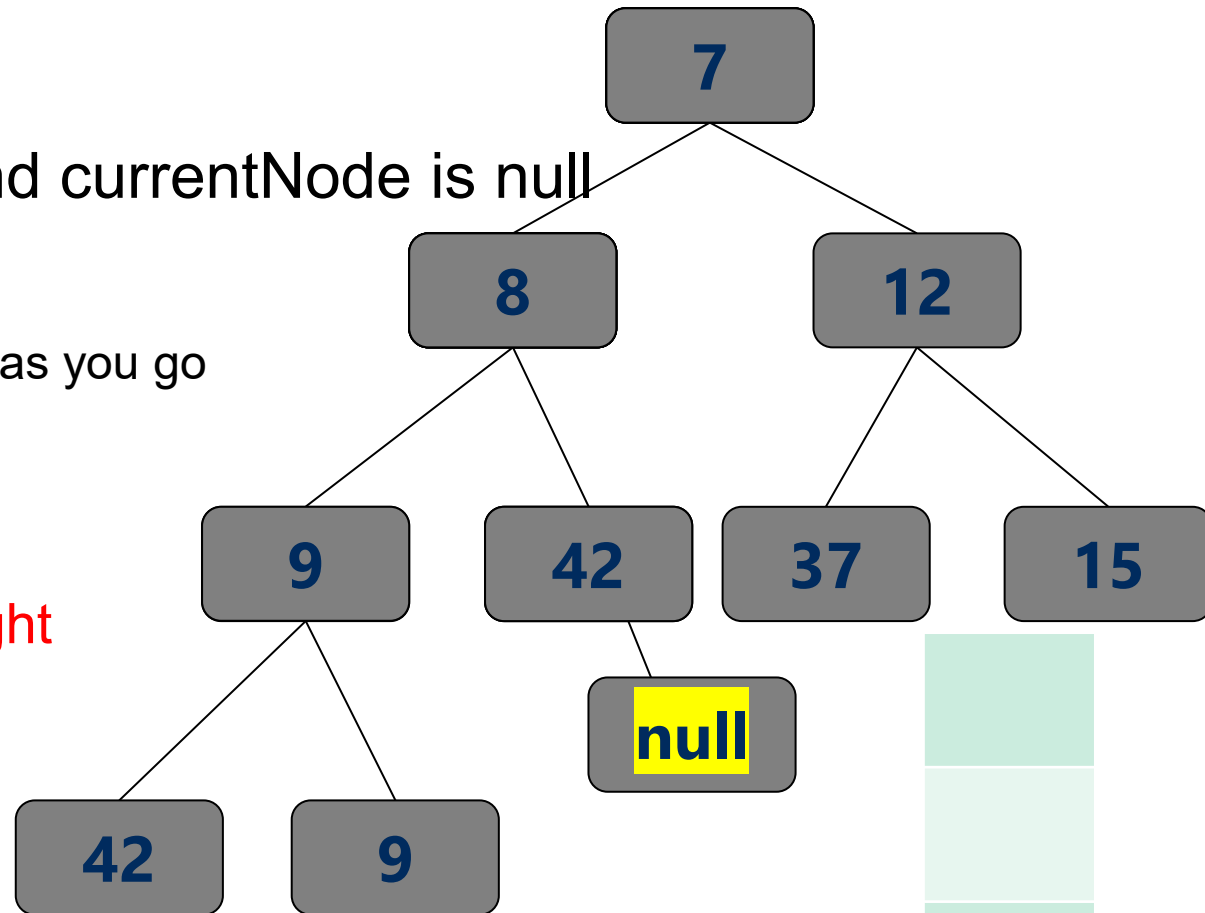


Stack

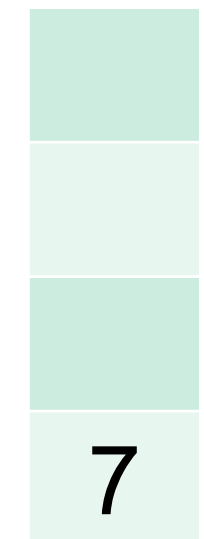
Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42, 9, 9, 8, 42

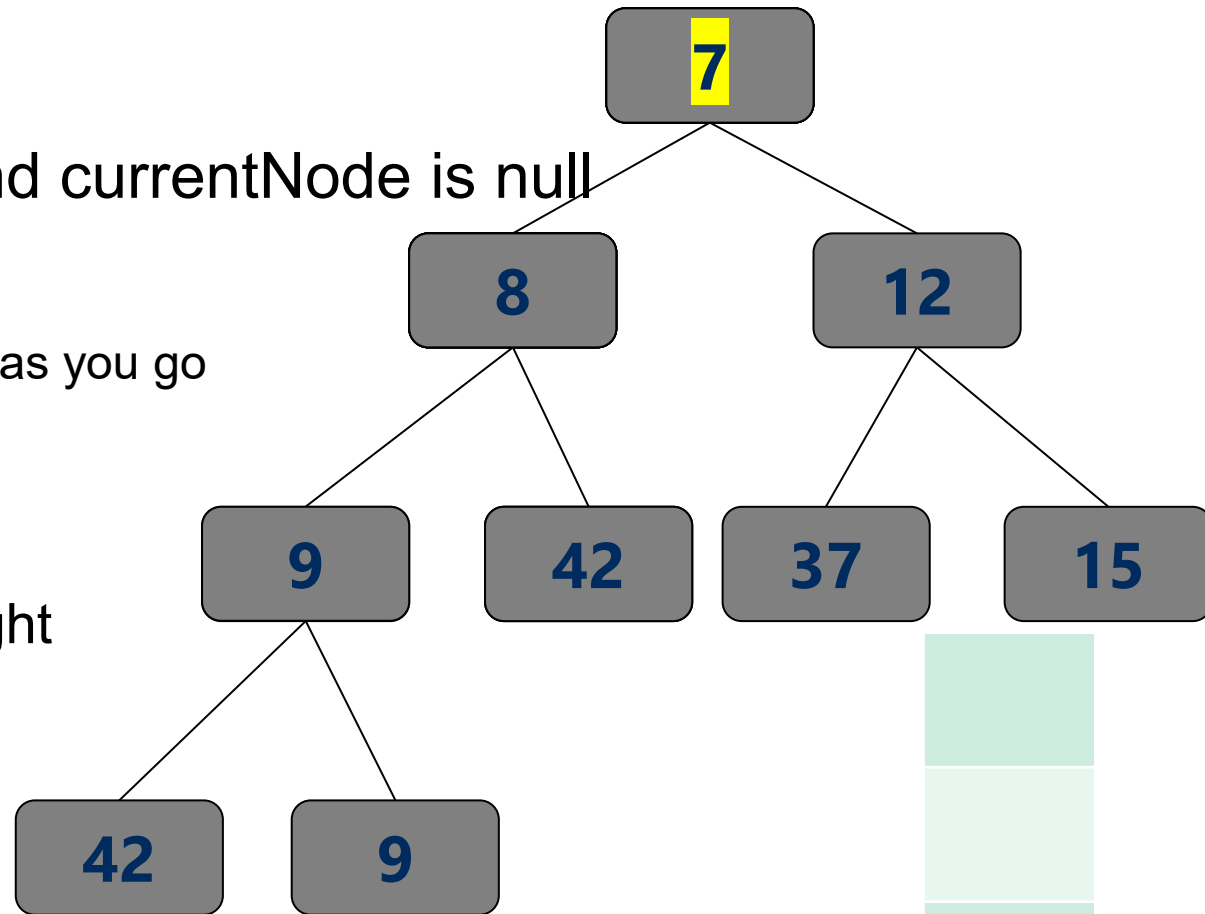


Stack

Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right

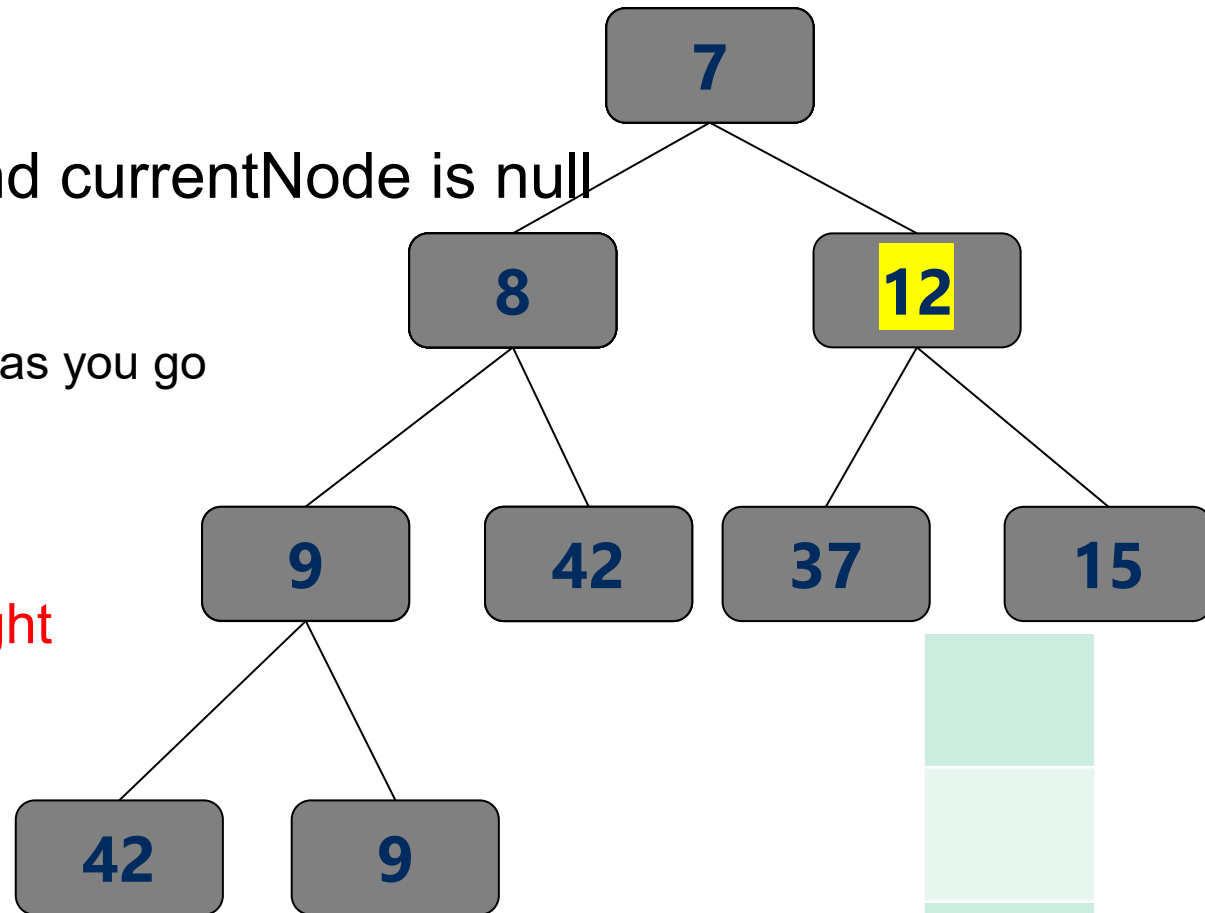


42, 9, 9, 8, 42, 7

Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right

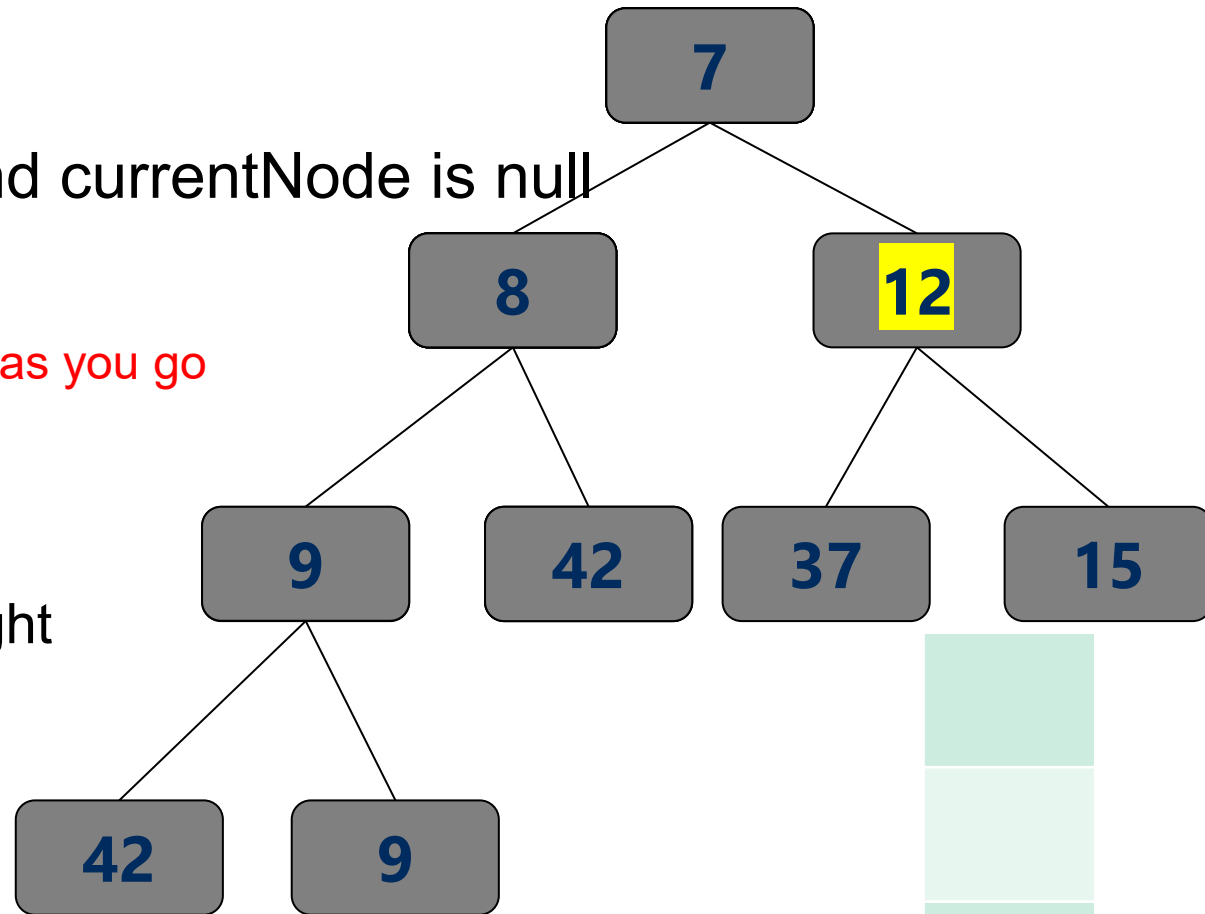


42, 9, 9, 8, 42, 7

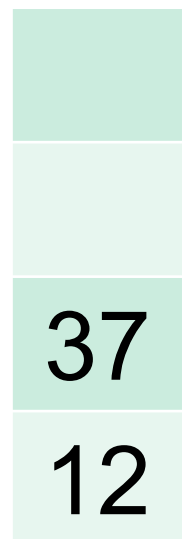
Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42, 9, 9, 8, 42, 7

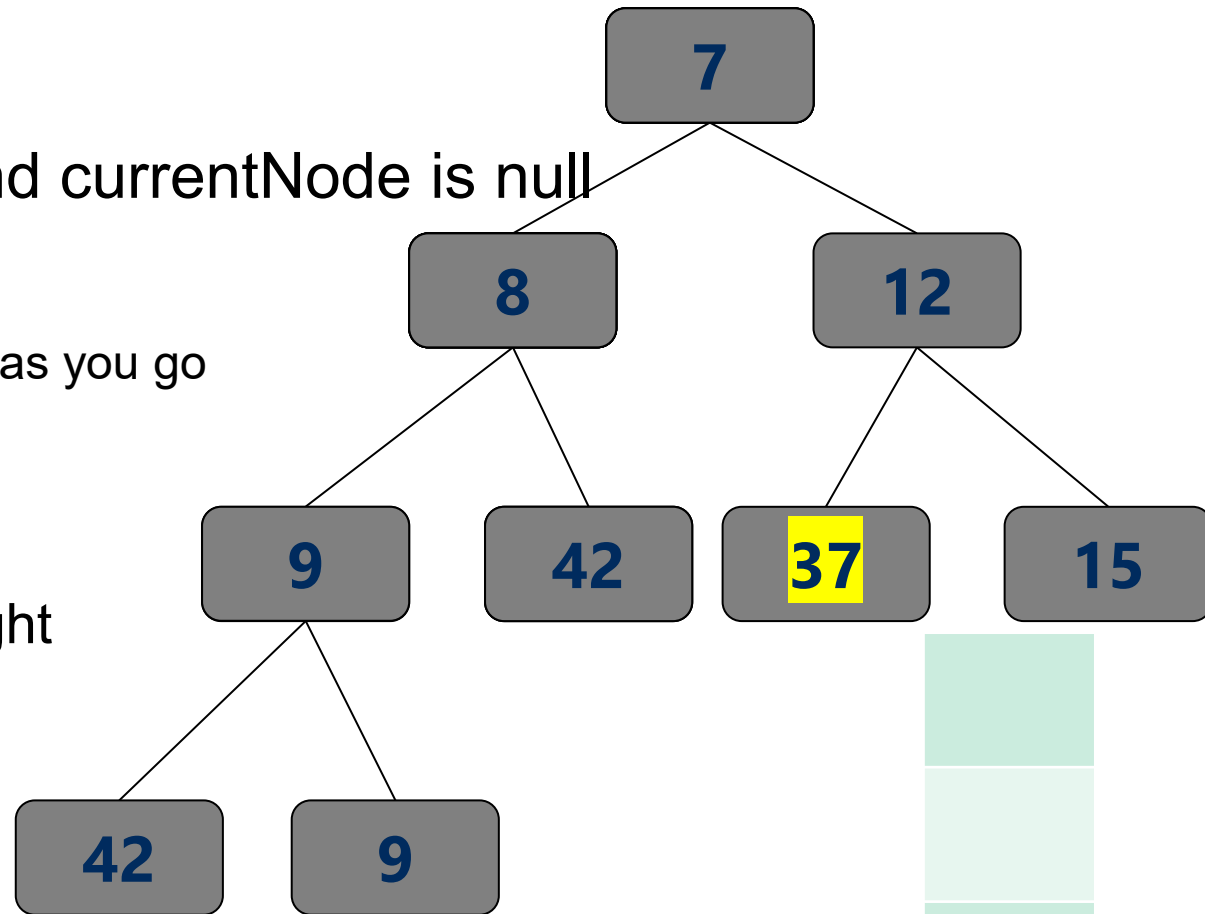


Stack

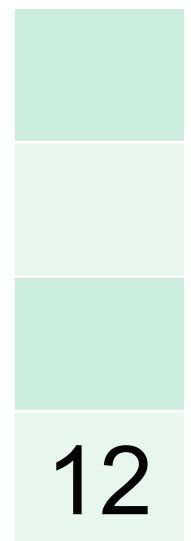
Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42, 9, 9, 8, 42, 7, 37

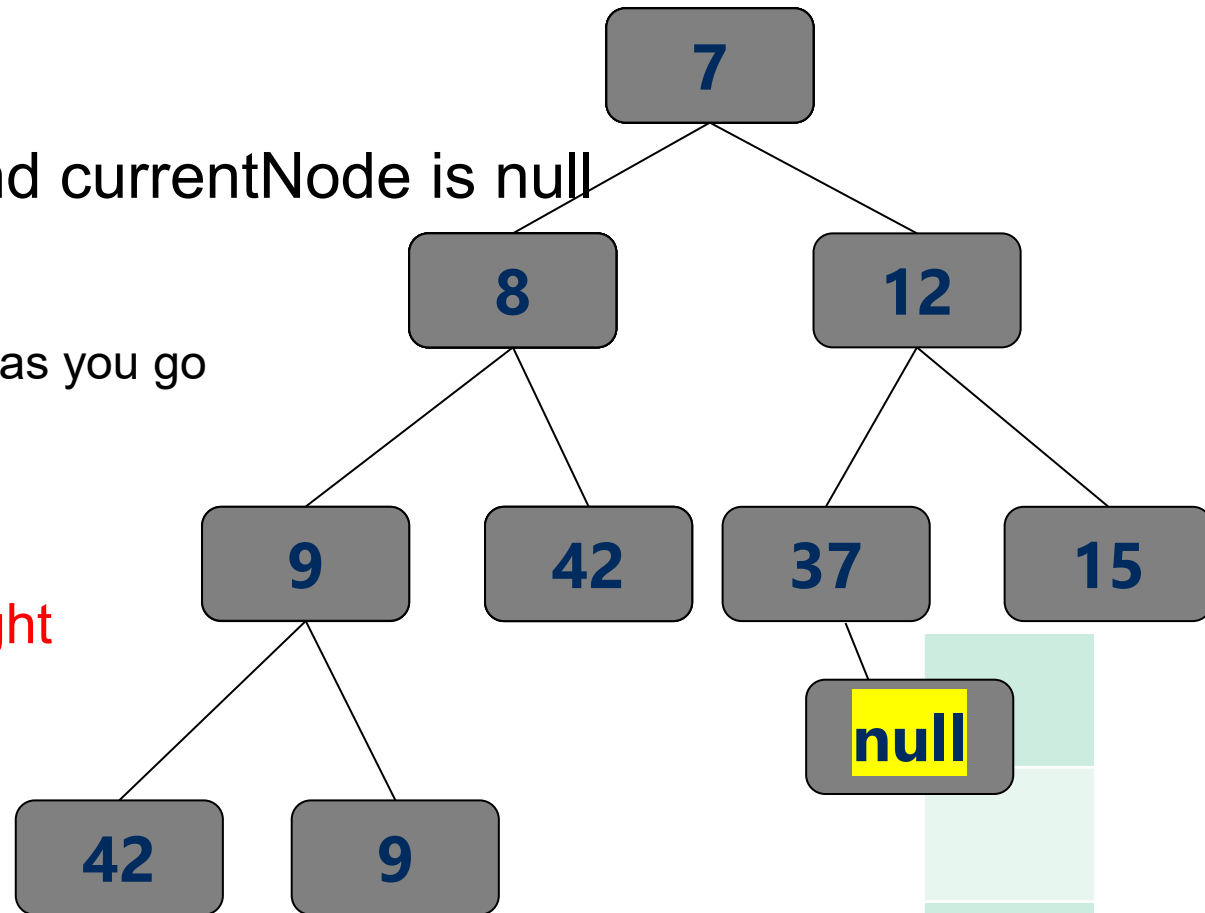


Stack

Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42, 9, 9, 8, 42, 7, 37

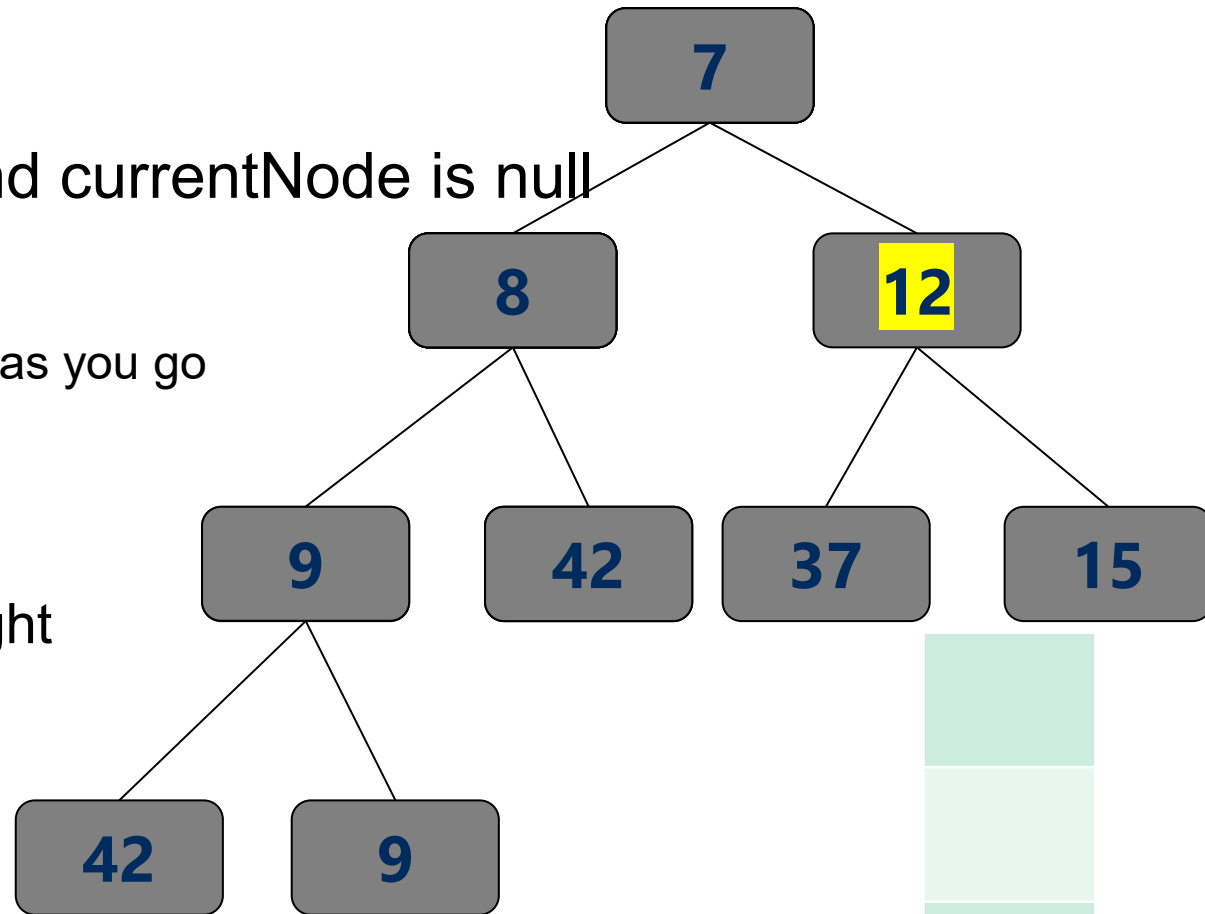
12

Stack

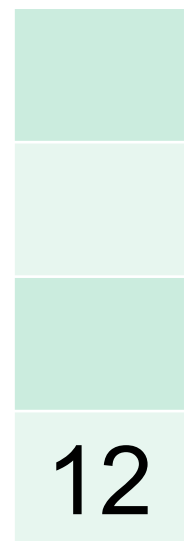
Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42, 9, 9, 8, 42, 7, 37

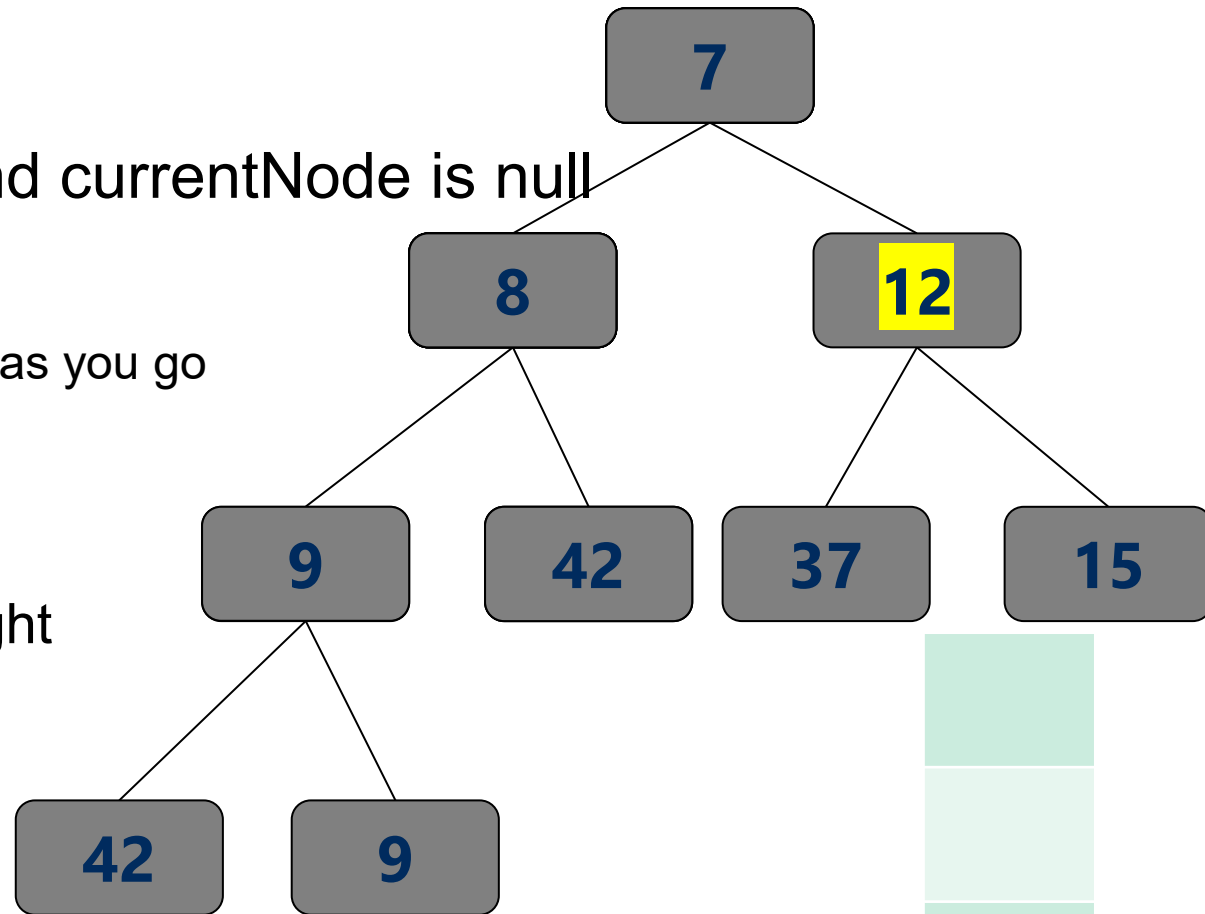


Stack

Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right

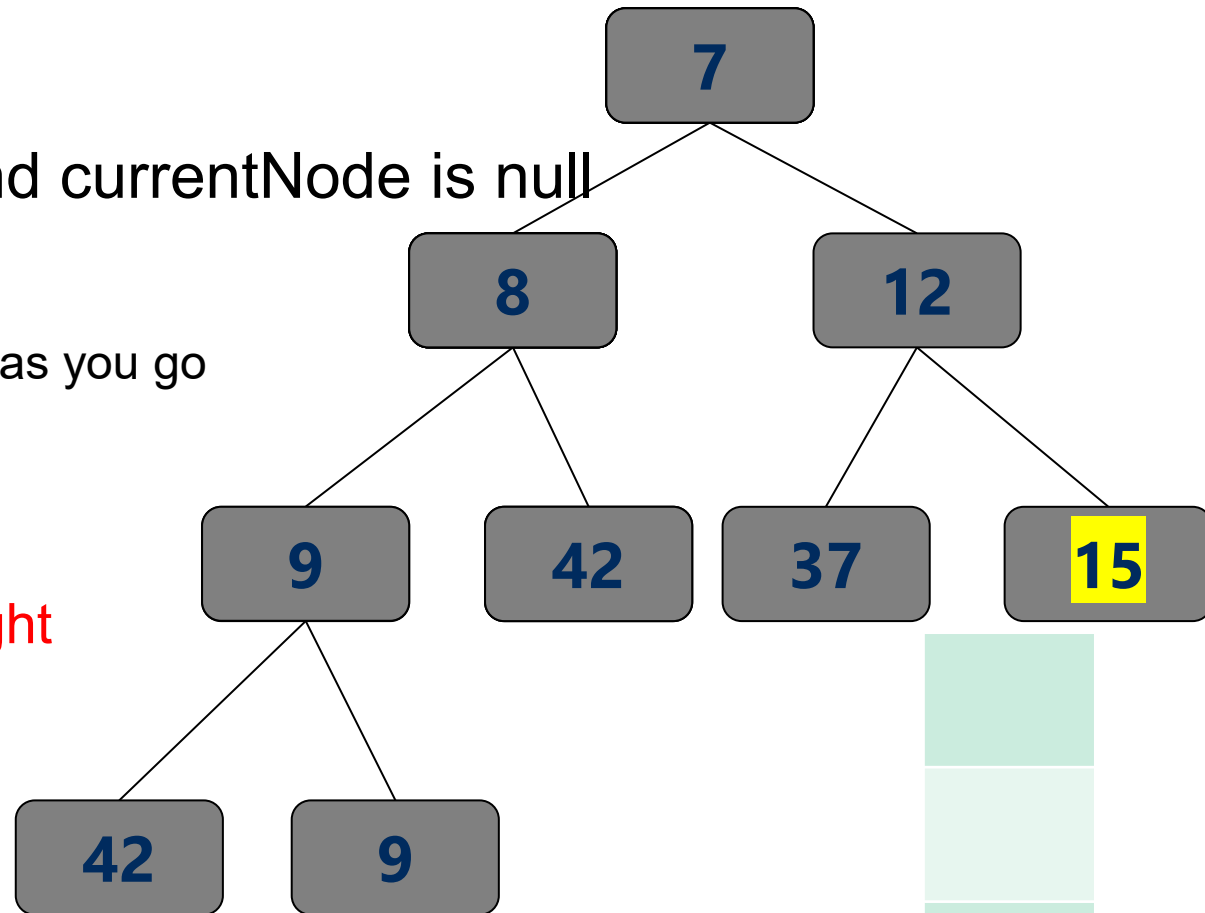


42, 9, 9, 8, 42, 7, 37, 12

Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right

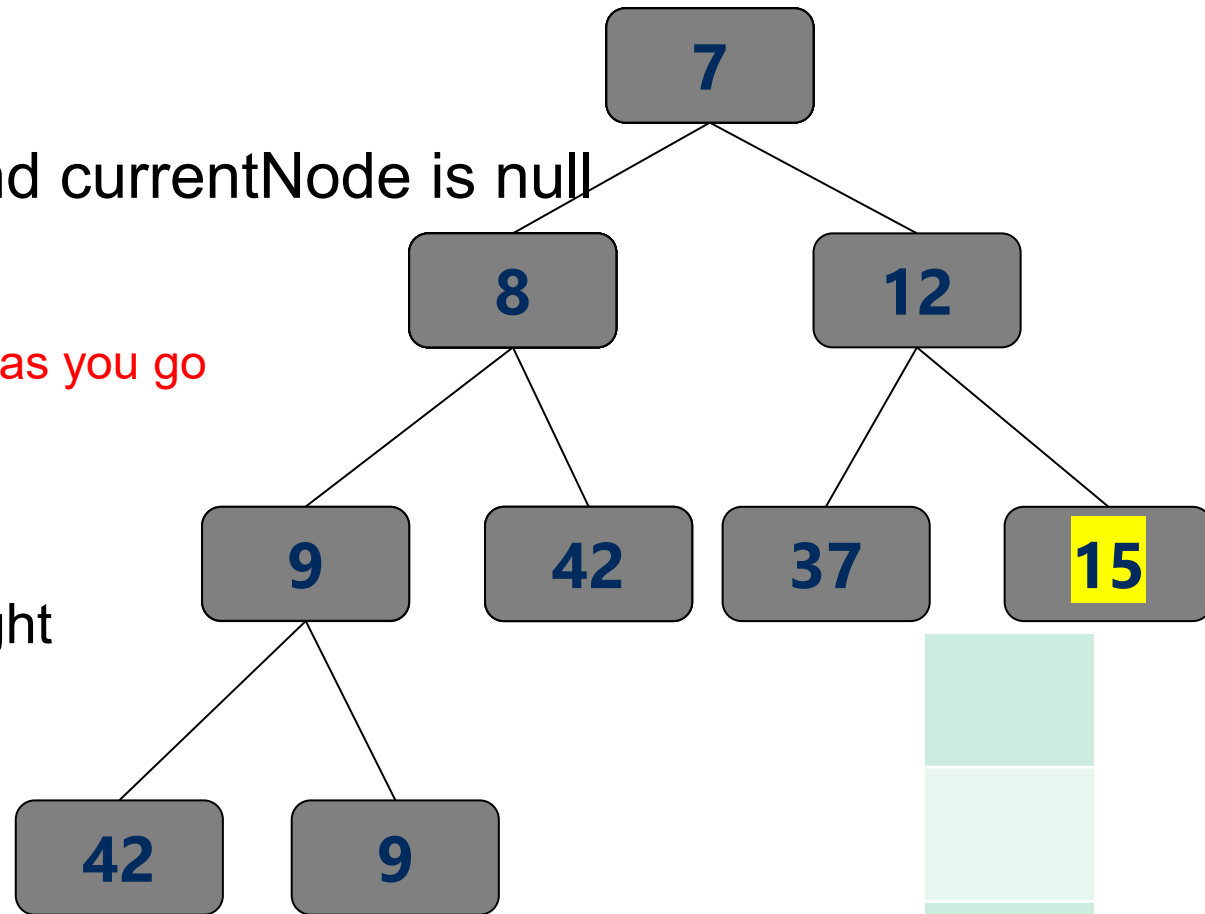


42, 9, 9, 8, 42, 7, 37, 12

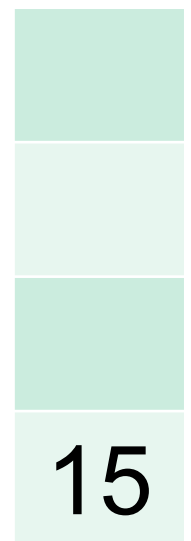
Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42, 9, 9, 8, 42, 7, 37, 12

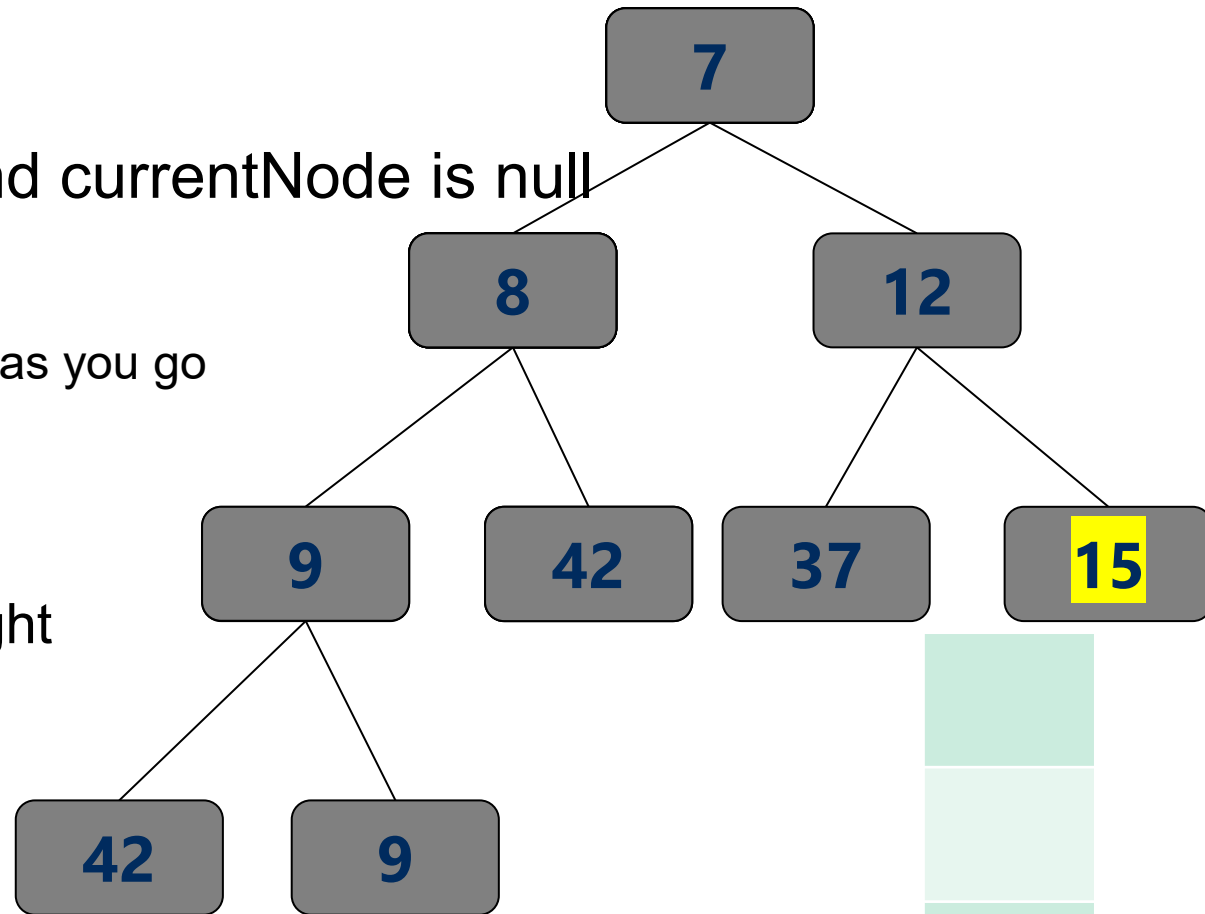


Stack

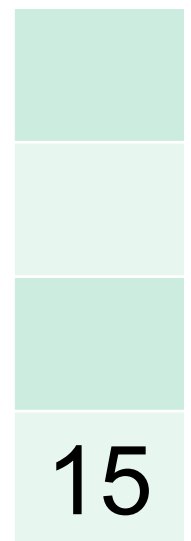
Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42, 9, 9, 8, 42, 7, 37, 12

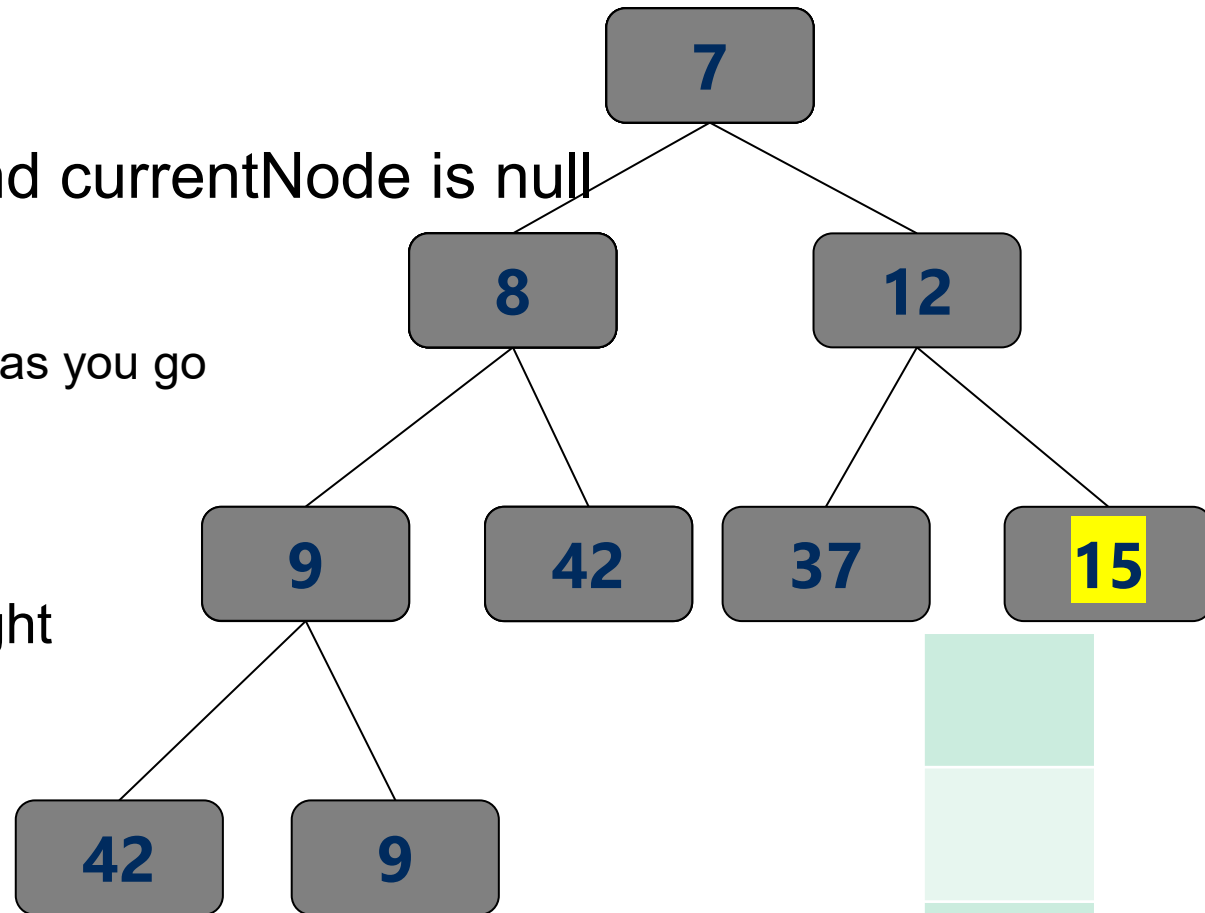


Stack

Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right

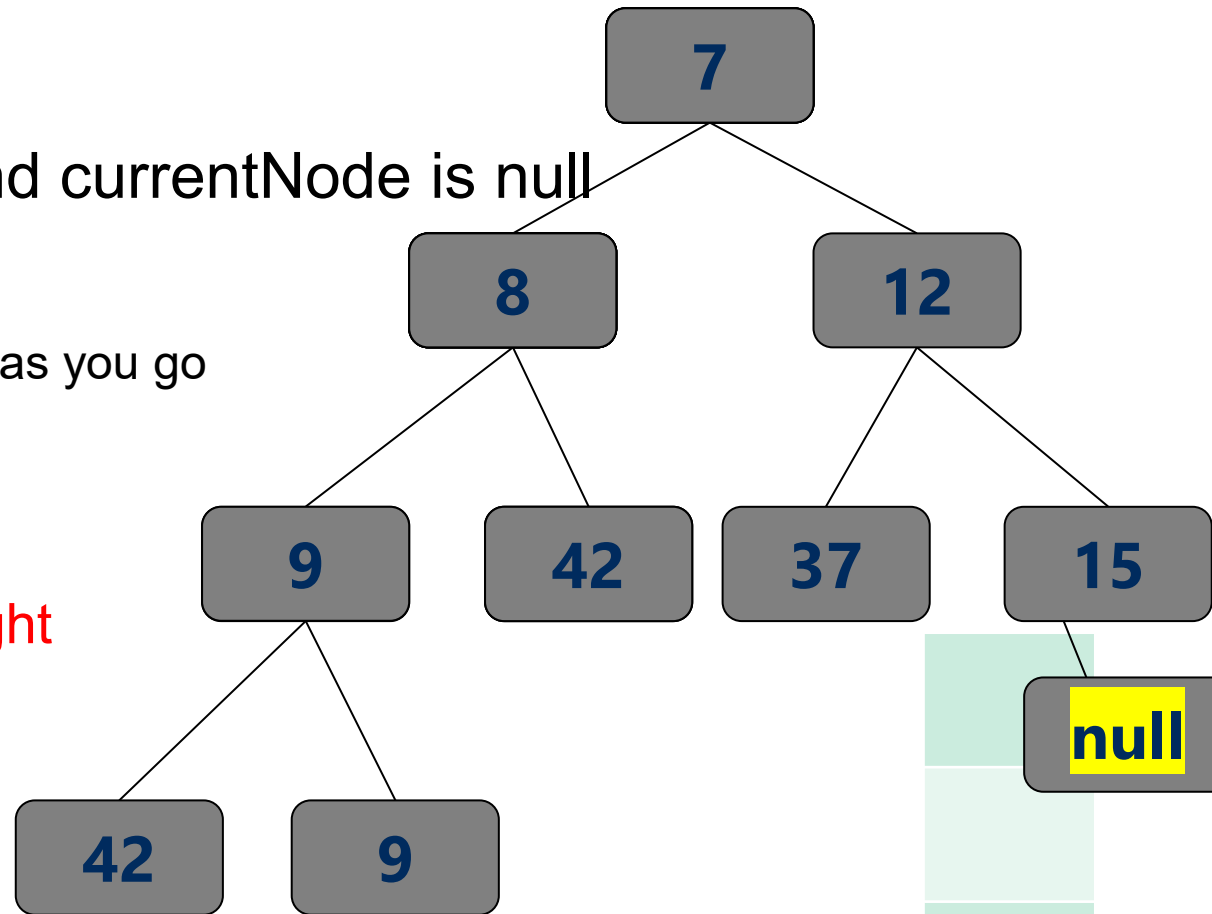


42, 9, 9, 8, 42, 7, 37, 12, 15

Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right

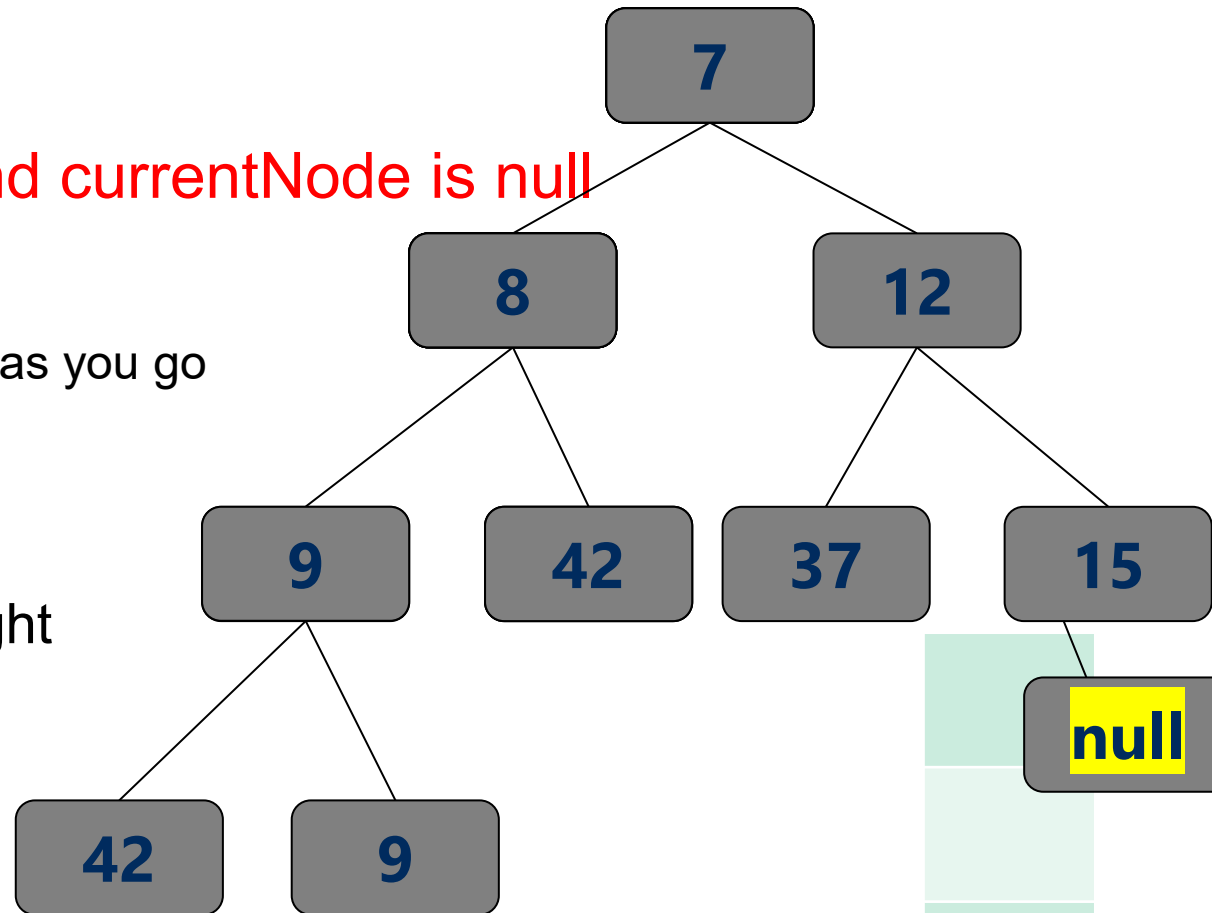


42, 9, 9, 8, 42, 7, 37, 12, 15

Muddiest Points

- Q: Review iterative inorder traversal

- currentNode = root
- repeat until stack is empty and currentNode is null
 - if currentNode != null
 - Move to leftmost node and push as you go
 - currentNode = pop()
 - visit currentNode
 - currentNode = currentNode.right



42, 9, 9, 8, 42, 7, 37, 12, 15

Burrows-Wheeler Transform Example

- ABABABA
- **Step 1: Build an array of 7 strings, each a circular rotation of the original by one character**
- ABABABA
- BABABAA
- ABABAAB
- BABAABA
- ABAABAB
- BAABABA
- AABABAB
- **Step 2: Sort the array alphabetically**
- **Notice that** the first column of the sorted array is the same as the last column of the original array
- all columns have the same set of letters
- **Step 3: Output the last column of the sorted array and the index of the input string in the sorted array**

original array

ABABABA
BABABAA
ABABAAB
BABAABA
ABAABAB
BAABABA
AABABAB

sorted array

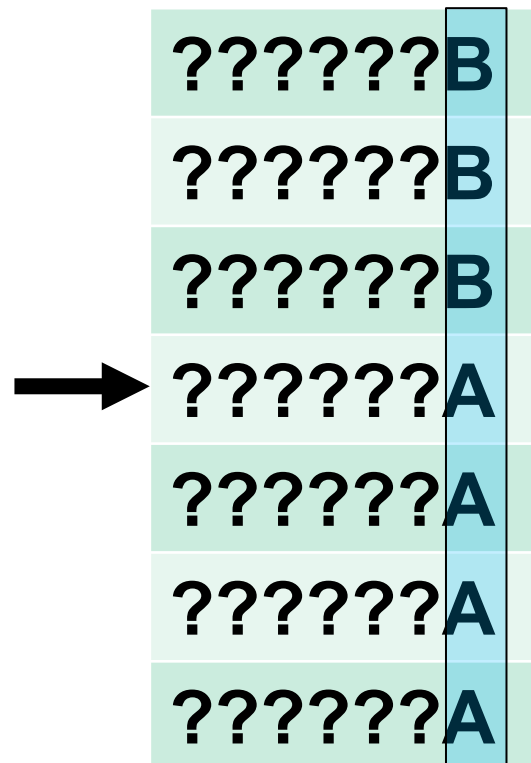
AABABAB
ABAABAB
ABABAAB
ABABABA
BAABABA
BABAABA
BABABAA



Output of BWT:
BBBAAAA and 3

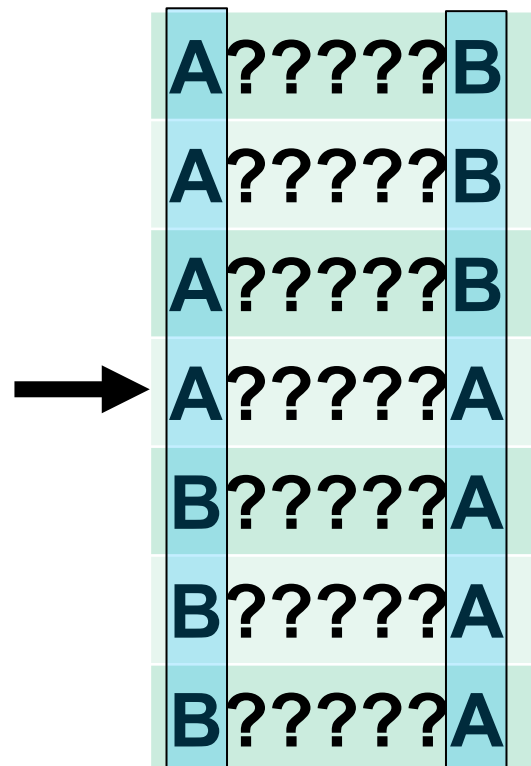
Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 1: Sort the encoded string**
 - BBBAAAA → AAAABBB
 - The first column of the sorted array has the same characters as the last column
 - but in sorted order



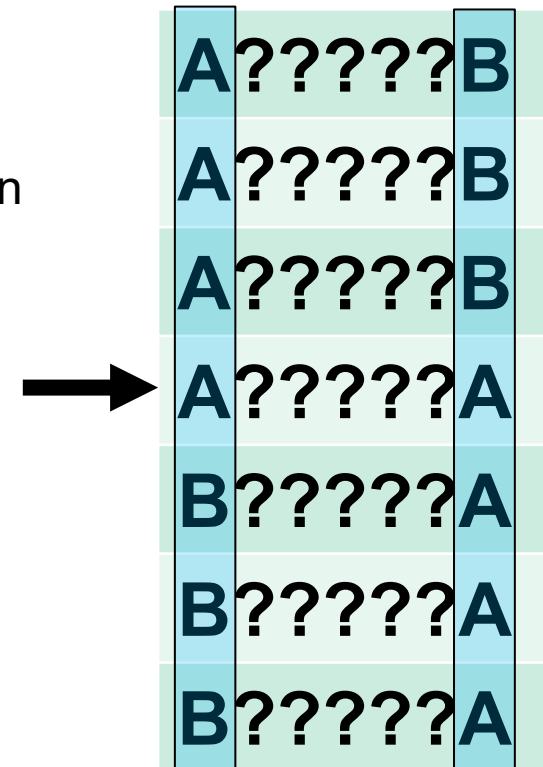
Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 1: Sort the encoded string**
 - BBBAAAA \rightarrow AAAABBB
 - This gives us the first column of the sorted array



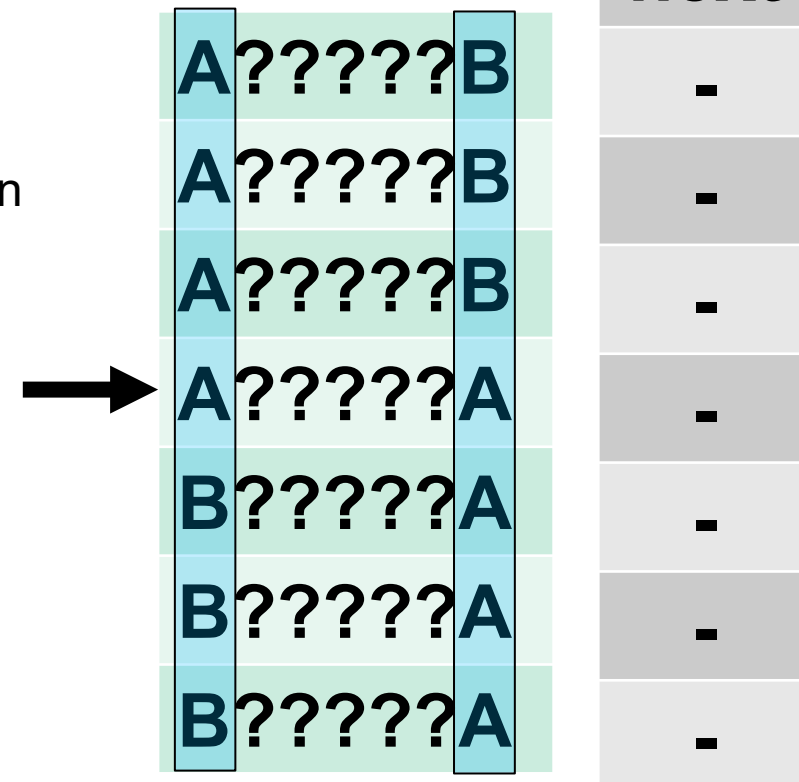
Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 2: Fill an array next[]**
 - defined for each entry in the sorted array
 - tells us the index in sorted array of the next string in the original array
 - Scan through the first column
 - for each row i holding character c
 - $\text{next}[i] = \text{first unassigned index of } c \text{ in the last column}$



Burrows-Wheeler Transform Decoding

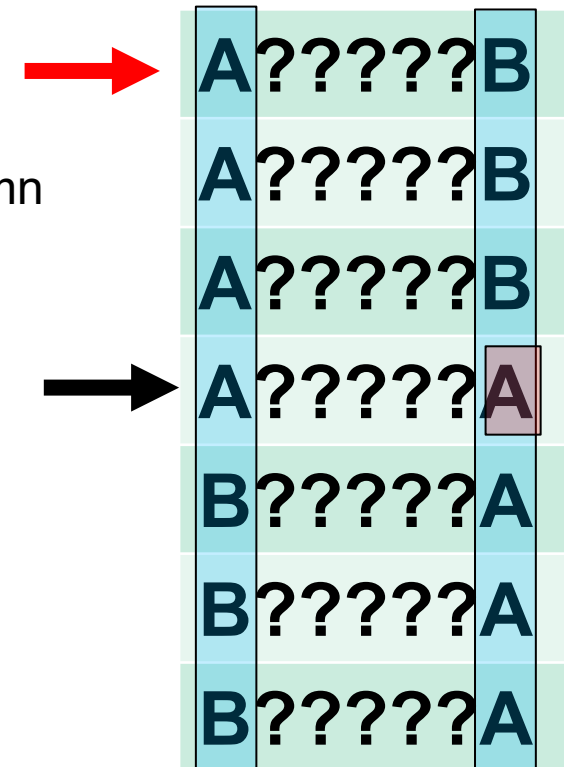
- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 2: Fill an array next[]**
 - defined for each entry in the sorted array
 - tells us the index in sorted array of the next string in the original array
 - Scan through the first column
 - for each row i holding character c
 - $\text{next}[i] = \text{first unassigned index of } c \text{ in the last column}$



	next
A?????B	-
A?????B	-
A?????B	-
A?????A	-
B?????A	-
B?????A	-
B?????A	-

Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 2: Fill an array next[]**
 - defined for each entry in the sorted array
 - tells us the index in sorted array of the next string in the original array
 - Scan through the first column
 - for each row i holding character c
 - $\text{next}[i] = \text{first unassigned index of } c \text{ in the last column}$



A	?????	B
A	?????	B
A	?????	B
A	?????	A
B	?????	A
B	?????	A
B	?????	A

next
3
-
-
-
-
-
-

Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 2: Fill an array next[]**
 - defined for each entry in the sorted array
 - tells us the index in sorted array of the next string in the original array
 - Scan through the first column
 - for each row i holding character c
 - $\text{next}[i] = \text{first unassigned index of } c \text{ in the last column}$

A	?????	B
A	?????	B
A	?????	B
A	?????	A
B	?????	A
B	?????	A
B	?????	A

next
3
4
-
-
-
-
-

Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 2: Fill an array next[]**
 - defined for each entry in the sorted array
 - tells us the index in sorted array of the next string in the original array
 - Scan through the first column
 - for each row i holding character c
 - $\text{next}[i] = \text{first unassigned index of } c \text{ in the last column}$

A	?????	B
A	?????	B
A	?????	B
A	?????	A
B	?????	A
B	?????	A
B	?????	A

next
3
4
-
-
-
-
-

Burrows-Wheeler Transform Decoding

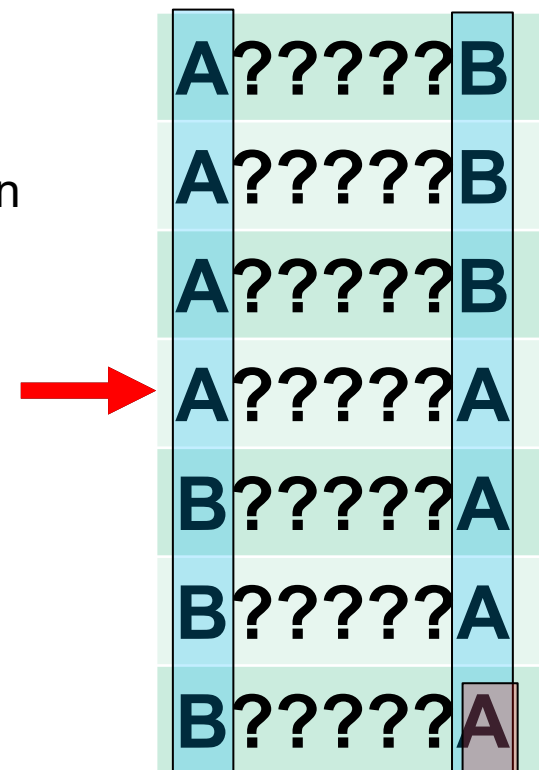
- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 2: Fill an array next[]**
 - defined for each entry in the sorted array
 - tells us the index in sorted array of the next string in the original array
 - Scan through the first column
 - for each row i holding character c
 - $\text{next}[i] = \text{first unassigned index of } c \text{ in the last column}$

A	?	?	?	?	?	B
A	?	?	?	?	?	B
A	?	?	?	?	?	B
A	?	?	?	?	?	A
B	?	?	?	?	?	A
B	?	?	?	?	?	A
B	?	?	?	?	?	A

next
3
4
5
-
-
-
-

Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 2: Fill an array next[]**
 - defined for each entry in the sorted array
 - tells us the index in sorted array of the next string in the original array
 - Scan through the first column
 - for each row i holding character c
 - $\text{next}[i] = \text{first unassigned index of } c \text{ in the last column}$



A	?????	B
A	?????	B
A	?????	B
A	?????	A
B	?????	A
B	?????	A
B	?????	A

next
3
4
5
6
-
-
-

Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 2: Fill an array next[]**
 - defined for each entry in the sorted array
 - tells us the index in sorted array of the next string in the original array
 - Scan through the first column
 - for each row i holding character c
 - $\text{next}[i] = \text{first unassigned index of } c \text{ in the last column}$

A	?	?	?	?	?	B
A	?	?	?	?	?	B
A	?	?	?	?	?	B
A	?	?	?	?	?	A
B	?	?	?	?	?	A
B	?	?	?	?	?	A
B	?	?	?	?	?	A

next
3
4
5
6
0
-
-

Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 2: Fill an array next[]**
 - defined for each entry in the sorted array
 - tells us the index in sorted array of the next string in the original array
 - Scan through the first column
 - for each row i holding character c
 - $\text{next}[i] = \text{first unassigned index of } c \text{ in the last column}$

A	?	?	?	?	?	B
A	?	?	?	?	?	B
A	?	?	?	?	?	B
A	?	?	?	?	?	A
B	?	?	?	?	?	A
B	?	?	?	?	?	A
B	?	?	?	?	?	A

next
3
4
5
6
0
1
-

Burrows-Wheeler Transform Decoding


- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 2: Fill an array next[]**
 - defined for each entry in the sorted array
 - tells us the index in sorted array of the next string in the original array
 - Scan through the first column
 - for each row i holding character c
 - $\text{next}[i] = \text{first unassigned index of } c \text{ in the last column}$

A	?????	B
A	?????	B
A	?????	B
A	?????	A
B	?????	A
B	?????	A
B	?????	A

next
3
4
5
6
0
1
2

Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 2: Fill an array next[]**
 - defined for each entry in the sorted array
 - tells us the index in sorted array of the next string in the original array
 - Scan through the first column
 - for each row i holding character c
 - $\text{next}[i] = \text{first unassigned index of } c \text{ in the last column}$
- Why does that work?
 - first character of a string becomes the last character in the next string in the original order



	A	?	?	?	?	B		next
	A	?	?	?	?	B		3
	A	?	?	?	?	B		4
	A	?	?	?	?	B		5
	A	?	?	?	?	A		6
	B	?	?	?	?	A		0
	B	?	?	?	?	A		1
	B	?	?	?	?	A		2

Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 3: Recover the input string using the next[] array**
- We can conclude that A is the first character in the input string
 - why?

A???????

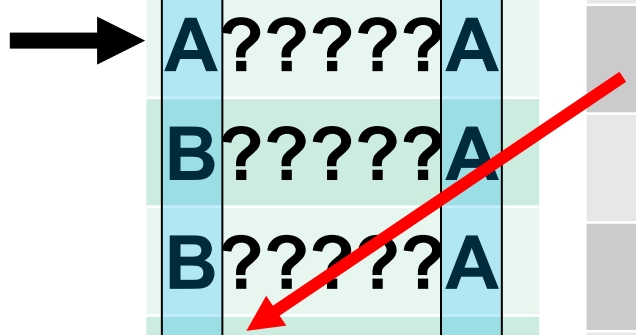


	next
A??????B	3
A??????B	4
A??????B	5
A??????A	6
B??????A	0
B??????A	1
B??????A	2

Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 3: Recover the input string using the next[] array**
- We can conclude that A is the first character in the input string
 - why?
- The next character is the first character of the next string in the original order
 - first character in string at next[3]

AB?????

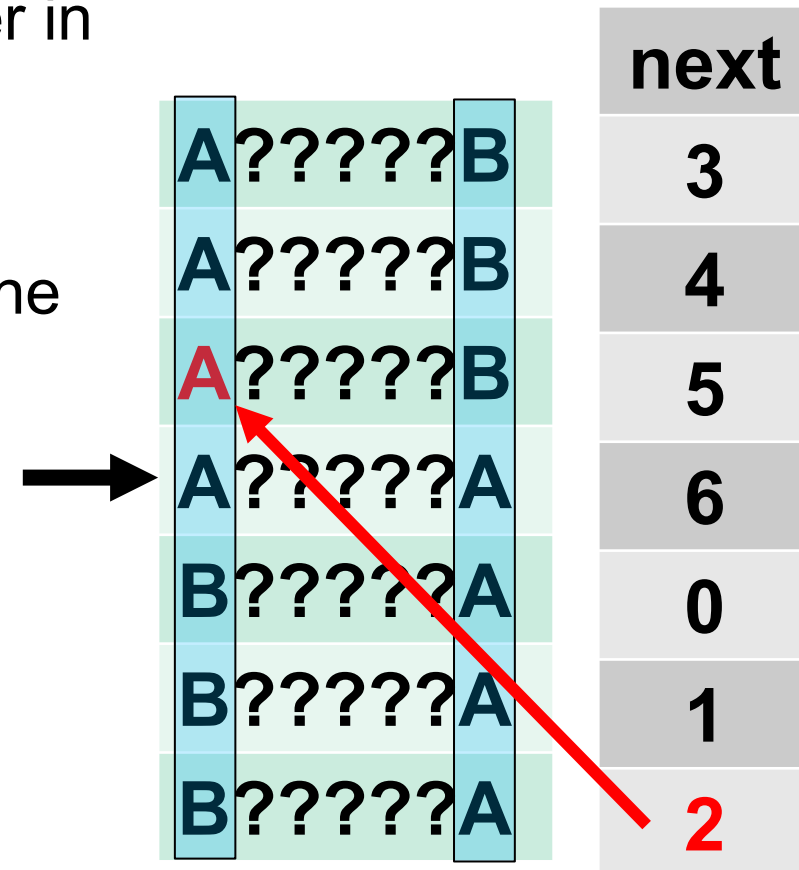


	next
A?????B	3
A?????B	4
A?????B	5
A?????A	6
B?????A	0
B?????A	1
B?????A	2

Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 3: Recover the input string using the next[] array**
- We can conclude that A is the first character in the input string
 - why?
- The next character is the first character of the next string in the original order
 - first character in string at next[6]

ABA????

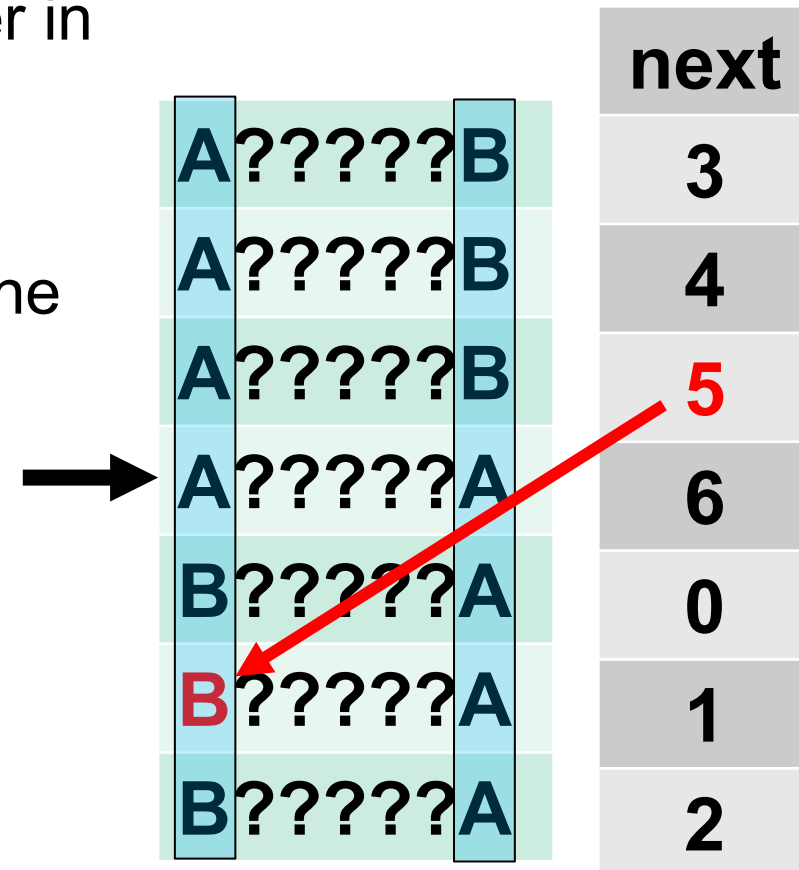


	next
A?????B	3
A?????B	4
A?????B	5
A?????A	6
B?????A	0
B?????A	1
B?????A	2

Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- Step 3: Recover the input string using the next[] array**
- We can conclude that A is the first character in the input string
 - why?
- The next character is the first character of the next string in the original order
 - first character in string at next[2]

ABAB???

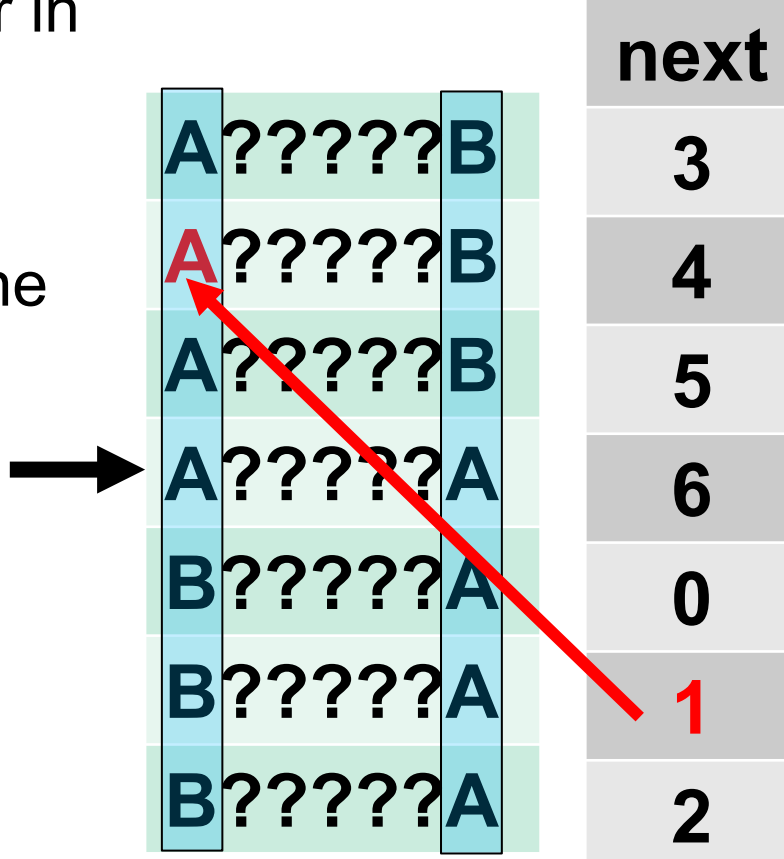


	next
A?????B	3
A?????B	4
A?????B	5
A?????A	6
B?????A	0
B?????A	1
B?????A	2

Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- Step 3: Recover the input string using the next[] array**
- We can conclude that A is the first character in the input string
 - why?
- The next character is the first character of the next string in the original order
 - first character in string at next[5]

ABABA??



	next
A?????B	3
A?????B	4
A?????B	5
A?????A	6
B?????A	0
B?????A	1
B?????A	2

Burrows-Wheeler Transform Decoding

- Output of BWT:
 - BBBAAAA and 3
- How can we recover ABABABA?
- **Step 3: Recover the input string using the next[] array**
- We can conclude that A is the first character in the input string
 - why?
- The next character is the first character of the next string in the original order
 - first character in string at next[5]

ABABABA



	next
A?????B	3
A?????B	4
A?????B	5
A?????A	6
B?????A	0
B?????A	1
B?????A	2