

ESE534: Computer Organization

Day 21: April 4, 2012
Lossless Data Compression



Today

- Basic Idea
- Example
- Systolic
- Tree
- Reclaiming space in tree
- CAMs

Lossless vs. Lossy

- Lossless – can reconstruct source perfectly (bit-identical)
 $\text{uncompress}(\text{compress}(x))=x$
- Lossy – capture important elements of original, but maybe not all bits
- Examples
 - Huffman
 - Run Length Coding
 - Lempel-Ziv
 - Unix compress/gzip
- Examples
 - MP3
 - JPEG
 - MPEG

Dictionary Idea

- Send id for long string rather than all the characters

Dictionary Example

- “the instruction controls the behavior of the ALU, data memory, and interconnect on each cycle.”
- Characters?
 - Bits at 8b/character
- Encoding with dictionary? Bits?

Code	Word
000	ALU
001	memory
010	interconnect
011	instruction
100	data
101	cycle
110	control
111	the

Dictionary Usability

- When can we do this?
- What might prevent us from pulling this trick?

Big Idea

- Use data already sent as the dictionary
 - Don't need to pre-arrange dictionary
 - Adapt to common phrases/idioms in a particular document

Example

- First line of Dr. Suess's *Green Eggs and Ham*
 - I AM SAM SAM I AM
- Recurring substrings?

Example

- An encoding:
- I AM S<2,3> <5,4><0,4>
- Decode.
- Characters in original?
 - Bits based on 8b characters?

Example

- An encoding:
- I AM S<2,3> <5,4><0,4>
- Encode:
 - Add 1 bit to identify character vs <x,y>
 - 9b characters
 - <x,y>: 1b says this + 4b for x, 4b for y
 - Also 9b
- How many bits?

Technical Issue

- How many bits assign to x and y?
- Issues?
- What if the document is huge?
 - What problems might that pose?

Windows

- Pragmatic solution
 - Only keep the last D characters
 - D is window size
 - Need $\log_2(D)$ bits to specify a position
 - Parameterize encoder based on D
 - Typically larger D \rightarrow Greater compression

Encoding

- Greedy simplification
 - Encode by successively selecting the longest match between the head of the remaining string to send and the current window

Penn ESE534 Spring2012 -- DeHon

13

Algorithm Concept

- While data to send
 - Find largest match in window
 - If length=1
 - Send character
 - Else
 - Send $\langle x, y \rangle = \langle \text{match-pos}, \text{length} \rangle$
 - Shift data encoded into window

Penn ESE534 Spring2012 -- DeHon

14

Run Algorithm

- Use D=8
- I AM SAM SAM I AM
- How many bits?

Penn ESE534 Spring2012 -- DeHon

15

What's challenging to implement?

- While data to send
 - Find largest match in window
 - If length=1
 - Send character
 - Else
 - Send $\langle x, y \rangle = \langle \text{match-pos}, \text{length} \rangle$
 - Shift data encoded into window

Penn ESE534 Spring2012 -- DeHon

16

Systolic Algorithm

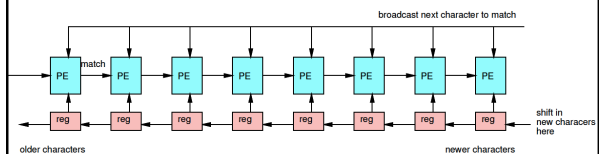
- Give character in window a PE
- Broadcast characters to PE
 - While (some PE has match-out=true)
 - match-out=new-search*match
 - OR cont-search*match*match-in
 - Len=len+1
 - Broadcast next character
 - Send $\langle \text{pos-last-match-out}, \text{len} \rangle$
 - Shift last set of characters into window

Penn ESE534 Spring2012 -- DeHon

17

Systolic Hardware

- While (some PE has match-out=true)
 - match-out=new-search*match
 - OR cont-search*match*match-in
 - Len=len+1
 - Broadcast next character



Simulate Systolic

- Each student is a PE in the window
 - Identify left and right neighbors
 - Raise right hand for match-out
 - Note left neighbor's hand at end of previous cycle to know match-in
- I AM SAM SAM I AM

Penn ESE534 Spring2012 -- DeHon

19

Contemplate Solution

- How complicated is each PE?
- How fast PE?
- How fast does encoding operate?
- How much area do we need?
- How much energy?

Penn ESE534 Spring2012 -- DeHon

20

Contemplate Solution

- What's inefficient or unsatisfying about this solution?

Penn ESE534 Spring2012 -- DeHon

21

Tree Based

Penn ESE534 Spring2012 -- DeHon

22

Idea

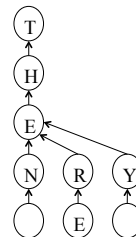
- Avoid need to track multiple substrings
 - Compress storage
- BY
- Storing common prefixes together in a tree

Penn ESE534 Spring2012 -- DeHon

23

Tree Example

- THEN AND THERE, THEY STOOD...



Penn ESE534 Spring2012 -- DeHon

24

Idea

- Avoid need to track multiple substrings
 - Compress storage
- BY
- Storing common prefixes together in a tree

Penn ESE534 Spring2012 -- DeHon

25

Tree Algorithm

Root for each character

- Follow tree according to input until no more match
- Send <name of last tree node>
- Extend tree with new character
- Start over with this character

Penn ESE534 Spring2012 -- DeHon

26

Run Algorithm

- I AM SAM SAM I AM

Penn ESE534 Spring2012 -- DeHon

27

Encoding

- Encoding bits assuming $D=512$
 - So, 9b to encode tree node

Penn ESE534 Spring2012 -- DeHon

28

Finite Window

- How can we maintain a finite window in this case?

Penn ESE534 Spring2012 -- DeHon

29

Finite Window

- Clear and start over
- LRU on tree nodes
- Maintain two areas
 - Encode from one (perhaps both)
 - Add to new
 - When new fills,
 - New->old, clear old
- Pick old leaf node to replace

Penn ESE534 Spring2012 -- DeHon

30

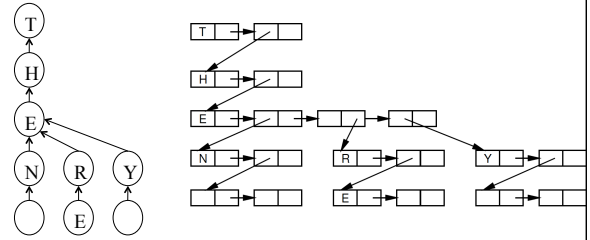
Complexity

- How much work per character to encode?

Penn ESE534 Spring2012 -- DeHon

31

Tree Node Representation

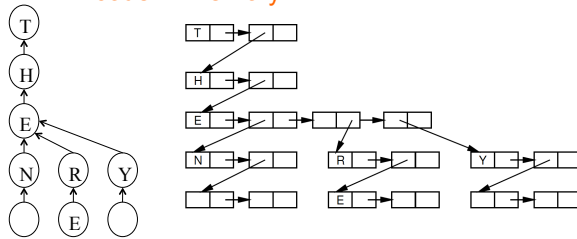


Penn ESE534 Spring2012 -- DeHon

32

Tree Node Representation

- Encode in memory



Penn ESE534 Spring2012 -- DeHon

33

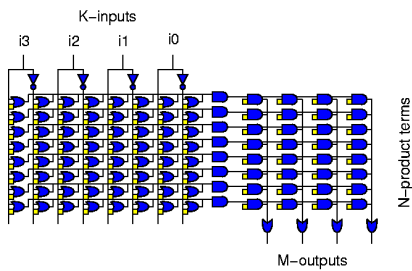
Content Addressable Memory

- What's a CAM?

Penn ESE534 Spring2012 -- DeHon

34

PLA



Penn ESE534 Spring2012 -- DeHon

35

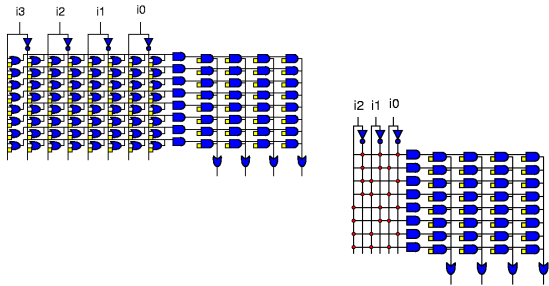
CAM

- Memory with Programmable Addresses
 - Capacity $< 2^{(\text{matchbits})}$
- PLA with both planes writeable

Penn ESE534 Spring2012 -- DeHon

36

PLA and Memory



Penn ESE534 Spring2012 -- DeHon

Contemplate

- What value do Bunton and Borriello get from using a CAM?

Penn ESE534 Spring2012 -- DeHon

38

Admin

- Reading for Monday on Web
- FM1 for Monday
 - Implement tree version on processor and estimate energy

Penn ESE534 Spring2012 -- DeHon

39

Big Ideas [MSB Ideas]

- Can often compress data without loss of information
- Exploit structure in data to encode
- Build dictionary based on data already sent
- Code repeating substrings compactly in terms of data already seen in recent past

Penn ESE534 Spring2012 -- DeHon

40