DRESS Code For The Storage Cloud

Distributed Replication-based

Exact Simple Storage

Salim El Rouayheb EECS Department University of California, Berkeley

Joint work with: Sameer Pawar Nima Noorshams Prof. Kannan Ramchandran

Talk Roadmap



Data Everywhere



- Increasing popularity of data-enabled devices and Hot spots, 4G, Wimax, Wi-fi etc, => "always online" infrastructure.
- Truly mobile users
- Want to be able to seamlessly access and share data anytime anywhere

Off To The Cloud...



- Want our files to "follow" us wherever we go
 - > Put them in the cloud
 - Get them from the cloud







Amazon Web Services

Amazon Simple Storage Service (Amazon S3)

Amazon S3 is storage for the Internet. It is designed to make web-scale computing easier for developers.









SmugMug 👸

Pricing

Region: US	- Standard	\$					
Storage (Designed for 99.999999999% Durability)		Reduced Redundancy Storage (Designed for 99.99% Durability)		Data Transfer**		Requests	
Tier	Pricing	Tier	Pricing	Tier	Pricing	Туре	Pricing
First 1 TB / month of Storage Used	\$0.140 per GB	First 1 TB / month of Storage Used	\$0.093 per GB	All data transfer in	\$0.100 per GB	PUT, COPY, POST, or LIST Requests	\$0.01 per 1,000 Requests

Wuala: P2P Cloud Storage

- Trade idle local disk space on your computer for online storage
- First 1GB is free



Cloud storage = local donated storage * percentage of online time

70 GB

100 GB

70%

Can We Trust the Cloud?

Main Challenge: Cloud is formed of unreliable and untrusted components

CLOUD COMPUTIN(

Business Communications

E-Commerce

E-Commerce Times > Enterprise IT > Software > Cloud Computi

Gmail Meltdown Casts Sha

Ga



By Richard Adhikari E-Commerce Times 02/28/11 12:53 PM PT

A Gmail disruption บร



Home Busi

Flickr

STORY HIGHLIGHTS

 Amazon.com was apparently unprepared for demand for Lady Gaga's new album

(Mashable) - Lady Gaga fans were delighted Monday to learn that they could download her new album, Born This Way, from Amazon for a mere \$0.99 - until, of course, technical difficulties set in.

 The online retailer is offering the new record for 99 cents

Downloads of the album are delayed, leaving folks unable to get the entire album immediately upon purchase. Amazon issued following statement

user's 4,000 photos

By Laurie Segall, staff reporter February 2, 2011: 3:15 PM ET

Lady Gaga deal overwhelms Amazon servers

Mashable By Brenna Ehrlich, Mashable May 24, 2011 6:06 a.m. EDT | Filed under: Web



Quote

Sn

Cloudy with a Chance of Failure



Component failures are the norm rather then exception.

- Nodes are unreliable by nature:
 - Hard disk failure
 - Network failure
 - Peer churning
 - Malicious nodes



Solution: store data redundantly

Multiple System Considerations



Different Codes for Different Folks



Reducing Bandwidth Overhead

- Repair BW can be reduced from 2 to 1.5 MB
- Helper nodes send random linear combinations: random linear network coding



Replacement node

Dimakis, Wu & Ramchandran '07

Fundamental Storage/Bandwidth Tradeoff



(Dimakis et al. '07)

Repair BW

We Want More: Exact Repair?



• Recent important progress: Rashmi et al. (Product-Matrix Codes), Suh & Ramchandran '10, Cadambe et al. '10

Storage/Bandwidth Tradeoff



Repair BW

Can we also reduce Disk I/O & Energy?

Bandwidth-efficiency comes with a price-tag.

- 1) Excessive data reading & computations
- 2) Read/write bandwidth of a hard disk is the bottleneck (100 MB/s vs. network speed ~1000 MB/s)

Can we have *exact* codes with minimum overheads of:

- 1) Repair bandwidth
- 2) Disk reads
- 3) Computations?



Contributions: Yes We Can

- DRESS codes:
 - ✓ Exact Repair
 - ✓ Minimum disk reads for repair
 - ✓ No data processing for repair: uncoded repair



- Surprise: no loss of bandwidth efficiency
- Bonus: well-suited for security applications
- System price: repair not as flexible w.r.t. who can help

Talk Roadmap



DRESS Code Example



Repairing DRESS Codes

(n,k)=(7,3) parity-check File 6 MB 1 2 3 4 5 6 7

- ✤ Repetition r=3→can tolerate up to 2 simultaneous failures
- Table based repair
- If 3 nodes fail?
 - Download complete file
 - Reconstruct the lost data
 - Rare bad even
 - ✓ Min repair bandwidth
 - ✓ Min disk I/O calls
 - ✓ Uncoded repair



Why "Traditional" Repetition Won't Work?

Max file size= 6MB: cut-set bound (Dimakis et al.)



(n,k)=(7,3) repetition factor r=3



Smart Packet Placement

How should we place the packet replicas? Goal: maximize file size





$$\frac{k(k-1)}{2} = \binom{k}{2}$$

Rule: Make sure that any two nodes have at most one packet in common

Where Did That Solution Come From?



Higher Order Projective Planes



m=4

 Projective planes are guaranteed to exist for any prime power order m.

Codes from Projective Planes



Theorem 1 (R. & Ramchandran Allerton'10) A projective plane of order m gives DRESS codes with $n=m^2+m+1$ nodes and repetition factor r=repair degree d = m+1.

Kirkman's Schoolgirls Problem

"Fifteen young ladies in a school walk out three abreast for

seven days in succession; it is required to arrange them daily so that no two shall walk twice abreast."

Kirkman, 1847

- Fifteen schoolgirls: A, B,..., O
- They always take their daily walks in groups of threes

	group 1	group 2	group 3	group 4	group 5
Sun	ABC	DEF	GHI	JKL	MNO
Mon	ADH	BEK	MOI	FLN	GJC
÷	?	?	?	?	?
Sat	?	?	?	?	?

• Arrange them such that any two schoolgirls walk together exactly once a week

Why Do We Care?



5 groups

- Think of schoolgirls as packets and groups as storage nodes
- Kirkman's solution gives an optimal DRESS code for a
 - ✓ System with n=35 nodes
 - ✓ Each node stores 3 packets
- What is the repetition factor here? answer=7

Not Your Average Puzzle...

- Kirkman's schoolgirls problem initiated a branch of combinatorics known now as *Design Theory*
- Deep connections to other branches in mathematics: Algebra, Geometry, matroid theory...
- Applications:
 - Coding Theory
 - Cryptography
 - Statistics



Steiner Systems

A Steiner system $S(t, \alpha, v)$ is a collection of points & lines such that:

- 1. There are v points in total
- 2. Each line contains exactly α points

3. There is exactly one line that contains any given t points





Fano plane = S(2,3,7) # pts on a line

Total # points

Codes from Steiner Systems

Theorem 2: (R. & Ramchandran Allerton'10)

A Steiner system $S(2, \alpha, v)$ gives a capacity-achieving DRESS code with parameters $r = \frac{v-1}{\alpha-1}$ $n = \frac{rv}{\alpha}$

using correspondence lines: nodes, points: packets.



Open problem I DRESS codes beyond Steiner systems?

Min Reads Capacity

Min Bandwidth Capacity

k



$$C = k \times d - \left(\right)$$



Open problem II

Capacity for minimum reads (and lookup repair)?

Example





- Users contacting any 3 nodes will observe either 9 or 7 packets
- This code guarantees a "rate" = 7 to any user
- Capacity when repairing from any d nodes = 6

We Want More

- Explicit constructions for Steiner systems exist for small values of r=3,4,5
- Existential guarantees for large number of nodes [Wilson 75]
- Want more flexibility
- Want to "grow". Nested property?
- Solution: randomized constructions



P2P network

Lazy Man's Codes



- Fill the bins randomly by throwing 3 independent colors in each
- User takes a hit on the file size if he contacts bins that have two many colors in common
- Replication factor is random and can be less than 3
- Hope the bad events are rare

It Pays To Be Lazy



Theorem 3:(Concentration)

Let F be the number of colors observed by a user contacting k nodes then

$$P(|F - \bar{F}| \ge B) \le 2 \exp(-\frac{B^2}{2kd})$$

with $\bar{F} = \theta \left(1 - (1 - \frac{1}{\theta})^{kd}\right), \theta = nd/r.$

Content Distribution with Helpers



- Examples of helper include:
 - ✓ idle Internet PC users: Thunder, PPStream, PPLive, Wuala etc...
 - ✓ set-top boxes, home gateways: Nano data centers
 - ✓ ISP P2P infrastructure nodes: Comcast
- Goal: minimize server load or maximize a utility function

Talk Roadmap



Security, Security, Security

- How to guarantee data privacy and integrity?
- What are the fundamental limits on cloud security?





- Computationally unbounded intruder
- Two classes of intruders:
 - Passive: eavesdropper
 - Active: Malicious attacks

Malicious Byzantine Nodes

- Hard problem related to open problems in network coding. (Kosut, Tong & Tse '10)
- Non-linear coding, secure capacity still an open problem...
- "Node vs. Edge Curse" known in the network coding literature
- How about the cloud?



Why Is It a Hard Problem?





- Even a single malicious node can contaminate the whole system
- Although controlling 3 bits, the malicious node can introduce 5 errors
- Can we at least store 1 bit securely?

No Bit Left Behind



- Repeat the 1 bit information on all the nodes
- 5 errors out of 9: simple majority decoding won't work.
- Idea: not any 5 errors can occur
- DRESS code forces errors to have certain patterns

Tetris Decoder



- Decoder checks the error pattern and decode
- It will always decodes to the correct bit
- Bonus: can catch the malicious nodes

Omniscient Adversary: Main Result



- + Linear outer code. Internal nodes in the cloud operate as usual.
- Decoder: exponential time in number of malicious nodes

Upen

Conclusion

- Cloud storage: guarantee reliability using unreliable nodes!
- What codes should we use? storage, bandwidth, Disk I/O, computational complexity, latency, energy, security...
- DRESS codes: uncoded repair and efficient codes
- Deterministic constructions from Proj. Spaces & Steiner sys.
- Randomized constructions using balls & bins
- Separation result for achieving security

 Open problems: Code constructions for more system parameters, "Smart" random constructions for better concentration, general secure capacity...