Scalable Data Management and Storage on the Cloud:

- State of the Art and Emerging Trends

Presented by

Anwitaman Datta & Frédérique Oggier NTU Singapore

ICDCN 2012, Hong Kong

Who are we?





Anwitaman Datta School of Computer Engineering NTU Singapore

Frédérique Oggier School of Physical and Mathematical Sciences NTU Singapore

Who are we? Disclaimer

A note from the trenches: "You know you have a large storage system when you get paged at 1 AM because you only have a few petabytes of storage left." – from Andrew Fikes' (Principal Engineer, Google) faculty summit talk ` Storage Architecture and Challenges `, 2010.

... and some ask/say: why do you care about efficient storage space utilization, it is so cheap ...

We two never get such calls!!





- Distributed storage systems
 - Scale how much? Scale how?
 - From P2P to Data centers
- Data storage & management
 - NoSQL systems
 - Detailed examples: Dynamo, GFS, ...
- Erasure coding based efficient fault-tolerant storage
 - Background: RAID, (traditional) erasure codes
 - Codes tailor-made for storage
 - Detailed examples: Regenerating codes, Self-repairing codes, ...
- Appendix: MapReduce

The "aging" state-of-the-art

> n emerging trend ...



• Rather big ...



© 2012 A. Datta & F. Oggier, NTU Singapore



- June 2011 **EMC²** study
 - world's data is *more than doubling* every 2 years
 - faster than Moore's Law
 - 1.8 *zettabytes* of data to be created in 2011

Zetta: 10²¹

Zettabyte: If you stored all of this data on DVDs, the stack would reach from the Earth to the moon and back.



* http://www.emc.com/about/news/press/2011/20110628-01.htm

Introduction NoSQL Codes for storage Appendix The data deluge: Some numbers

- Facebook "currently" (in 2010) stores over 260 billion images, which translates to over 20 petabytes of data. Users upload one billion new photos (60 terabytes) each week and Facebook serves over one million images per second at peak. [Beaver et al. in "Haystack" paper]
- On "Saturday", photo number four billion was uploaded to photo sharing site Flickr. This comes *just five and a half months after the 3 billionth* and nearly 18 months after photo number two billion. Mashable (13th October 2009)

[http://mashable.com/2009/10/12/flickr-4-billion/]



Introduction NoSQL Codes for storage Appendix The data deluge: Some numbers

Caffeine lets us index web pages on an enormous scale. In fact, every second Caffeine processes hundreds of thousands of pages in parallel. If this were a pile of paper it would grow three miles taller every second. Caffeine takes up nearly 100 million gigabytes of storage in one database and adds new information at a rate of hundreds of thousands of gigabytes per day. You would need 625,000 of the largest iPods to store that much information; if these were stacked end-to-end they would go for more than 40 miles.



http://googleblog.blogspot.com/2010/06/o ur-new-search-index-caffeine.html

Introduction NoSQL Codes for storage Appendix The data deluge: Some numbers



Introduction NoSQL Codes for storage Appendix Scale how?

To scale vertically (or **Scale up**) means to add resources to a single node in a system*

To scale horizontally (or **scale out**) means to add more nodes to a system, such as adding a new computer to a distributed software application*



Scale out

© 2012 A. Datta & F. Oggier, NTU Singapore

* Definitions from Wikipedia

Scale up

Introduction NoSQL Codes for storage Appendix Distribution is essential Appendix

- Scaling up
 - May just not even be feasible
 - Even if feasible, it will be very expensive
 - What happens when "the" machine fails?
- Scaling out => distributed storage
 - Distribution => added complexity and vulnerabilities
 - latency, consistency, faults, ...
 - CAP theorem
 - Consistency, Availability, Partition tolerance choose any two
 - but, not distributing is not a choice!

Introduction NoSQL Codes for storage Appendix

- Different flavors of functionality
 - File level data storage: Network-attached storage
 - Block level data storage: Storage area network
 - Distributed databases
 - Caching/CDNs
- Different flavors of architecture
 - Centralized, decentralized/peer-to-peer, hybrid ...
- Different flavors of interfaces
 - depends on/constrains functionality



- Popularized by distributed hash tables
 - Basic operations: Get, Put, ...
 - Easy to distribute (partition the key space)





• DHT itself as the storage layer • DHT as a directory service



Introduction NoSQL Codes for storage Appendix Wuala's 3-tier architecture

- Complete disentanglement of indexing and storage
- Many (encoded) fragments per object
 - Suitable for sharing very *large* but *static* files
 - Parallel download
- Piggy-backed, large DHT routing states
 - So very few hops needed, gives high through-put

Source: Google tech talk on Wuala: <u>http://www.youtube.com/watch?v=3xKZ4KGkQY8</u>

© 2012 A. Datta & F. Oggier, NTU Singapore

Superpeer

DHT

Wuala's dedicated

storage resource as

a fall-back option

sers

Storage peers



- We are not here to talk about P2P systems, the title mentioned "on the cloud"!
 - **Don't panic**! The fundamental things apply (even as time goes by ...)



© 2012 A. Datta & F. Oggier, NTU Singapore

Introduction NoSQL Codes for storage Appendix But, what is the cloud?

- At least, we can all agree
 - Cloud is something "big" and happening!
 - It's all of these ...
 - ... and some more!



Systems of elephantine proportion needed



Introduction NoSQL Codes for storage Appendix NIST definition for cloud computing

• Cloud computing is a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction

Introduction NoSQL Codes for storage Appendix Two sides of the cloud coin

- Outside view
 - A "single/exclusive" entity
 - Access through a "demilitarized zone" ...
 - API based
 - Agnostic to multi-tenancy
 - Infinite/elastic resources
 - Pay per use, on-demand, ...
 - Browser based access (often)
 - Anytime, anywhere, any device ...



Introduction NoSQL Codes for storage Appendix Two sides of the cloud coin

- Inside view
 - Pool of resources



- In flux: New compute units joining, old ones retiring
- Self-*: Load-balancing, fault-tolerance, auto-configuration, ...
- Multi-tenancy
 - Virtualization, transparent migration, ...
- Distributed file system, data-management, data processing
 - Google's GFS, Amazon's Dynamo, Facebook's Cassandra, Yahoo!'s Pnuts
 - Map Reduce/Hadoop, Pig, Chubby, ...



Reliable storage service

Distributed Physical Infrastructure: Storage & Compute Nodes, Interconnect, ...

Disclaimer: This "stack" is a personal "view", and is not universal © 2012 A. Datta & F. Oggier, NTU Singapore



Source of topology: http://www.cisco.com/en/US/docs/solutions/Enterprise/Data_Center/vmware/VMware.html Note: More recently, the trend is to deploy "Fat Tree" interconnects 22 © 2012 A. Datta & F. Oggier, NTU Singapore

Data center design evolution



Slide courtesy **Roger Barga** (Microsoft) from his P2P 2009 Keynote talk

Introduction NoSQL Codes for storage Appendix Amazon's AWS: Availability zones

Amazon Web Services Overview of Security Processes

May 2011



Build physical infrastructure to store immense data ...



Still, how to manage so much data and distribution???





NO

geek & poke

DOESN

Leverage the NoSQL boom

- Neither always needed
- Nor scales (out)
- Various workload specific custom storage & data management solutions

 NoSQL





- Document store
- Graph
- Key-value store
- Multivalue databases
- Object database
- Tabular
- Tuple store



Dynamo: Amazon's Highly Available Key-Value Store Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, Gunavardhan Kakulapati, Avinash Lakshman, Alex Pilchin, Swaminathan Sivasubramanian, Peter Vosshall and Werner Vogels SOSP 2007

Introduction NoSQL Codes for storage Appendix Fast and furious

- Latency sensitivity
 - Shopping Cart Service
 - 10s of millions of requests per day
 - Millions of checkouts each day
 - Hundreds of thousands of concurrent activities
- Stringent SLAs for each service
 - 99.9th percentile < 300ms
 - Mean/std. dev. inadequate
- One page request
 - 100s of services
 - Multiple service dependencies



Introduction NoSQL Codes for storage Appendix O The show must go on O

- Extremely high availability
 - Down time => lost business
- "customers should be able to view and add items to their shopping cart even if disks are failing, network routes are flapping, or data centers are being destroyed by tornados"
- an infrastructure comprised of millions of components
 - tens of thousands of servers located across many data centers world-wide
 - a small but significant number of server and network components that are failing at any given time



The dreaded parthyuake - enhanced tornado being hit by a radioactive meteor event.



- Infrastructure comprised of millions of components
 - tens of thousands of servers located across many data centers world-wide
 - a small but significant number of server and network components that are failing at any given time
- Redundancy needed for fault tolerance
 - Trade-off: Availability is more important than consistency
- An "always writeable" data store
 - Conflict resolution complexity at "reads"
 - Unlike most traditional data stores
 - Can be handled at clients (using application logic)
 - Data store provides some default fall-back option "last write wins"

Introduction NoSQL Codes for storage Appendix KISS: Keep it simple, Stupid!

- Both stateless and stateful (needing persistent storage) services
- Most services access/store data using "primary keys"
 - No need for complex queries
 - No operation span multiple data items
 - Relatively small objects (<1MB)
- RDBMS is an overkill!
 - Also, difficult/impossible to scale-out
 - Needs much more expensive hardware/administration



- Incremental scalability
- Symmetric
 - Simpler system provisioning & maintenance
- Decentralized
 - Self-*, no single point of failure, ...
- Heterogeneity friendly



36
Introduction NoSQL Codes for storage Appendix O Dynamo architecture: Interface

- **get**(key)
 - Locate object replicas associated with the key
 - Return object/list of objects, with context
- put(key, context, object)
- Context encodes system metadata
 - E.g., version
- MD5[Caller key] \rightarrow 128 bit identifier

Over the second second

Codes for storage

• Zero-hop DHT

Introduction

- Consistent hashing based data partitioning
- All nodes know all other nodes

NoSQL

- Multiple "tokens" per node
 - Virtual node instances
 - Easy to handle heterogeneity
 - Node departure/arrival
 - Load is distributed
- Replication (configurable)
 - For any key: *preference list*
 - Ensure distinct physical nodes
 - Across multiple data centers



Appendix

Introduction NoSQL Codes for storage Appendix O Dynamo architecture: Data versioning



Introduction NoSQL Codes for storage Appendix O Dynamo architecture: Executing get()/put()

- Symmetry
 - Client can send get/put requests for any key to any Dynamo node
- Sloppy-quorum
 - First N healthy nodes in the preference list
 - R+W > N quorum \circ
- Upon receiving put() request
 - coordinator generates vector clock, writes locally
 - sends to N highest-ranked reachable nodes
 - W-1 acks implies a successful write
- Upon receiving get() request
 - coordinator requests for all existing versions to N highest-ranked nodes

object

availability, consistency

æ

durability

trade-offs

 waits for R responses, gathers all versions, and sends all causally unrelated versions

Introduction NoSQL Codes for storage Appendix O Dynamo architecture: Other aspects

- Hinted handoff: Always writeable
- Gossip based
 - failure detection
 - membership information propagation
- Buffered writes:
 - Writes stored in main memory buffer, periodically written to storage
 - Improves latency, risks durability

Introduction NoSQL Codes for storage Appendix O Dynamo in December 2006



Figure 4: Average and 99.9 percentiles of latencies for read and write requests during our peak request season of December 2006. The intervals between consecutive ticks in the x-axis correspond to 12 hours. Latencies follow a diurnal pattern similar to the request rate and 99.9 percentile latencies are an order of magnitude higher than averages

Figure 5: Comparison of performance of 99.9th percentile latencies for buffered vs. non-buffered writes over a period of 24 hours. The intervals between consecutive ticks in the x-axis correspond to one hour.

Introduction NoSQL Codes for storage Appendix

Dynamo: Summary

Table 1: Summary of techniques used in Dynamo and their advantages.

Problem	Technique	Advantage
Partitioning	Consistent Hashing	Incremental Scalability
High Availability for writes	Vector clocks with reconciliation during reads	Version size is decoupled from update rates.
Handling temporary failures	Sloppy Quorum and hinted handoff	Provides high availability and durability guarantee when some of the replicas are not available.
Recovering from permanent failures	Anti-entropy using Merkle trees	Synchronizes divergent replicas in the background.
Membership and failure detection	Gossip-based membership protocol and failure detection.	Preserves symmetry and avoids having a centralized registry for storing membership and node liveness information.

Introduction NoSQL Codes for storage Appendix

- Dynamo provides a bare-bone storage service
 - using the key-value data structure
 - how about something more sophisticated, i.e., a file system?

Dynamo@ SOSP'07



GFS@ SOSP'03





Introduction NoSQL Codes for storage Appendix The big fat Gogle computations

- Throughput for bulk data processing in batches
 - High sustained bandwidth is more important
 - Than low latency
- Few million files
 - Mostly > 100 MB
 - Multi-GB files very common
 - Small files must be supported
 - But no need for optimization



Introduction NoSQL Codes for storage Appendix O Workload is primarily two kinds of reads

- Large streaming reads
 - $100s \text{ of KBs}, > 1 \text{MB} \dots$
- Small random reads

- Successive operations from a client *read contiguous region* of a file
- Applications can sort small reads to advance steadily through a file
 - To avoid going back & forth



Introduction NoSQL Codes for storage Appendix Many large sequential writes/appends ...

- Mainly when a file is being created
 - Once written, it is seldom modified
- Small writes at arbitrary position needs to be supported
 - But do not have to be efficient!

- Multiple clients that concurrently append
 - Producer-consumer queues, many-way merging
 - e.g., MapReduce operations
 - Needs atomicity



Introduction NoSQL Codes for storage Appendix KISS: Keep it simple, Stupid!

- Basic operations
 - E.g., Not standard (POSIX) compliant
 - create, delete, open, close, read, write files
 - snapshot, record append



- Simplifies design
 - Can carry out sophisticated data placement and replication decisions using global knowledge
- What about fault-tolerance, bottleneck, ...?
- Multiple chunk servers & multiple clients
 - Could be run on same machines

• Note: HDFS (Hadoop[™] Distributed File System) follows a similar architecture

Introduction NoSQL Codes for storage Appendix Files \rightarrow Fixed-sized chunks Maximize: • Data reliability Global 64 bit uid for each chunk - e.g., not all replicas on the same rack - Assigned by master at chunk creation Network utilization Chunkservers store chunks on local disksaggregate BW of racks * during reads - As Linux files * but multi-rack write - Chunks replicated for reliability traffic • Default: 3 replicas - Collosus (GFS v 2) uses erasure codes: 1.5x Clients interact with master for metadata - Interacts with chunkservers for actual data manipulation • Caching: Client "cache" metadata Chunkservers have automatic Linux caching • No other caching at either clients/chunkservers: - Neither meaningful nor feasible ... (for the involved workloads)



Introduction NoSQL Codes for storage Appendix Chunk size and system scalability

• 64MB

+ Reduces clients' need to interact with master

- Many operations are contiguous/sequential on a file
- involves same chunkserver
- Client can cache chunk locations for even a multi-TB working set
- + Client likely to carry out more operations on same chunk
 - Since chunks are large
 - Amortizes network connection costs (persistent TCP connection)
- + Reduces the size of metadata stored on the master
 - So that it *fits in memory* (<64 bytes metadata per chunk at Master)
 - Significant performance boost!
- Possibility of hot spot
- Fragmentation/poor space utilization
 - Most files are large, so only last chunk partially filled
 - Further mitigated using lazy space allocation



Design choices: *Reduce the master's involvement* as much as possible

- Cache chunkserver meta-info at clients
- Obtain/send extra meta-info in advance (amortizing cost of communication)
- "Heavy" tasks directly between clients/chunkservers



- GFS namespace
 - No per directory data structures
 - Lookup table: Mapping full pathnames to metadata
 - Use prefix compression for storage in memory
- Each node in namespace tree has associated read/write lock
 - To manipulate /d1/d2/.../dn/leaf
 - Obtain read locks for /d1, /d1/d2, ... /d1/d2/.../dn
 - Obtain read or write lock for /d1/d2/.../dn/leaf
- Notes:
 - File creation does not require a write lock on the parent directory because there is no "directory", or inode-like, data structure to be protected from modification, thus allowing multiple simultaneous mutations in a directory
 - In that order (consistent total order) to avoid deadlocks
 - Allocated lazily, removed once not in use

Leases for mutations

NoSQL

Introduction

- One chunk replica is given a lease to act as the "primary" replica
 - Typically for 60 seconds
 - Renewable repeatedly if mutation continues
 - over HeartBeat messages
 - Master can revoke lease (e.g., for snapshot)
 - Primary determines a serial order for all mutations in a chunk
 - Other replicas follow this order
- Global mutation order is defined by
 - The lease grant order
 - Serial order inside a lease

Codes for storage

Appendix



Data is pipelined linearly:

- disentangled from control messages
- optimized for network topology
- leveraging on full-duplex links
- providing1MB Tx in ~80ms+



- Snapshot
 - Using copy on write
- Atomic record appends
 - Special/different from writes, since Google had append heavy workload
- Consistency model
 - Complicated semantics
 - both for defining consistency as well as replicas
- Stale replica detection, data integrity (using check sums)
- Garbage collection
- replica creation, re-creation, rebalancing

Introduction NoSQL Codes for storage Appendix In the meanwhile ... Google grew

- Mix of applications ...
 - Gmail, docs, buzz, ...
 - Google earth
 - Youtube ...

- Latency sensitivity
- Diverse "sizes" of data
- Both read/write intensive



Introduction NoSQL Codes for storage Appendix O BigTable: Support for structured data



- Data organized into three dimensions
 - rows, columns, timestamps
 - Not a full fledged RDMBS
- It is an "application" on top of GFS
 - but, is part of the "infrastrcuture"
 - For other applications
 - relies on Paxos algo. based Chubby
 - A highly available and persistent lock service



Introduction NoSQL Codes for storage Appendix Google's next generation FS: Collossus Details are "sketchy"

- Infrastructure for diverse mix of workload ...
- Cluster-level file system
- Data typically written using Reed-Solomon (1.5x) erasure codes (ECs)*
- Support for smaller "chunks"
 - More flexibility for diverse application
- Distribution of Master functionality
- Why ECs? Can we do better than the state of the art?
 - Second half!

* With Microsoft's Azure also using ECs, they are now pretty much a part of mainstream deployment







2012 Data Center Expansion Plans Cloud Computing Prompts Worldwide Expansion



Introduction NoSQL Codes for storage Appendix O Failure Is Inevitable

- But, failure of the system is not an option!
 - Failure is the pillar of rivals' success ...
- Solution: Redundancy & Distribution





NoSQL Codes for storage Appendix Introduction Is the Danger Real? Yes

Online Backup Company Carbonite Loses Customers' Data, Blames And Sues Suppliers (Updated)

Robin Wauters Mar 23, 2009

TechCrunch

Hotmail Data Loss Reveals Cloud Trust Issues **PCWorld**

By Keir Thomas, PCWorld

Hot on the heels of what was possibly the first major cloud data leak a few weeks ago, as the new year got underway Microsoft followed up by appearing to wipe the e-mails of a significant number of Hotmail users.

Gmail outage passes 24 hours for some (updated)

By Seth Weintraub February 28, 2011: 1:01 PM ET

Amazon's Cloud Crash Disaster Permanently Destroyed Many **Customers' Data**

Amazon is currently experiencing a degradation. They are working on it. We are still waiting on them to get to our volumes. Sorry,

Henry Blodget | Apr. 28, 2011, 7:10 AM



Introduction NoSQL Codes for storage Appendix O Data Center Fault-Tolerance

- Faults are *omnipresent*
 - Hardware, network, software, human, misconfiguration, ...
- *Cascade of failures* in interdependent networks
 - Power failure => Network switches stop working
 - Network failure => Control system for power system ineffective





- Data from Los Alamos National Laboratory (DSN 2006), gathered over 9 years, 4750 machines, 24101 CPUs.
- Distribution of failures:
 - Hardware 60%
 - Software 20%
 - Network/Environment/Humans 5%
- Failures occurred between once a day to once a month.

Redundancy Based Fault Tolerance

Codes for storage

Replicate data

Introduction

- e.g., 3 or more copies
- In nodes on different racks

NoSQL

• Can deal with switch failures



Appendix

Power back-up using battery between racks (Google)



Redundancy Based Fault Tolerance

NoSQL Codes for storage Appendix

- Using "independent" physical infrastructure
 - Over different availability zones (Amazon AZ)
 - How independent are components in a complex network?
 - Over multiple geographical regions



Note: The recent (April 2011) AWS outage was the first *region-wide failure*

© 2012 A. Datta & F. Oggier, NTU Singapore

Introduction

Introduction NoSQL Codes for storage Appendix

- Physical
- Virtual resource
- Availability zone
- Region
- Cloud



From: http://broadcast.oreilly.com/2011/04/the-aws-outage-the-clouds-shining-moment.html



- Failure is not an option, but ...
 - ... are the overheads acceptable?



Reducing the Overheads of Redundancy

Appendix

Codes for storage

• Erasure codes

Introduction

- Much lower storage overhead
- High level of fault-tolerance

NoSOL

- In contrast to replication or RAID based systems
- Has the potential to significantly improve the "bottomline"
 - Can it however match the performance needs?
 - An open question*

Does erasure coding have a role to play in my data center? Zhe Zhang, Amey Deshpande, Xiaosong Ma, Eno Thereska, Dushyanth Narayanan MSR TR 2010

* Note: Both Google's new DFS Collossus, as well as Microsoft's Azure now use ECs

Introduction NoSQL Codes for storage Appendix O Erasure Codes (ECs)

- An (n,k) erasure code = a map that takes as input k blocks and outputs n blocks, thus introducing n-k blocks of redundancy.
- 3 way replication is a (3,1) erasure code!



• An erasure code such that the k original blocks can be recreated out of any k encoded blocks is called MDS (maximum distance separable).
Appendix

Reed-Solomon Codes

NoSQL

Introduction

(named after Irving S. Reed and Gustave Solomon)



- Reed-Solomon Codes are well-known erasure codes.
- Encoding of (o_1, \dots, o_k) is done by polynomial evaluation:

Codes for storage

$$p(X) = \sum_{i=0}^{k-1} p_i X^i, \ p_i = o_{i+1}.$$

• The encoding blocks are then $p(\alpha_1), \dots, p(\alpha_n)$.





- Originally designed for communication
 - EC(n,k)







- Storage efficiency: k/n
- Access: Find k good blocks

Blk • Assumption: Peer failure is i.i.d. with failure probability f

For f=0.1

its **10-3**

replica

its ~

۶lk

BV

For f=0.1

3 of 9 code

Blk

Blk

object

replica

object

replica

Blk



• Repairs are expensive!

© 2012 A. Datta & F. Oggier, NTU Singapore



• Erasure codes tailor-made for distributed networked storage.



Introduction NoSQL Codes for storage Appendix O What are Tailored-Made Codes?

Desired code properties include:

- Low repair bandwidth
- Low storage overhead
- Good fault tolerance



But also:

- Repair time
- I/O



(redundant array of independent disks; originally redundant array of inexpensive disks)

- RAID 0 = store data across multiple drives so that the disk head can read more data in a single move (striping)
- RAID 1 = replication is introduced (mirroring)
- RAID 2,3,4 ,5 = parity bit , sum of all bits across one drive (Hamming parity code)





Pyramid & Hierarchical Codes

NoSQL Codes for storage

- If `small' number of faults
 - Communication restricted within the `hierarchy' suffice

Appendix

Pros

Cons

- Progressively go higher-up for larger number of faults
- Isolated faults can be repaired independently
- Naturally maps to hierarchical data-center design?
- Asymmetry

Introduction

- Different encoded blocks have different importance
- Difficult to analyze
- Complex algorithm (for decoding/repair) and system design

Pyramid Codes: Flexible Schemes to Trade Space for Access Efficiency in Reliable Data Storage Systems

Cheng Huang, Minghua Chen, and Jin Li

NCA 2007

Another (essentially identical, but independent) proposal: **Hierarchical Codes**: How to Make Erasure Codes Attractive for Peer-to-Peer Storage Systems *from Eurecom*



• Network information flow based arguments to determine "optimal" trade-off of storage/repair-bandwidth





- Example code (w/ Functional Repair)
 - Based on random linear network coding





- Example code (w/ Exact Repair)
 - Construction for one failure, contacting all live nodes for repair





Introduction NoSQL Codes for storage Appendix

(Collaborative) Regenerating Codes

- Pros
 - Network information flow analysis determines optimal (w.r.to repair bandwidth)
 - Catch: Subject to MDS property of code (more on this, later)
 - Some proposed codes apply network coding on top of ECs
 - Inherit the properties for EC for de/coding
- Cons
 - Codes for only specific points on trade-off curve
 - Information flow analysis itself does not suggest any code
 - Restrictive
 - Some proposed codes can carry out repair only for one fault
 - Needs to contact all live nodes for repair to be optimal
 - Not simple: Algorithmic as well as system design

Introduction NoSQL Codes for storage Appendix

Design Space For Cheaper Repairs

- What is the best one can do (w.r.to repairs)?
 - Minimize bandwidth usage per repair
 - Regenerating codes
 - Minimize number of live nodes used per repair
 - Self-repairing codes

Introduction NoSQL Codes for storage Appendix O Self-repairing Codes (SRC)

- Self-repairing codes are (n, k) codes s.t.
 - encoded fragments can be repaired directly from other subsets of encoded fragments.
 - a fragment can be repaired from a fixed number of encoded fragments, independently of which specific blocks are missing
 - Analogous to erasure codes supporting reconstruction using any n k losses, independently of which
 - number of live nodes contacted for repair is minimized.



Self-repairing Homomorphic Codes for Distributed Storage Systems Frédérique Oggier and Anwitaman Datta Infocom 2011





• There is at least one pair to repair a node, for up to (n-1)/2 simultaneous failures (Parallel & fast repair of multiple faults)



One Example: Reconstruction

NoSQL Codes for storage

node	basis vectors	data stored
N_1	$v_1 = (1000), \ v_2 = (0110)$	$\{o_1, o_2 + o_3\}$
N_2	$v_3 = (0100), \ v_4 = (0011)$	$\{o_2, o_3 + o_4\}$
N ₃	$v_5 = (0010), \ v_6 = (1101)$	$\{o_3, o_1 + o_2 + o_4\}$
N_4	$v_7 = (0001), \ v_8 = (1010)$	$\{o_4, o_1 + o_3\}$
N_5	$v_9 = (1100), \ v_{10} = (0101)$	$\{o_1 + o_2, o_2 + o_4\}$

Appendix



Reconstruction, say using N_3 , N_4 and N_5

© 2012 A. Datta & F. Oggier, NTU Singapore

Introduction

Maximum Distance Separable (MDS)?

- SRC is not MDS (and can not be!)
 - Does it matter?

Introduction

- Not much
- In practice, access will be "planned" ...
- PSRC needs less bandwidth than `optimal' RGC!



© 2012 A. Datta & F. Oggier, NTU Singapore



Practical properties

- (Current) SRCs are not systematic
 - PSRC is "like systematic"
 - Need to contact more nodes (than k)
 - To obtain systematic `pieces'
 - Same total bandwidth usage
 - Parallel download for access can even be an `advantage'
 - `mixed' strategies for access, i.e. get some systematic pieces, and some others ...
 - Power saving (by switching off nodes) strategies possible
- Coding/decoding in PSRC are both using
 - XOR operations only

node	basis vectors	data stored
N_1	$v_1 = (1000), \ v_2 = (0110)$	$\{o_1, o_2 + o_3\}$
N_2	$v_3 = (0100), \ v_4 = (0011)$	$\{o_2, o_3 + o_4\}$
N ₃	$v_5 = (0010), \ v_6 = (1101)$	$\{o_3, o_1 + o_2 + o_4\}$
N_4	$v_7 = (0001), \ v_8 = (1010)$	$\{o_4, o_1 + o_3\}$
N_5	$v_9 = (1100), \ v_{10} = (0101)$	$\{o_1 + o_2, o_2 + o_4\}$





- Interested to
 - Follow:

http://sands.sce.ntu.edu.sg/CodingForNetworkedStorage/

- Get involved: {anwitaman,frederique}@ntu.edu.sg







- Introduced by Google in 2004
 - Also available as open source Hadoop (w/ HDFS, etc)



Distributed Processing of Big Data

Systems of elephantine proportion needed



- Map: Transforms input data to intermediate (key, value) pair
- Reduce: Transforms all values for given key to final output





• Canonical example: Count



© 2012 A. Datta & F. Oggier, NTU Singapore





© 2012 A. Datta & F. Oggier, NTU Singapore