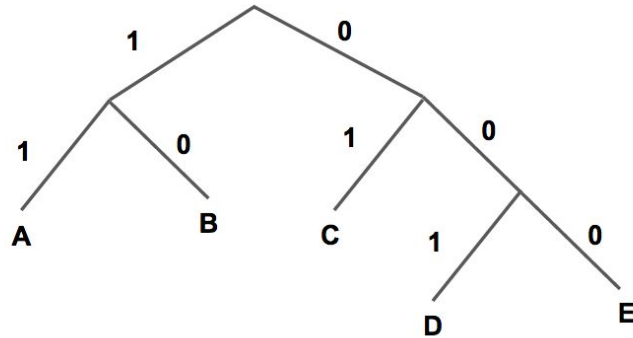


# Parallelizing Huffman Encoding / Decoding



Malcolm Mashig (mjm2396)

# Huffman Coding Overview

- 1) Creation of a code tree based on character likelihood/frequency
- 2) Encode each character in the text via traversal of the tree
- 3) Decode bit sequences via traversal of the tree

**Purpose:** optimize compression ratio by giving small codes to frequent characters

**Compression Ratio:**  $(\# \text{ of bits in compressed text}) / (\# \text{ of bits in original text})$

*...as I define it*

# Parallelization

1. Split input text into batches, then encode in parallel
2. Split encoded bit sequence into batches, then decode in parallel

## Caveats:

- What about creating the code tree?
  - Should we create a separate code tree per batch?
  - How do we store the code tree so that it can be independently known during decoding
- How do we split encoded bit sequence into batches?
  - **Character codes vary in length, and so our batches may split up individual characters**

# Approach

## Example Parallel Huffman Encoding File:

Number of Batches: 4

First Index: 828

Second Index: 1657

Third Index: 2469

Encoding: 0000000011000111101111101110001110110101000  
1111000011001101000101100110111101110101011  
0001011001100111111110000100110001001001010  
0010001000110101010100010000111011010100001  
1111110...

*Indices are determined by  
generating a cumulative sum of the  
bit sequence lengths during  
encoding*

# Implementation

```
let decInds = getInds dec_instr
let decBits = getBits dec_instr
let decBatches = batchBits decBits decInds
let decodings = Prelude.map (decode tree) decBatches `using` parList rdeepseq
let decoded = concat decodings
```

# Small-Scale Test

ORIGINAL TEXT:

But there is a problem. To decode the message, you need the coding tree for that message. I'm sure there are a variety of techniques to transmit the tree as part of the message. One of the things Apelfeld did that intrigued me was develop a stack based language and interpreter to represent the morse code tree as a string of characters. I think this would make a fine way to encode the tree into the front of the message. We should be able to extract a string of characters from the Huffman tree for a particular message and use that string of characters to reconstruct the tree on the other side. Now this will add some significant overhead to the output for a given message, but I'm not necessarily concerned with space efficiency here as much as learning about these kinds of transformations.

ENCODED:

[illegible]

DECODED:

But there is a problem. To decode the message, you need the coding tree for that message. I'm sure there are a variety of techniques to transmit the tree as part of the message. One of the things Apelfeld did that intrigued me was develop a stack based language and interpreter to represent the morse code tree as a string of characters. I think this would make a fine way to encode the tree into the front of the message. We should be able to extract a string of characters from the Huffman tree for a particular message and use that string of characters to reconstruct the tree on the other side. Now this will add some significant overhead to the output for a given message, but I'm not necessarily concerned with space efficiency here as much as learning about these kinds of transformations.

## LOSSLESS?

True

# Large-Scale Test

Shakespeare Works: <http://www.gutenberg.org/files/100/100-0.txt>

Length: 5,716,512 characters

Compression Ratio: ~0.60

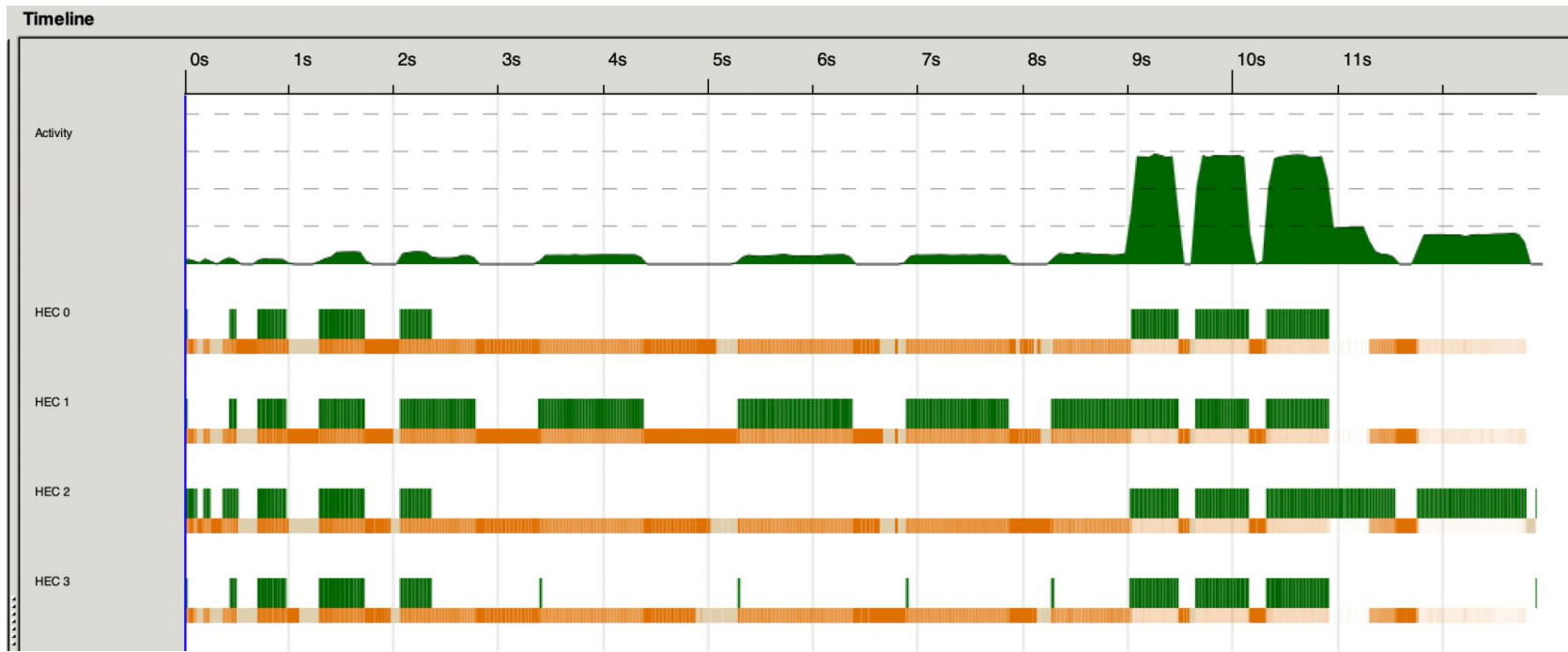
real 18.236s (100 batches)

real 16.432s (500 batches)

real 16.064s (2000 batches)

real 13.146s (10000 batches)

# Threadscope





# What I've Learned + Next Steps

- Still investigating improvements to achieve speedup
- Compression is sequential by nature, and parallelization is difficult

Next steps:

- Attempt other parallelization strategies and compare