

ESE532: System-on-a-Chip Architecture

Day 16: October 25, 2017
Deduplication and Compression Project

Midterm: average 56, std. dev 16



Penn ESE532 Fall 2017 -- DeHon

Midterm

- Solution ... (not out, next few days...)
- Suspect bit time constrained
- Biggest role prepare you for final
 - Know what these exams look like
 - Don't Panic – but take as serious diagnostic
 - 10% of grade
 - Will replace midterm grade with final exam grade if that is higher

Penn ESE532 Fall 2017 -- DeHon

2

Today

- Motivation
- Project
- Content-Defined Chunking
- Hashing / Deduplication
- LZW Compression

- Exam in after us today
 - Try to finish 4:20pm
 - ...and we should clear room for them

Penn ESE532 Fall 2017 -- DeHon

3

Message

- Can reduce data size by identifying and reducing redundancy
- Can spend computation and data storage to reduce communication traffic

Penn ESE532 Fall 2017 -- DeHon

4

Problem

- Always want more
 - Bandwidth
 - Storage space
- Carry data with me (phone, laptop)
- Backup laptop, phone data
 - Maybe over limited bw links
- Never delete data
- Download movies, books, datasets
- Make most use of space, bw given

Penn ESE532 Fall 2017 -- DeHon

5

Opportunity

- Significant redundant content in our raw data streams (data storage)
- **More formally:**
 - Information content < raw data
- Reduce the data we need to send or store by identifying redundancies

Penn ESE532 Fall 2017 -- DeHon

6

Example

- Two identical files
 - Different parts of my file systems
- Don't store separate copies
 - Store one
 - And the other says "same as the first file"
 - e.g. keep a pointer

Penn ESE532 Fall 2017 -- DeHon

7

Why Identical?

- Eniac file system (common file server)
 - Multiple students have copies of assignment(s)
 - Snapshots (.snapshot)
 - Has copies of your directory an hour ago, days ago, weeks ago
 - ...but most of that data hasn't changed

Penn ESE532 Fall 2017 -- DeHon

8

Broadening

- History file systems
 - snapshot, Apple Time Machine
- Version Control (git, svn)
- Manually keep copies
- Download different software release versions
 - With many common files

Penn ESE532 Fall 2017 -- DeHon

9

Cloud Data Storage

- E.g. Drop Box, Apple Cloud
- Saves data for large class of people
 - Want to only store one copy of each
- Synchronize with local copy on phone/laptop
 - Only want to send one copy on update
 - Only want to send changes
 - Data not already known on other side
 - (or, send that data compactly by just naming it)

Penn ESE532 Fall 2017 -- DeHon

10

Placement

- At file server
 - Deduplicate/compress data as stored
- In client
 - Dedup/compress to send to server
- In data center network
 - Dedup/compress data to send between server
- Network infrastructure
 - Dedup/compress from central to regional server

Penn ESE532 Fall 2017 -- DeHon

11

Optimizing the Bottleneck

- Saving data (transmitted, stored)
- By spending compute cycles
 - And storage database
- When communication (storage) is the bottleneck
 - We're willing to spend computation to better utilize the bottleneck resource

Penn ESE532 Fall 2017 -- DeHon

12

Project

Penn ESE532 Fall 2017 -- DeHon

13

Project

- Perform deduplication/compression at network speeds (1Gb/s, 10Gb/s)
- Use “chunks” instead of files
- Turn a raw/uncompressed data stream into one that exploits
 - Duplicate chunks
 - Redundancies within chunks

Penn ESE532 Fall 2017 -- DeHon

14

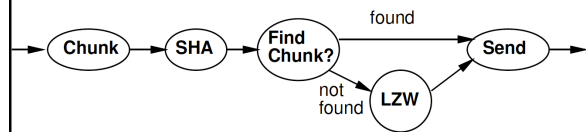
Project Context

- File server input link from network
 - Compress data before sending to disk
- Network link in data center or infrastructure
 - Compress data that goes over network

Penn ESE532 Fall 2017 -- DeHon

15

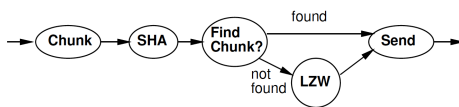
Project Task



Penn ESE532 Fall 2017 -- DeHon

16

Content-Defined Chunking



Penn ESE532 Fall 2017 -- DeHon

17

Files or chunks?

- Why files might be wrong granularity?

Penn ESE532 Fall 2017 -- DeHon

18

Blocks

- We regularly cut files into fixed-sized blocks
 - Disk sectors or blocks
 - inodes in File systems
- Why might fixed-sized blocks not be right division for deduplication?

Common Modifications

- Add a line of text
- Remove a line of text
- Fix a typo
- Rewrite a paragraph
- Trim or compose a video sequence

Content-Define Chunking

- Would like to re-align pieces around unchanged/common sequences
 - Around the content
- Break up larger thing (file) into pieces based on features of content

Preclass 1 and 2

- How much duplication opportunity in
 - Preclass 1 blocks?
 - Preclass 2 chunks?
- Why chunks able to do better?

Chunks

- Pieces of some larger file (data stream)
- Variable size
 - Over a limited range
- Discretion in how formed / divided

Chunk Creation

- How do we identify chunks?

Signature or Hash Digest

- A short, deterministic value generated from a set of data bytes
 - A document, chunk, block, or object
- Use for
 - Detecting equality (or likely equality)
 - Or, at least, detecting equivalence classes
 - Something must at least have the same signature to possibly be equal
- Hash should be short
 - Cannot be a 1:1 mapping from a large file (or chunk) to a short hash value

Penn ESE532 Fall 2017 -- DeHon

25

Example Hashes

- Sum up the bytes (or words) modulo some value
 - Variant: weighted sum
- XOR together the bits in some way
 - Variant: lots of different ways to shuffle bits for xor

Penn ESE532 Fall 2017 -- DeHon

26

Hashes and Chunk Creation

- Compute a hash on a window of values
 - Window: sequence of N-bytes
- Scan window over the input
- When hash has some special value (like 0)
 - Declare separate off a new chunk

Penn ESE532 Fall 2017 -- DeHon

27

Hashes as Chunk Cut Points

- What does this do?
- Guarantees that each chunk begins (or ends) at some fixed hash
- For a particular substring that matches the target hash
 - Always occurs at beginning (or end) of chunk
- If have a large body of repeated text
 - Will synchronize cuts at the same points based on the content

Penn ESE532 Fall 2017 -- DeHon

28

Chunk Size

- Assume hash is uniformly random
- The likelihood of each window having a particular value is the same
- So, if hash has a range of N, the probability of a particular window having the magic "cut" value is $1/N$
- ...making the average chunk size N
- So, we engineer chunk size by selecting the range of the hash we use
 - E.g. 12b hash for $2^{12} = 4\text{KB}$ chunks

Penn ESE532 Fall 2017 -- DeHon

29

Chunking Design

- Raises questions
 - How big should chunks be?
 - Apply maximum and minimum size beyond content definition?
 - How big should hash window be?
- Discuss
 - What forces drive larger chunks, smaller?
 - What forces drive larger windows, smaller?

Penn ESE532 Fall 2017 -- DeHon

30

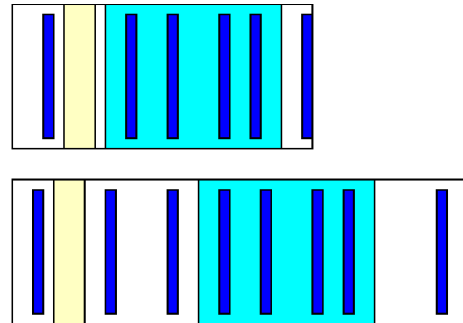
Example Text

- Beginning of repeated block of text.
- This stuff is has already been seen.
- But, we are only matching on something that has a hash of zero.
- Maybe this line has a hash of zero.
- But, our repeated text is before and after the magic window with the matched hash value.

Penn ESE532 Fall 2017 -- DeHon

31

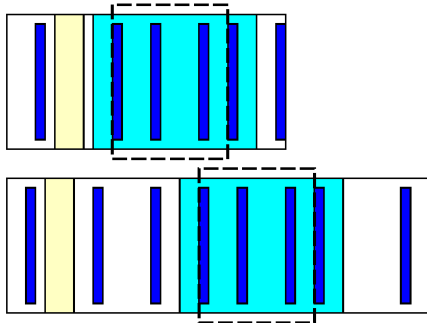
Example Data Stream



Penn ESE532 Fall 2017 -- DeHon

32

Example Data Stream



Penn ESE532 Fall 2017 -- DeHon

33

Rolling Hash

- A Windowed hash that can be computed incrementally
- $\text{Hash}(a[x+0], a[x+1], \dots, a[x+W-1]) = \text{Hash}(a[x-1], a[x+0], \dots, a[x+W-2]) - F(a[x-1]) + F(a[x+W-1])$
- i.e., hash computation is associative
- (+, - used abstractly here, could be in some other domain than modulo arithmetic)

Penn ESE532 Fall 2017 -- DeHon

34

Rabin Fingerprinting

- Particular scheme for *rolling hash* due to Michael Rabin based on polynomial over a finite field
- Commonly used for this chunking application

Penn ESE532 Fall 2017 -- DeHon

35

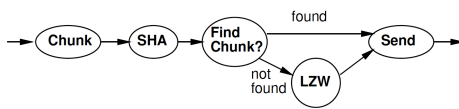
Content-Defined Chunking

- Compute rolling hash (Rabin Fingerprint) on input stream
- At points where hash value goes to 0, create a new chunk

Penn ESE532 Fall 2017 -- DeHon

36

Hashing Deduplication



Penn ESE532 Fall 2017 -- DeHon

37

Hashes for Equality

- We can also (separately) take the hash signature of an entire chunk
- The longer we make the hash, the lower the likelihood two *different* chunks will have the same hash
- If hash is perfectly uniform,
 - N-bit hash, two chunks have a 2^{-N} chance of having the same hash.

Penn ESE532 Fall 2017 -- DeHon

38

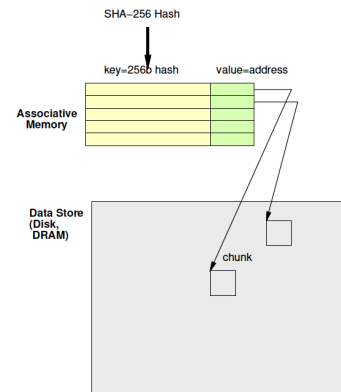
Deduplicate

- Compute chunk hash
- Use chunk hash to lookup known chunks
 - Data already have on disk
 - Data already sent to destination, so destination will know
- If lookup yields a chunk with same hash
 - Check if actually equal (maybe)
- If chunks equal
 - Send (or save) pointer to existing chunk

Penn ESE532 Fall 2017 -- DeHon

39

Deduplication Architecture



Penn ESE532 Fall 2017 -- DeHon

40

Associative Memory

- Maps from a key to a value
- Key not necessarily dense
 - Contrast simple RAM
- Talk about options to implement next week

Penn ESE532 Fall 2017 -- DeHon

41

Secure Hash

- We regularly use signatures to identify if a file has been tampered with
- Again, hashes are same, mean data might be the same
- For security, we would like additional property
 - not easy to make the anti-tamper signature match

Penn ESE532 Fall 2017 -- DeHon

42

Cryptographic Hash

- One-way functions
- Easy to compute the hash
- Hard to invert
 - Ideally, only way to get back to input data is by brute force
- Key: someone cannot change the content (add a backdoor to code) and then change some further to get hash signature to match original

Penn ESE532 Fall 2017 -- DeHon

43

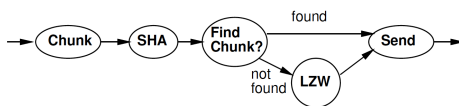
SHA-256

- Standard secure hash with a 256b hash digest signature
- Heavily analyzed
- Heavily used
 - TLS, SSL, PGP, Bitcoin, ...

Penn ESE532 Fall 2017 -- DeHon

44

LZW Compression



Penn ESE532 Fall 2017 -- DeHon

45

Preclass 3, 4, 5

- Message?
- Bits in unencoded (decoded) message?
- Bits for encoded message?

Penn ESE532 Fall 2017 -- DeHon

46

Idea

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

Penn ESE532 Fall 2017 -- DeHon

47

Encoding

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

Penn ESE532 Fall 2017 -- DeHon

48

Algorithm Concept

- While data to send
 - Find largest match in window of data sent
 - If length too small (length=1)
 - Send character
 - Else
 - Send $\langle x, y \rangle = \langle \text{match-pos}, \text{length} \rangle$
 - Add data encoded into sent window

Penn ESE532 Fall 2017 -- DeHon

49

Idea

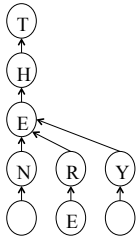
- Represent all strings as prefix tree
- Share prefix among substrings

Penn ESE532 Fall 2017 -- DeHon

50

Tree Example

- THEN AND THERE, THEY STOOD...



Penn ESE532 Fall 2017 -- DeHon

51

Tree Algorithm

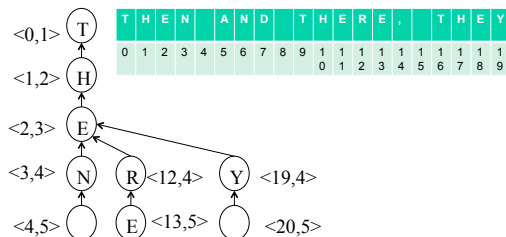
- Root for each character
- Follow tree according to input until no more match
- Send $\langle \text{name of last tree node} \rangle$
 - An $\langle x, y \rangle$ pair
- Extend tree with new character
- Start over with this character

Penn ESE532 Fall 2017 -- DeHon

52

Tree Example

- Label with $\langle \text{lastpos}, \text{len} \rangle$ pair
- THEN AND THERE, THEY STOOD...



Penn ESE532 Fall 2017 -- DeHon

53

Large Memory

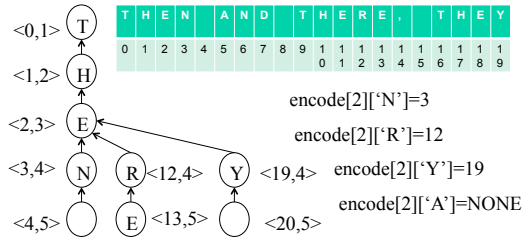
- `int encode[SIZE][256];`
- Name tree node by position in chunk
 - lastpos
- `c` is a character
- `Encode[lastpos][c]` holds the next tree node that extends tree node lastpos by `c`
 - Or NONE if there is no such tree node

Penn ESE532 Fall 2017 -- DeHon

54

Tree Example

- Label with <lastpos,len> pair
- THEN AND THERE, THEY STOOD...



Penn ESE532 Fall 2017 -- DeHon

55

Memory Tree Algorithm

```
curr – pointer into input chunk
// follow tree
y=0;
while(encode[x][input[curr]]!=NONE)
    x=encode[x][c]; y++;
if (y>0)
    send <x,y>
    send input[curr]
    encode[x][input[curr]]=curr
```

Penn ESE532 Fall 2017 -- DeHon

56

Complexity

- How much work per character to encode?

Penn ESE532 Fall 2017 -- DeHon

57

Compact Memory

- int encode[SIZE][256];
- How many entries in this table are not NONE?

Penn ESE532 Fall 2017 -- DeHon

58

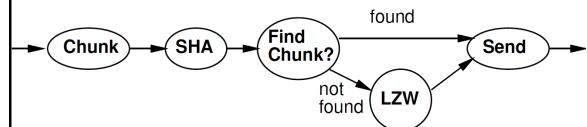
Compact Memory

- int encode[SIZE][256];
- Table is very sparse
- Store as associative memory
 - At most SIZE entries
- Look at how to implement associative memories next time

Penn ESE532 Fall 2017 -- DeHon

59

Project Task



Penn ESE532 Fall 2017 -- DeHon

60

Big Ideas

- Can reduce data size by identifying and reducing redundancy
- Can spend computation and data storage to reduce communication traffic

Admin

- HW7 due Friday
- Project assignment out
- Shuffling schedule a bit to deal with project needs
 - Monday → (near) associative memories
 - (...more shuffling to come...)
- First project milestone due next Friday
 - Including teaming