

Coding for Distributed Storage

Alex Dimakis (UT Austin)

Joint work with

Mahesh Sathiamoorthy (Google)

Megasthenis Asteris, Dimitris Papailipoulos, Karthik Shanmugam, Sriram Vishwanath, Ankit Rawat (UT Austin)

Overview

- How distributed file systems work
- Three repair metrics
- Part 1: Regenerating Codes
- Part 2: Locally Repairable Codes
- Part 3: Availability of Codes
- Open problems

current hadoop architecture

file 1



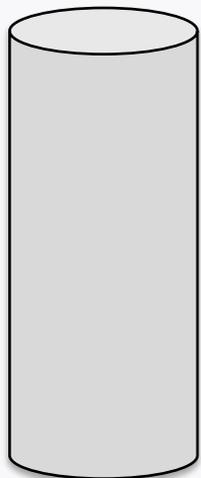
file 2



file 3



NameNode



DataNode 1



DataNode 2



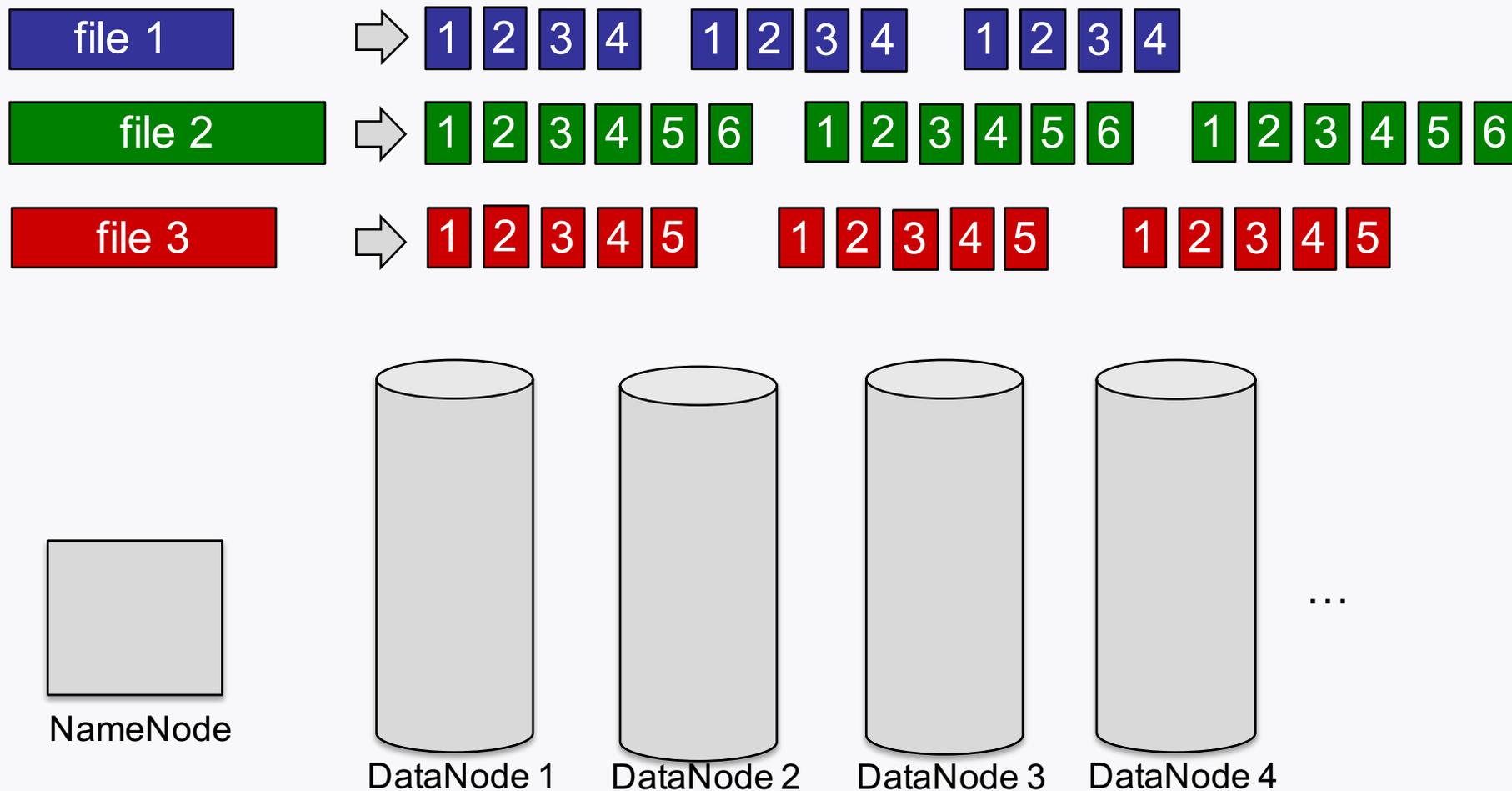
DataNode 3



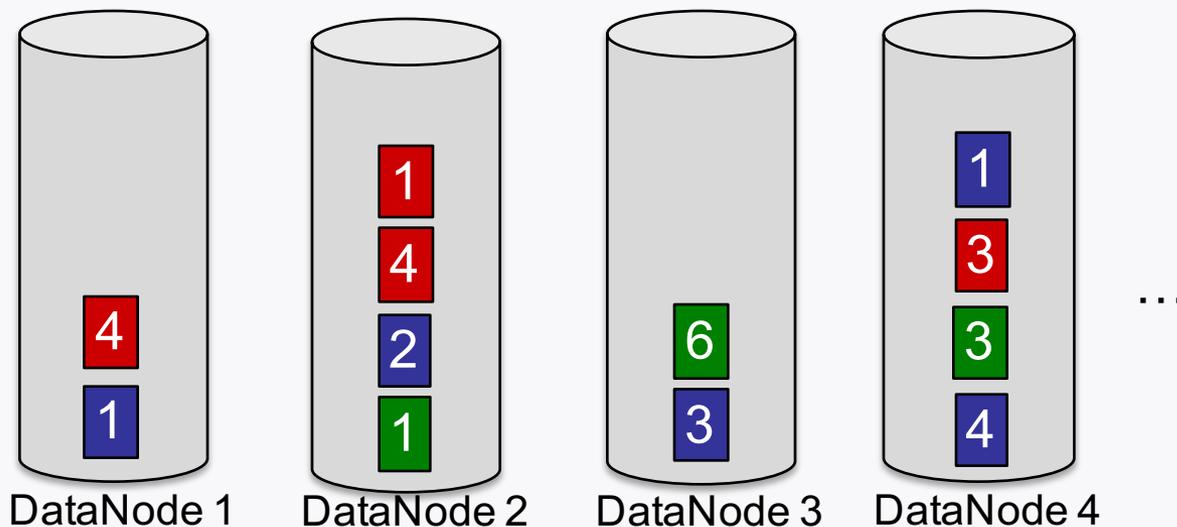
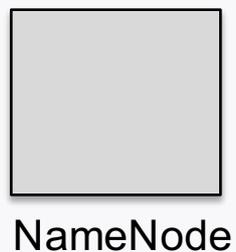
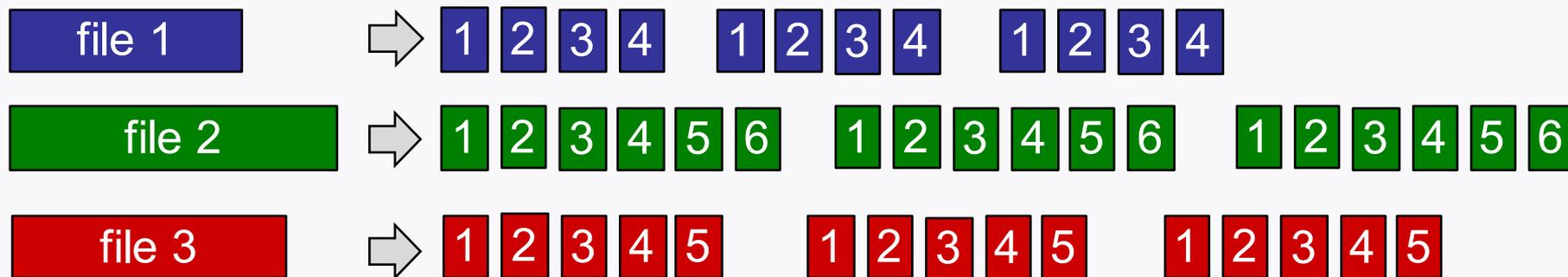
DataNode 4

...

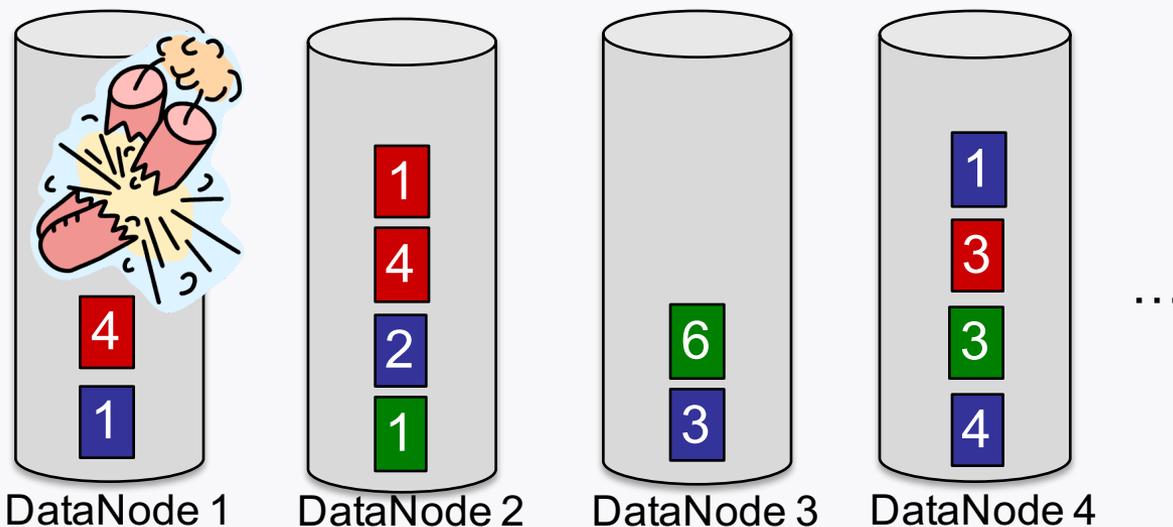
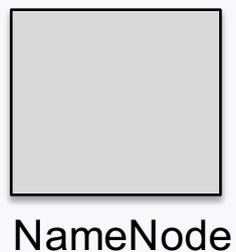
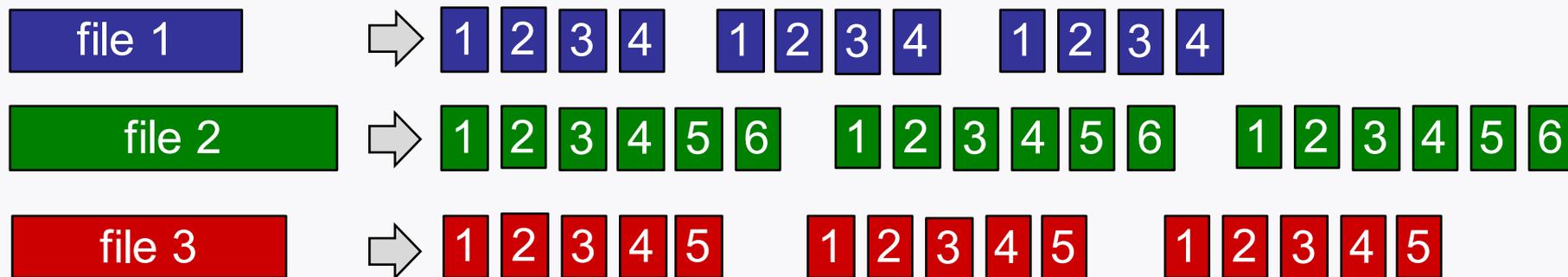
current hadoop architecture



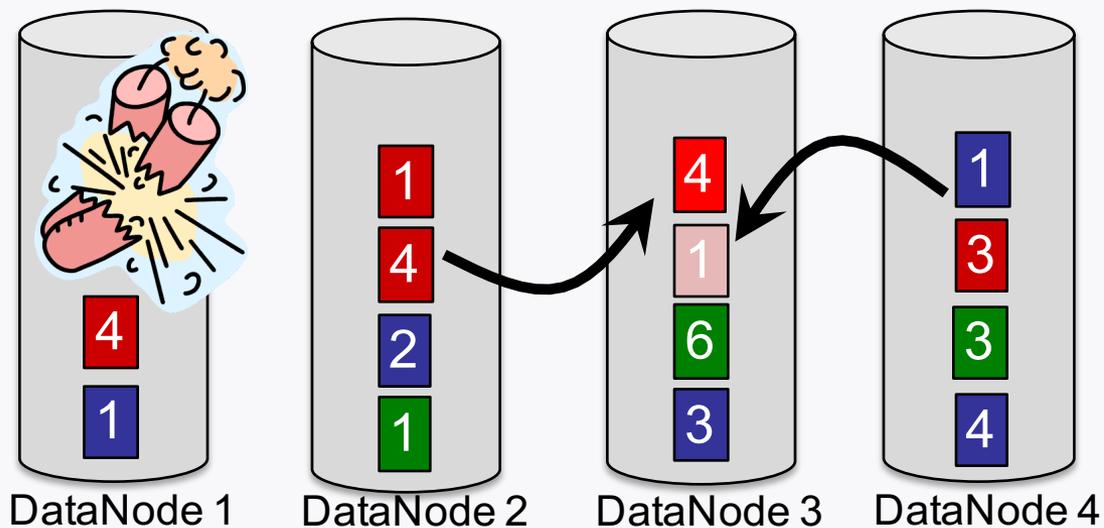
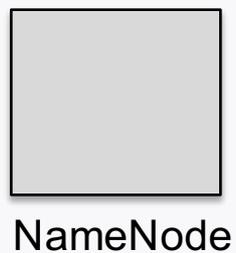
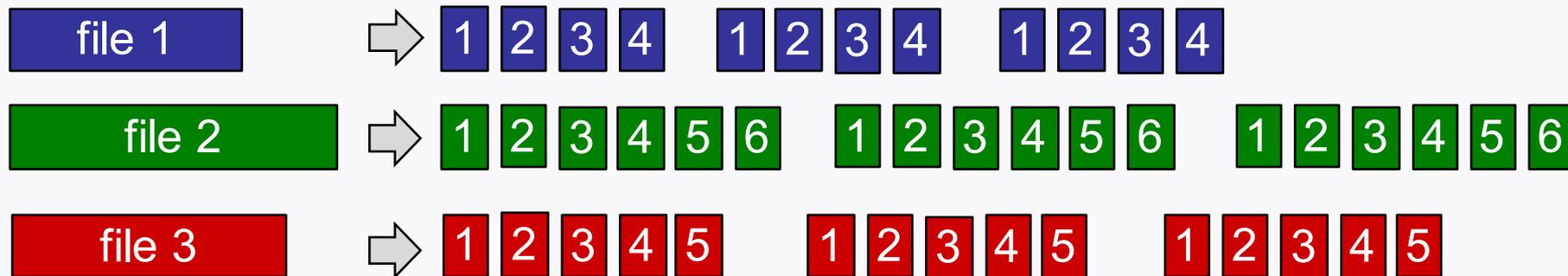
current hadoop architecture



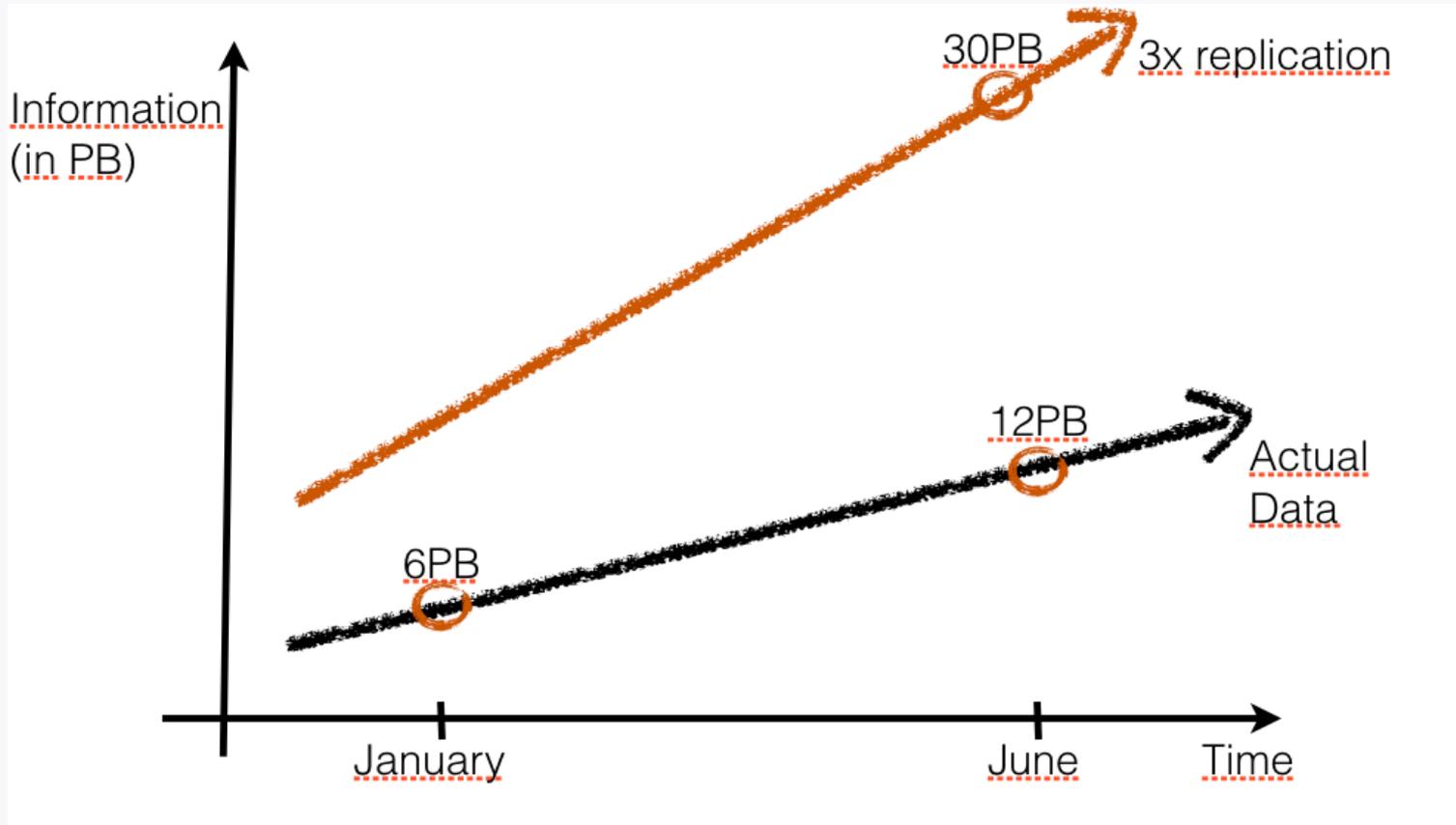
current hadoop architecture



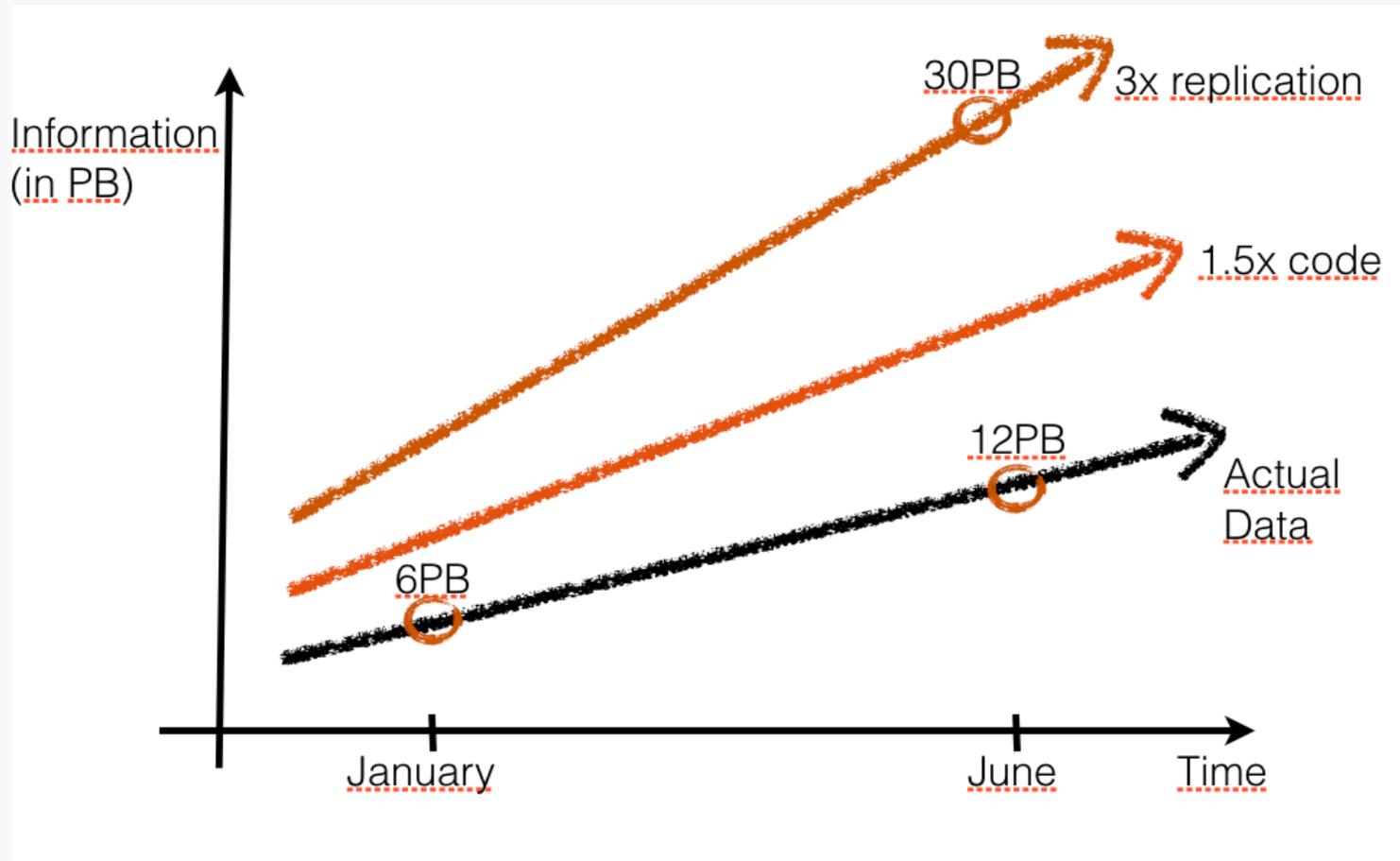
current hadoop architecture



erasure codes save space



erasure codes save space



Real systems that use distributed storage codes

- Windows Azure, (Cheng et al. USENIX 2012) (LRC Codes)
- Ships in Azure, Microsoft Server 2012 R2 and Windows 8.1

Real systems that use distributed storage codes

- Windows Azure, (Cheng et al. USENIX 2012) (LRC Codes)
- Ships in Azure, Microsoft Server 2012 R2 and Windows 8.1
- CORE (PPC Li et al. MSST 2013) (Regenerating EMSR Code)
- NCCloud (Hu et al. USENIX FAST 2012) (Regenerating Functional MSR)
- ClusterDFS (Pamies Juarez et al.) (SelfRepairing Codes)
- StorageCore (Esmaili et al.) (over Hadoop HDFS)

Real systems that use distributed storage codes

- Windows Azure, (Cheng et al. USENIX 2012) (LRC Codes)
- Ships in Azure, Microsoft Server 2012 R2 and Windows 8.1
- CORE (PPC Li et al. MSST 2013) (Regenerating EMSR Code)
- NCCloud (Hu et al. USENIX FAST 2012) (Regenerating Functional MSR)
- ClusterDFS (Pamies Juarez et al.) (SelfRepairing Codes)
- StorageCore (Esmaili et al.) (over Hadoop HDFS)
- HACFS (Xia, Saxena, Blaum) (IBM) FAST 2015
- HDFS Xorbas (Sathiamoorthy et al. VLDB 2013) (over Hadoop HDFS) (LRC code on Facebook clusters)
- Facebook F4 uses local parities in production [OSDI 2014]

Coded hadoop

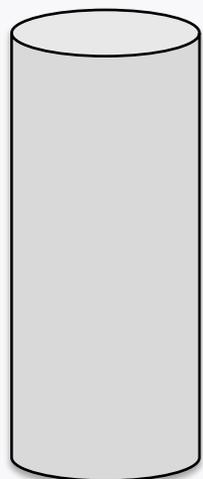
file 1



file 2



NameNode



DataNode 1



DataNode 2



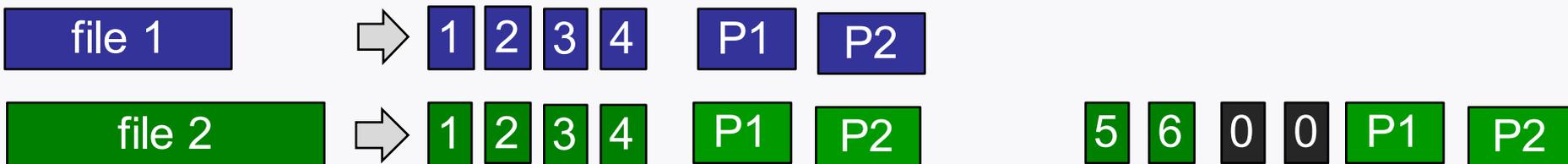
DataNode 3



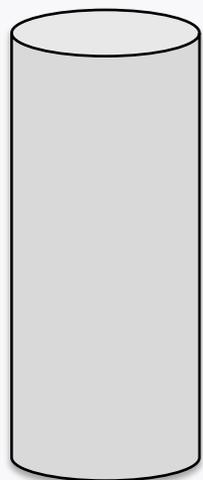
DataNode 4

...

Coded hadoop



NameNode



DataNode 1



DataNode 2



DataNode 3



DataNode 4

...

Code repair

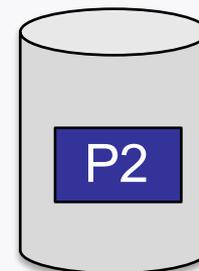
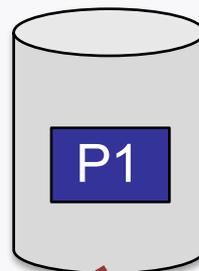
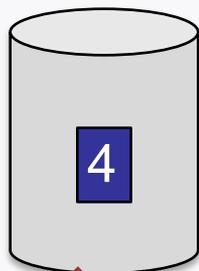
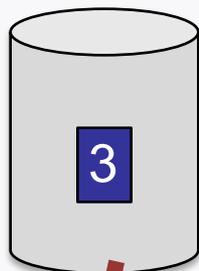
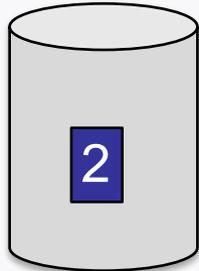
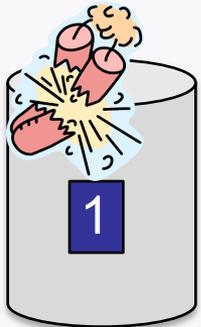
file 1



1 2 3 4

P1

P2



...

DataNode 1

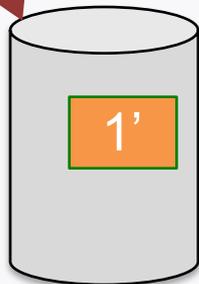
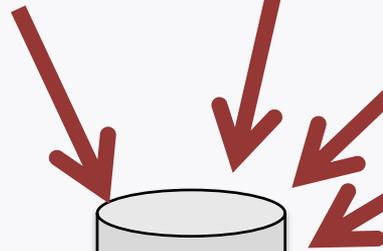
DataNode 2

DataNode 3

DataNode 4

DataNode 6

DataNode 6



DataNode 7 'newcomer'

Three repair metrics of interest

1. Number of bits communicated in the network during single node failures (**Repair Bandwidth**)

1. The number of bits read from disks during single node repairs (**Disk IO**)

3. The number of nodes accessed to repair a single node failure (**Locality**)

Three repair metrics of interest

1. Number of bits communicated in the network during single node failures (**Repair Bandwidth**)

Capacity known for two points only. My 3-year old conjecture for intermediate points was disproved. [ISIT13]

2. The number of bits read from disks during single node repairs (**Disk IO**)

3. The number of nodes accessed to repair a single node failure (**Locality**)

Three repair metrics of interest

1. Number of bits communicated in the network during single node failures (**Repair Bandwidth**)

Capacity known for two points only. My 3-year old conjecture for intermediate points was disproved. [ISIT13]

2. The number of bits read from disks during single node repairs (**Disk IO**)

Capacity unknown.

Only known technique is bounding by Repair Bandwidth

3. The number of nodes accessed to repair a single node failure (**Locality**)

Three repair metrics of interest

1. Number of bits communicated in the network during single node failures (**Repair Bandwidth**)

Capacity known for two points only. My 3-year old conjecture for intermediate points was disproved. [ISIT13]

2. The number of bits read from disks during single node repairs (**Disk IO**)

Capacity unknown.

Only known technique is bounding by Repair Bandwidth

3. The number of nodes accessed to repair a single node failure (**Locality**)

Capacity known for some cases.

Practical LRC codes known for some cases.

General constructions open

Three repair metrics of interest

1. Number of bits **communicated** in the network during single node failures (**Repair Bandwidth**)

Capacity known for two points only. My 3-year old conjecture for intermediate points was disproved. [ISIT13]

2. The number of bits **read** from disks during single node repairs (**Disk IO**)

Capacity unknown.

Only known technique is bounding by Repair Bandwidth

3. The number of nodes **accessed** to repair a single node failure (**Locality**)

Capacity known for some cases.

Practical LRC codes known for ~~some cases~~. Almost all cases

~~General constructions open~~—[Tamo-Barg!]

Code repair bandwidth

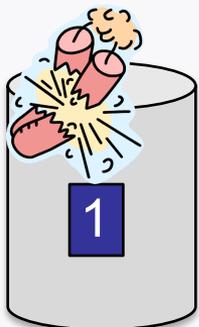
file 1



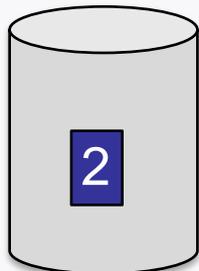
1 2 3 4

P1

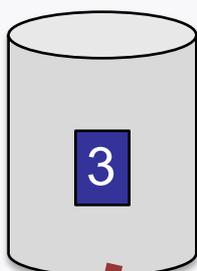
P2



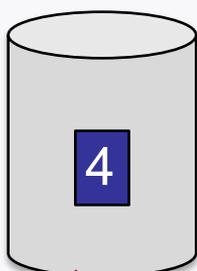
1



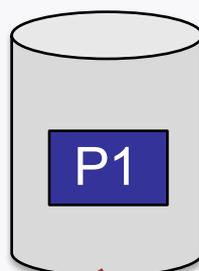
2



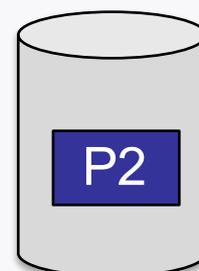
3



4



P1



P2

DataNode 1

DataNode 2

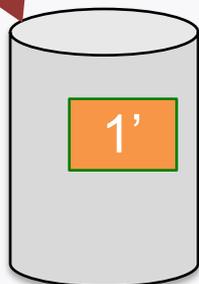
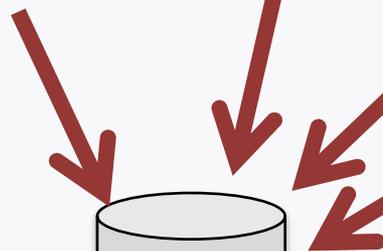
DataNode 3

DataNode 4

DataNode 6

DataNode 6

...

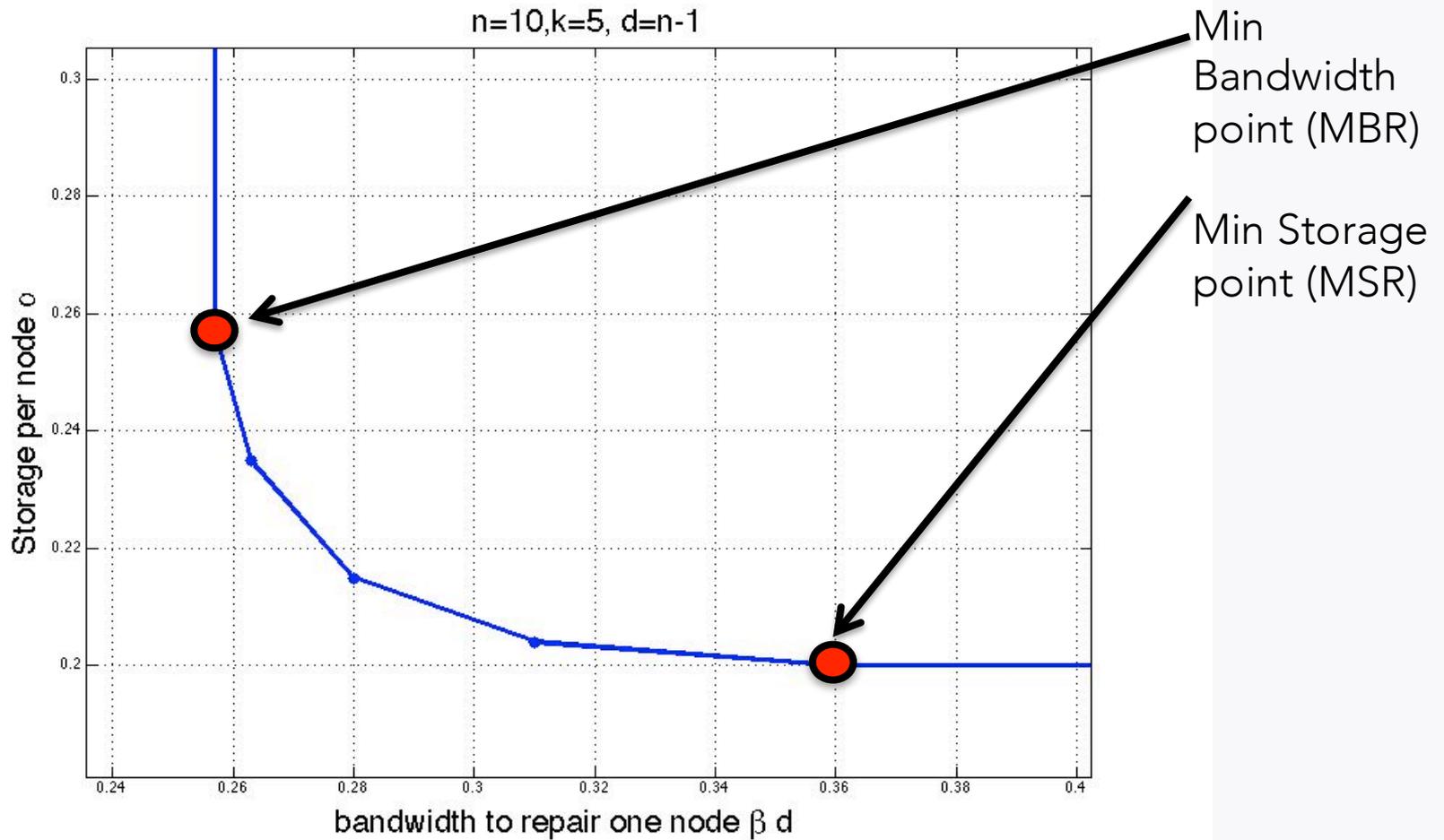


DataNode 7 'newcomer'

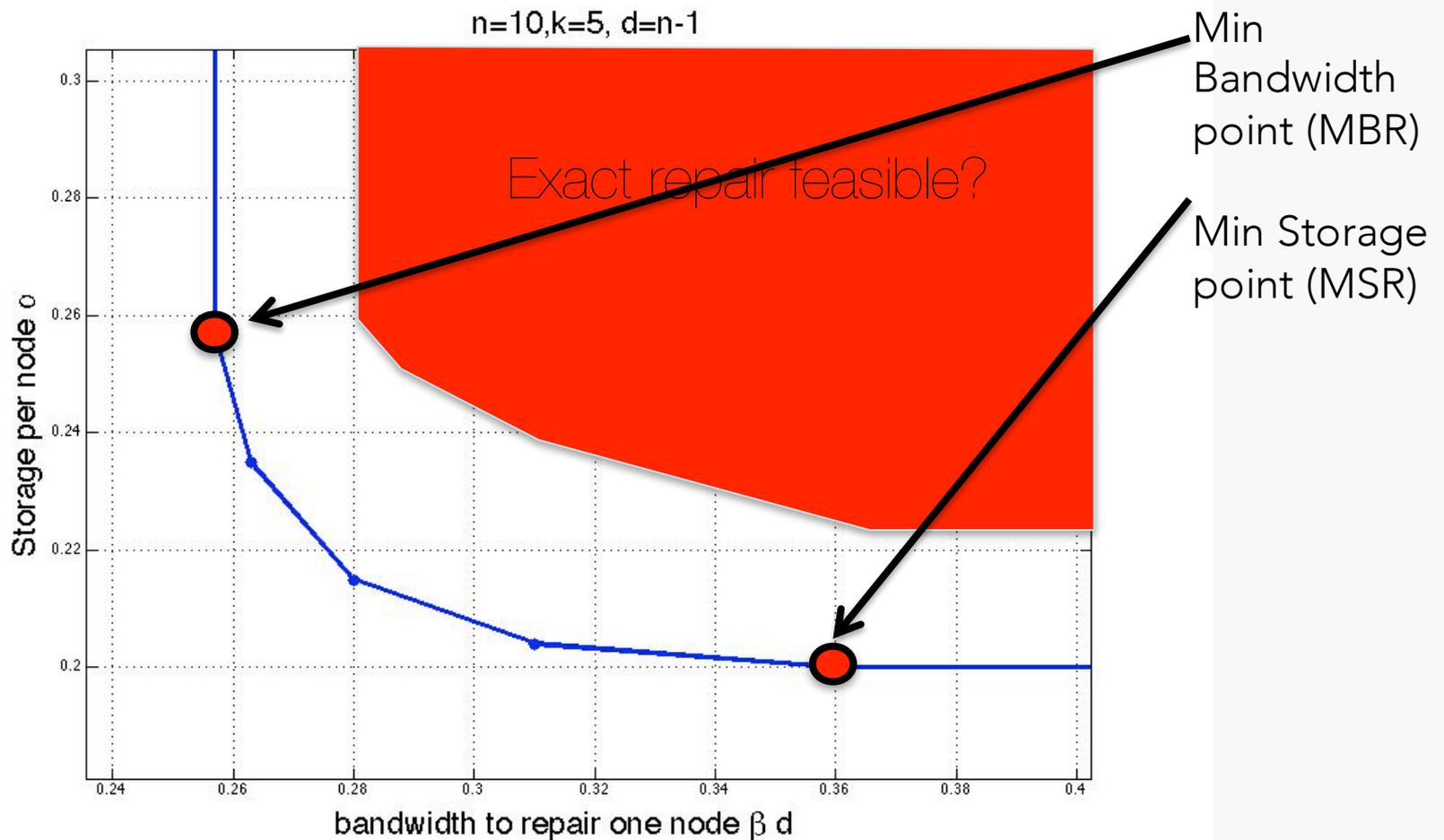
Functional repair: $1' \neq 1$
(but MDS distance maintained)

Exact repair: $1' = 1$

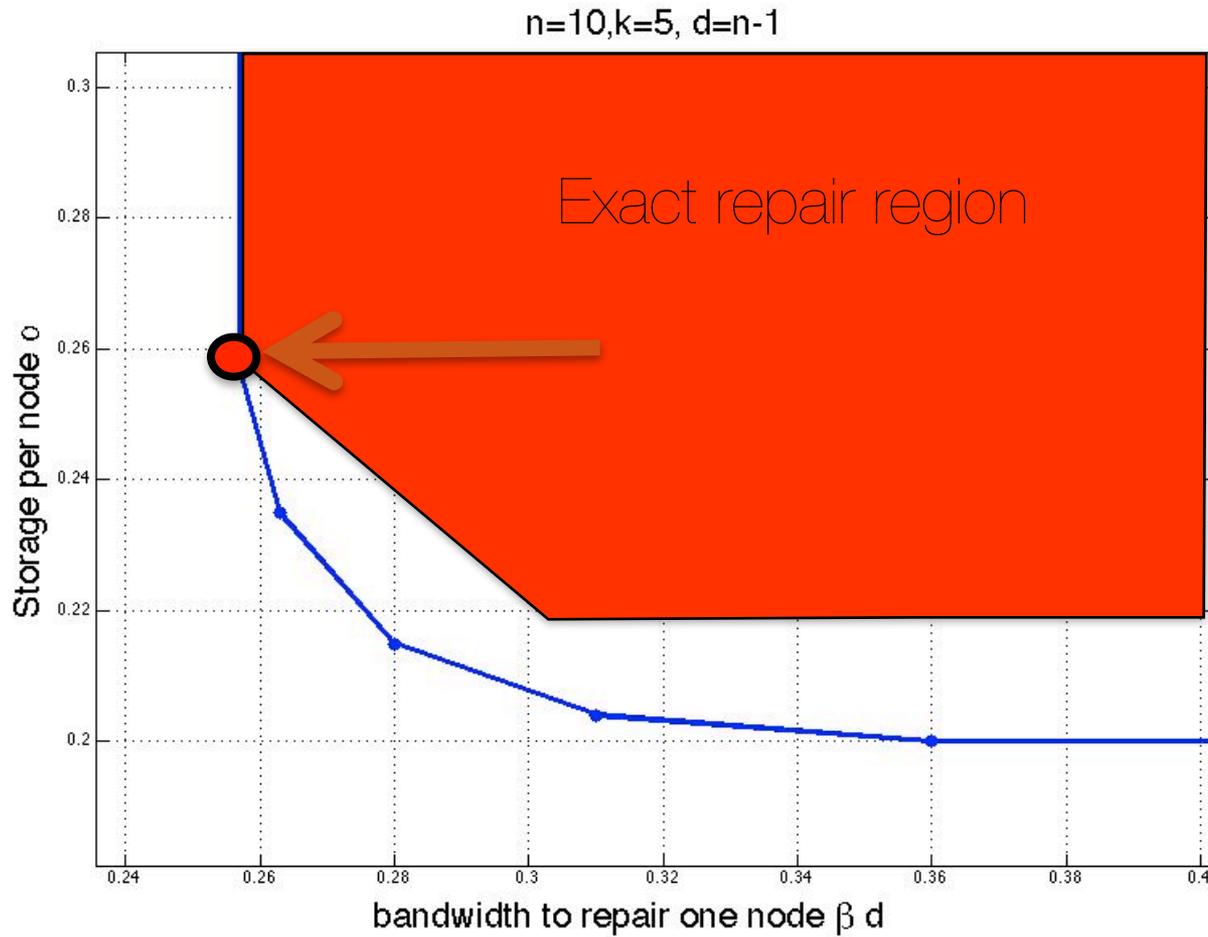
Repair Bandwidth Tradeoff Region



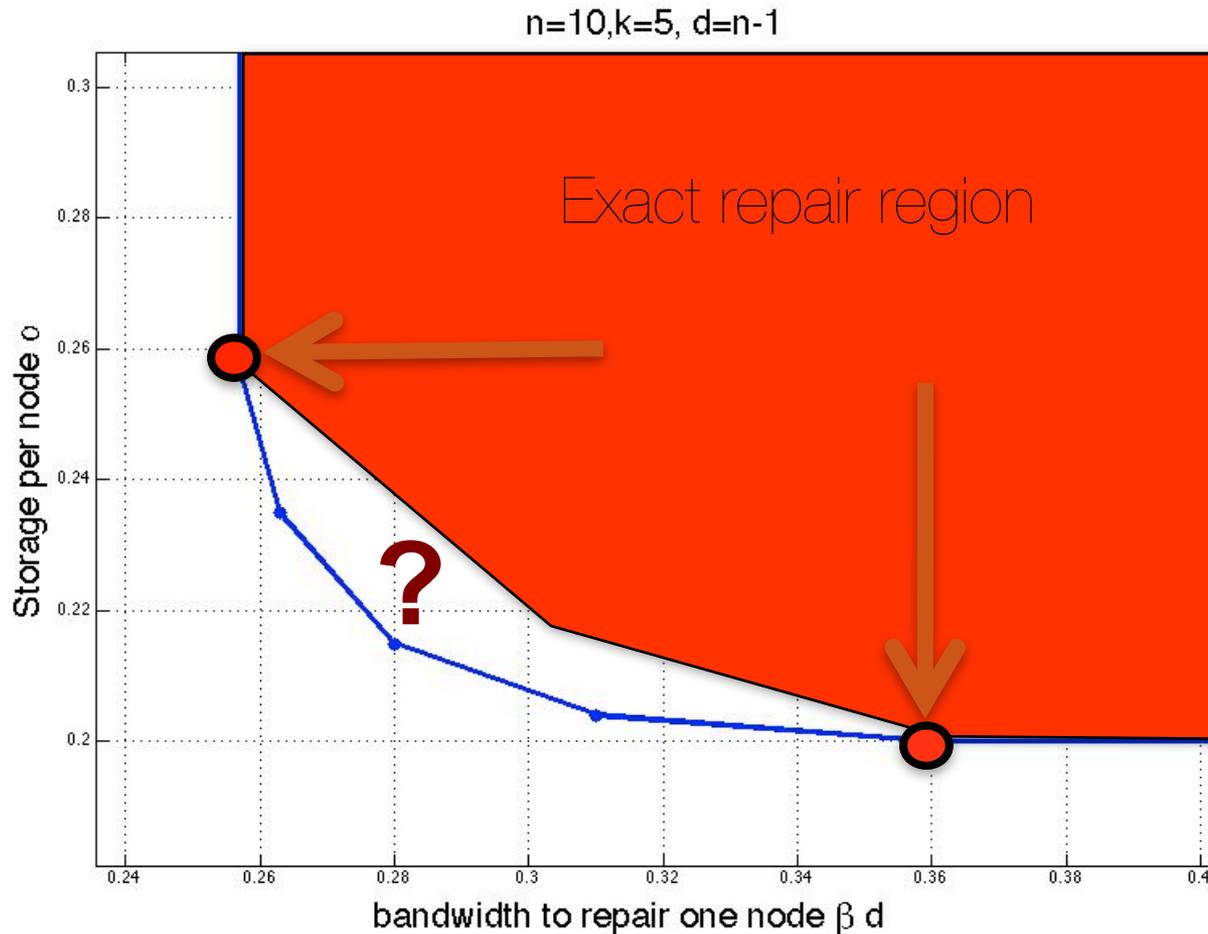
Exact repair region?



Status in 2011

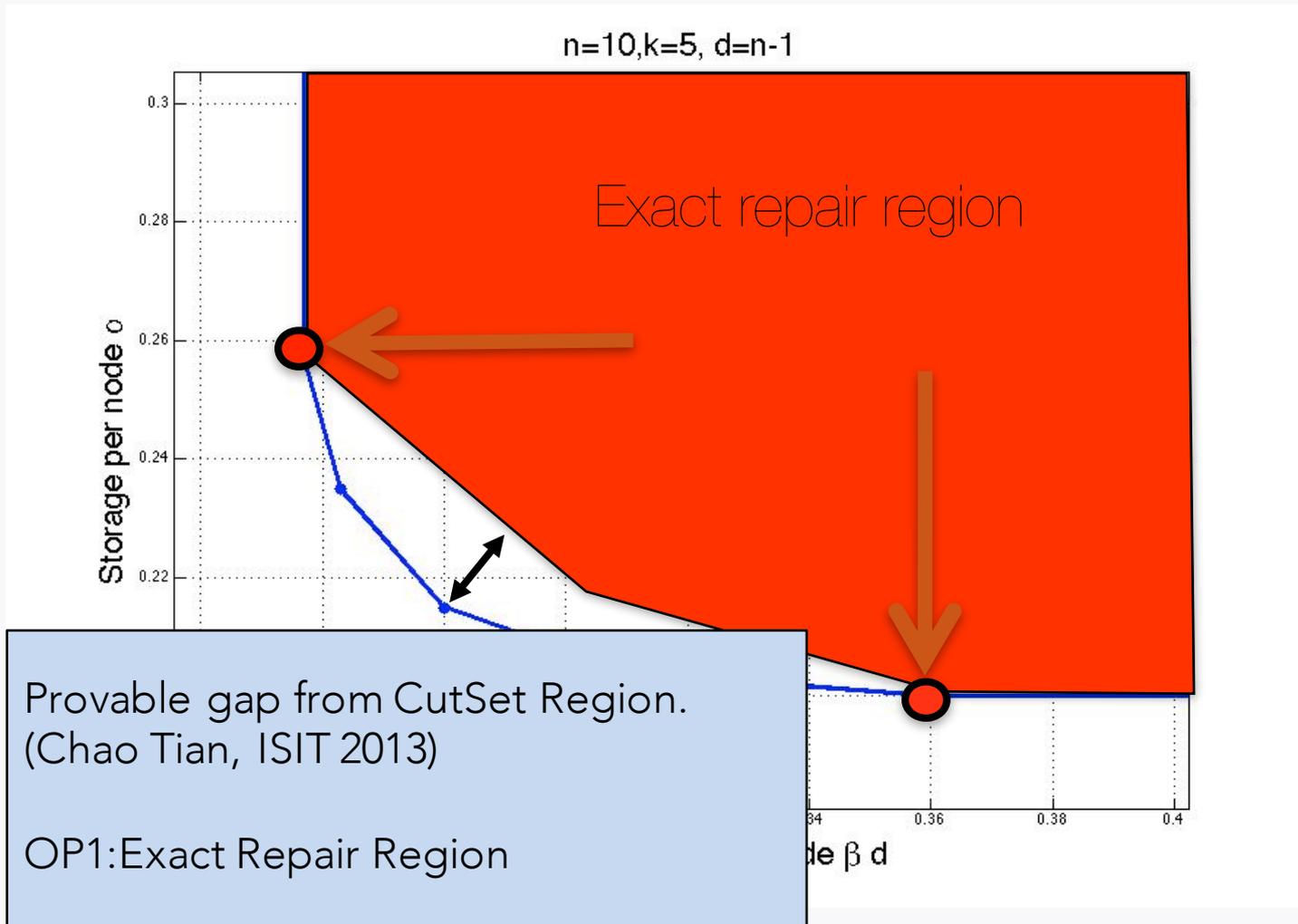


Status in 2012



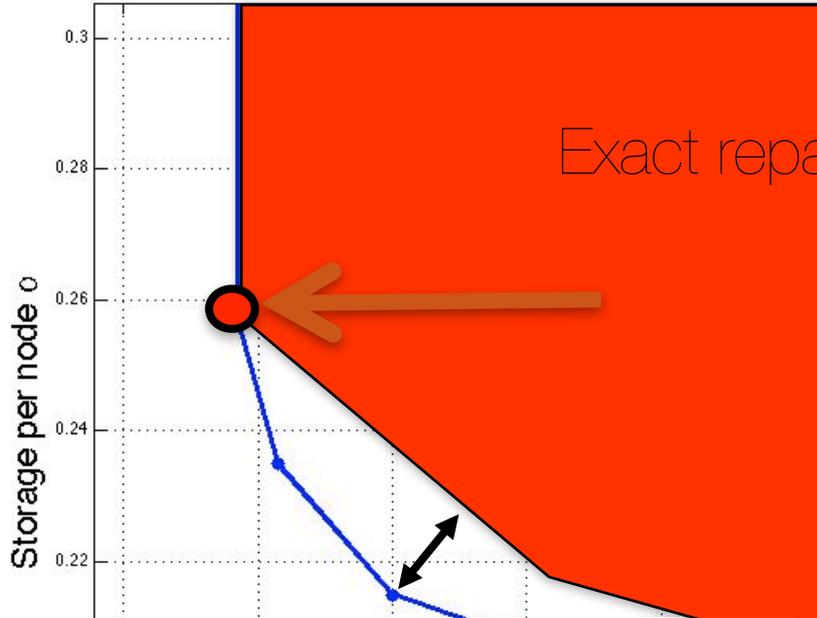
Code constructions by:
Rashmi, Shah, Kumar
Suh, Ramchandran
El Rouayheb, Shum,
Oggier, Datta
Silberstein, Viswanath
et al.
Cadambe, Maleki,
Jafar
Le Scouarnec et al.
Papailiopoulos, Wu,
Dimakis
Wang, Tamo, Bruck
Tamo, Barg

Status in 2013



Status in 2014

$n=10, k=5, d=n-1$



Provable gap from CutSet Region.
(Chao Tian, ISIT 2013)

OP1: Exact Repair Region



Taking a step back

- Finding exact regenerating codes is still an open problem in coding theory
- What can we do to make progress ?

Taking a step back

- Finding exact regenerating codes is still an open problem in coding theory
- What can we do to make progress ?



Taking a step back

- Finding exact regenerating codes is still an open problem in coding theory
- What can we do to make progress ?



or change the question

Changing the question: Locally Repairable Codes

Three repair metrics of interest

1. Number of bits communicated in the network during single node failures (**Repair Bandwidth**)

Capacity known for two points only. My 3-year old conjecture for intermediate points was just disproved. [ISIT13]

2. The number of bits read from disks during single node repairs (**Disk IO**)

Capacity unknown.

Only known technique is bounding by Repair Bandwidth

3. The number of nodes accessed to repair a single node failure (**Locality**)

Capacity known for some cases.

Practical LRC codes known and used!

Minimum Distance

- The distance of a code d is the minimum number of erasures after which data is lost.
- Reed-Solomon (10,14) ($n=14, k=10$). $d= 5$
- R. Singleton (1964) showed a bound on the best distance possible:

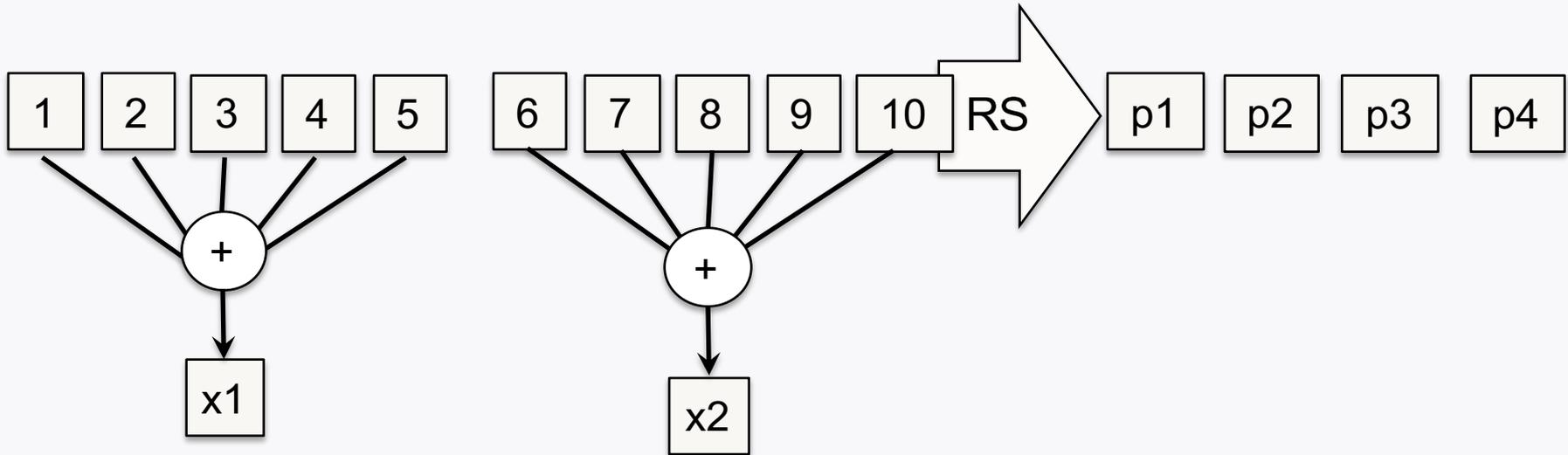
$$d \leq n - k + 1$$

- Reed-Solomon codes achieve the Singleton bound (hence called MDS)

Locality of a code

- A code symbol has locality r if it is a function of r other codeword symbols.
- A systematic code has **message locality** r if all its systematic symbols have locality r
- A code **has all-symbol locality** r if all its symbols have locality r .
- In an MDS code, all symbols have locality at most $r \leq k$
- **Easy lemma:** Any MDS code must have trivial locality $r=k$ for every symbol.

Example: code with message locality 5



All $k=10$ message blocks can be recovered by reading $r=5$ other blocks.

A single parity block failure requires still 10 reads.

Best distance possible for a code with locality r ?

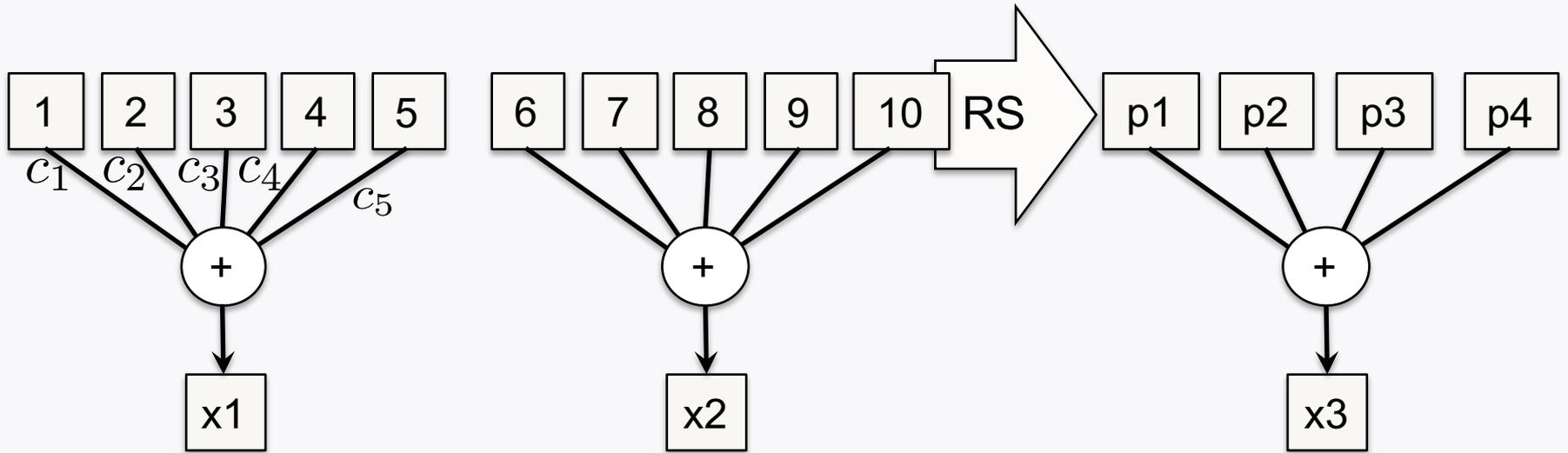
Locality-distance tradeoff

Codes with all-symbol locality r can have distance at most:

$$d \leq n - k - \left\lceil \frac{k}{r} \right\rceil + 2$$

- Shown by Gopalan et al. for scalar linear codes (Allerton 2012)
- Papailiopoulos et al. information theoretically (ISIT 2012)
- $r=k$ (trivial locality) gives Singleton Bound.
- Any non-trivial locality will hurt the fault tolerance of the storage system
- Pyramid codes (Huang et al) achieve this bound for message-locality

All-symbol locality



The coefficients need to make the local forks in general position compared to the global parities.

Random works whp in exponentially large field. Checking requires exponential time.

OP2: General Explicit LRCs that are maximally recoverable (MR) are open.

Three repair metrics of interest

1. Number of bits communicated in the network during single node failures (**Repair Bandwidth**)

Capacity known for two points only. My 3-year old conjecture for intermediate points was just disproved. [ISIT13]

2. The number of bits read from disks during single node repairs (**Disk IO**)

Capacity unknown.

Only known technique is bounding by Repair Bandwidth

3. The number of nodes accessed to repair a single node failure (**Locality**)

Capacity known for some cases.

Practical LRC codes known and used!

Three repair metrics of interest

1. Number of bits communicated in the network during single node

RC Capacity open but some results [Mohajer, Tandon, Tian]
No practical High-rate MSR codes known. Would be useful if we can find some.

2. The number of bits read from disks during single node repairs
(Disk IO)

Capacity unknown.

Only known technique is bounding by Repair Bandwidth

3. The number of nodes accessed to repair a single node failure
(Locality)

Capacity known for some cases.

Practical LRC codes known and used!

Three repair metrics of interest

1. Number of bits **communicated** in the network during single node

RC Capacity open but some results [Mohajer, Tandon, Tian]
No practical High-rate MSR codes known. Would be useful if we can find some.

2. The number of bits **read** from disks during single node repairs

(D)

Not much activity for Disk IO repair.
Maybe not a bottleneck for current technology

3. The number of nodes accessed to repair a single node failure

(Locality)

Capacity known for some cases.
Practical LRC codes known and used!

Three repair metrics of interest

1. Number of bits **communicated** in the network during single node

RC Capacity open but some results [Mohajer, Tandon, Tian]
No practical High-rate MSR codes known. Would be useful if we can find some.

2. The number of bits **read** from disks during single node repairs

(D)

Not much activity for Disk IO repair.
Maybe not a bottleneck for current technology?

3. The number of nodes **accessed** to repair a single node failure

(L)

Useful constructions and bounds. Influencing real systems.
MR LRCs would be directly useful.

Dealing with Hot data

Codes used for cold data, i.e. data not read very frequently
(Data Analytics clusters, Huge text file logs, Offline queries).

Dealing with Hot data

Codes used for cold data, i.e. data not read very frequently
(Data Analytics clusters, Huge text file logs, Offline queries).

Warm and Hot data (Haystack for photo storage, Video caching and delivery, analytics in interactive time, adaptive training of big machine learning models)

Dealing with Hot data

Codes used for cold data, i.e. data not read very frequently
(Data Analytics clusters, Huge text file logs, Offline queries).

Warm and Hot data (Haystack for photo storage, Video caching and delivery, analytics in interactive time, adaptive training of big machine learning models)

Multiple jobs or threads concurrently reading the same data blocks.

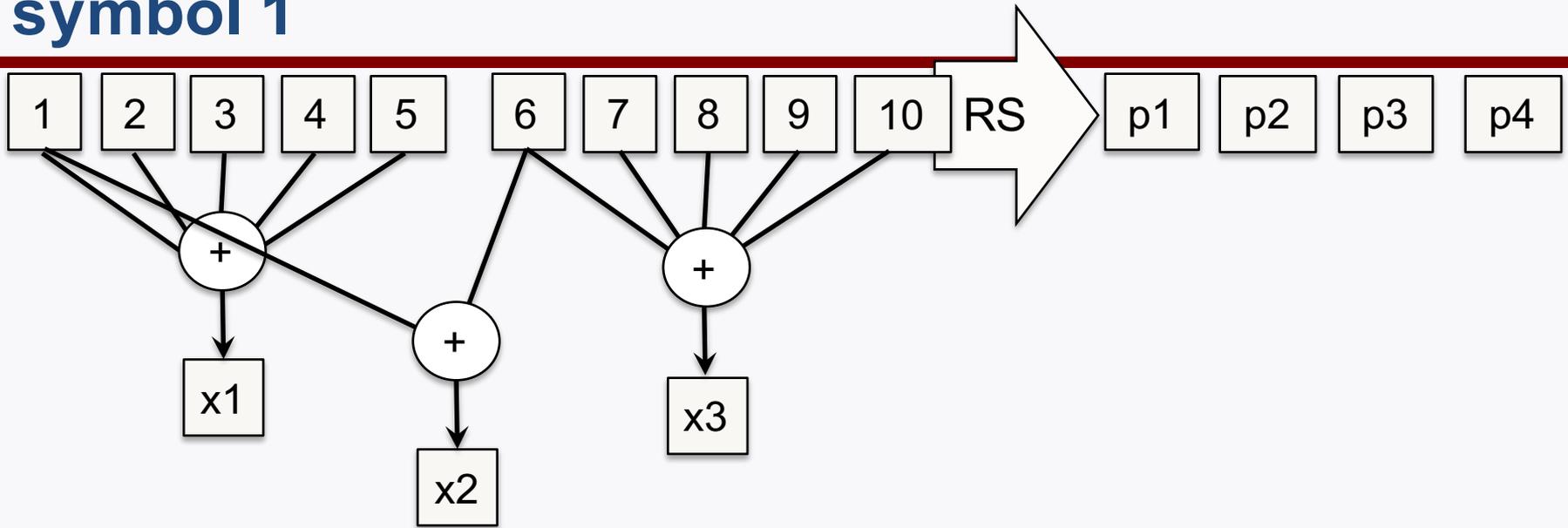
Some disks or servers become hot- use coding to relieve this.

Provide similar performance as replication with smaller cost.

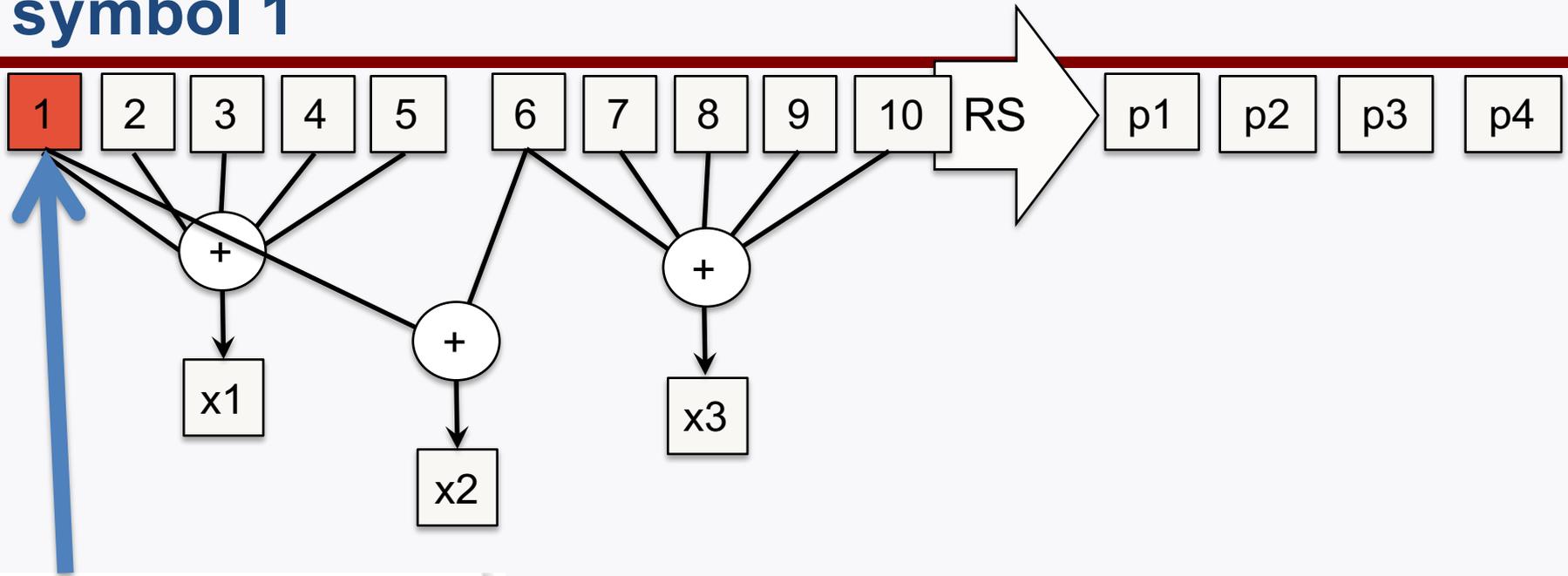
Code Locality r , Code Availability t

- A symbol has locality r if it is a function of r other codeword symbols.
- A code **has all-symbol locality r** if all its symbols have locality r .
- A symbol has **availability t** if it can be read in parallel by $t+1$ disjoint groups of symbols.
- These t reads have locality r if they involve up to r symbols each.

Example of Locality r and availability t for symbol 1



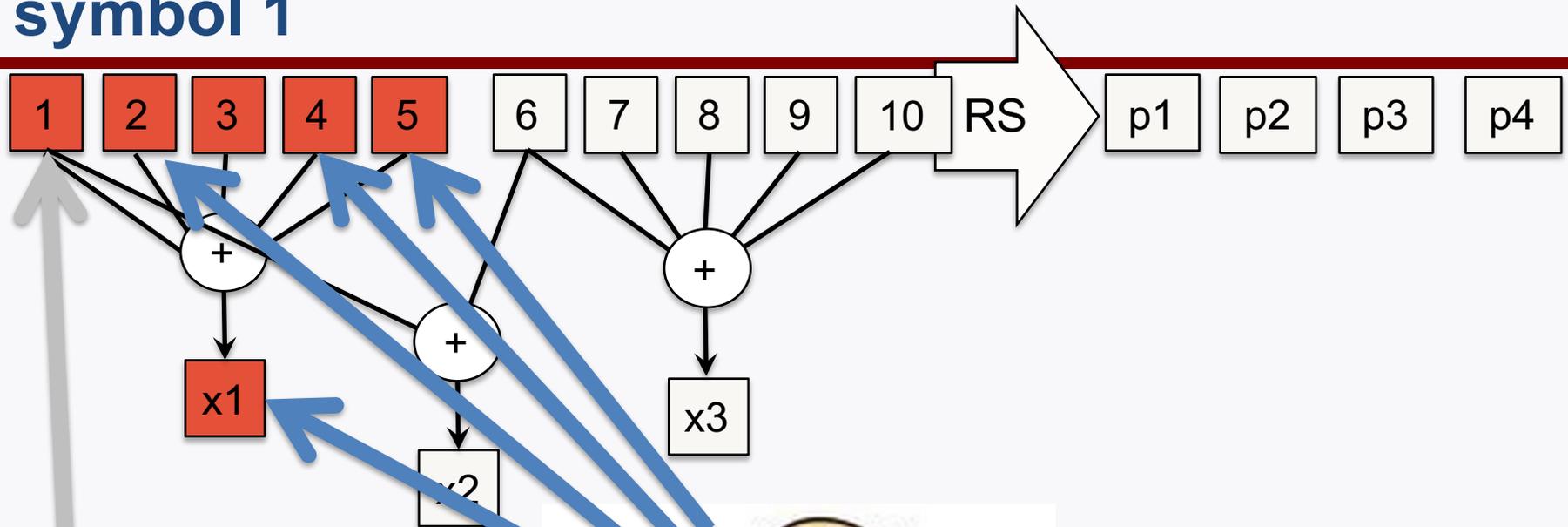
Example of Locality r and availability t for symbol 1



Want to read
Block 1

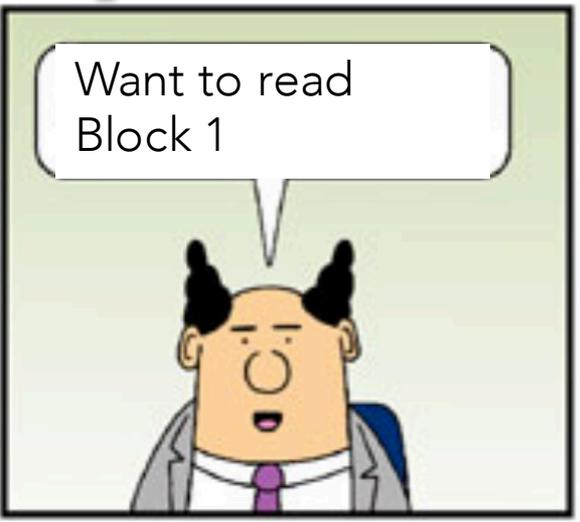
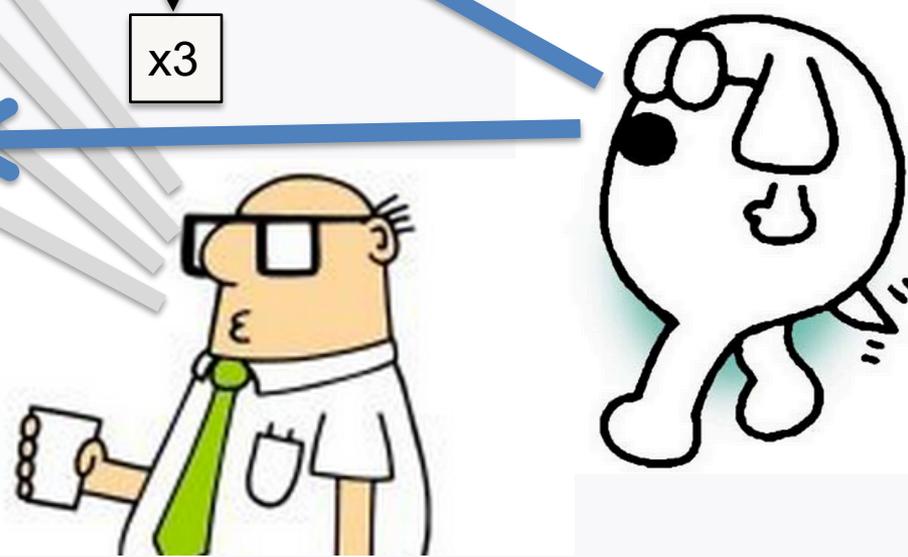
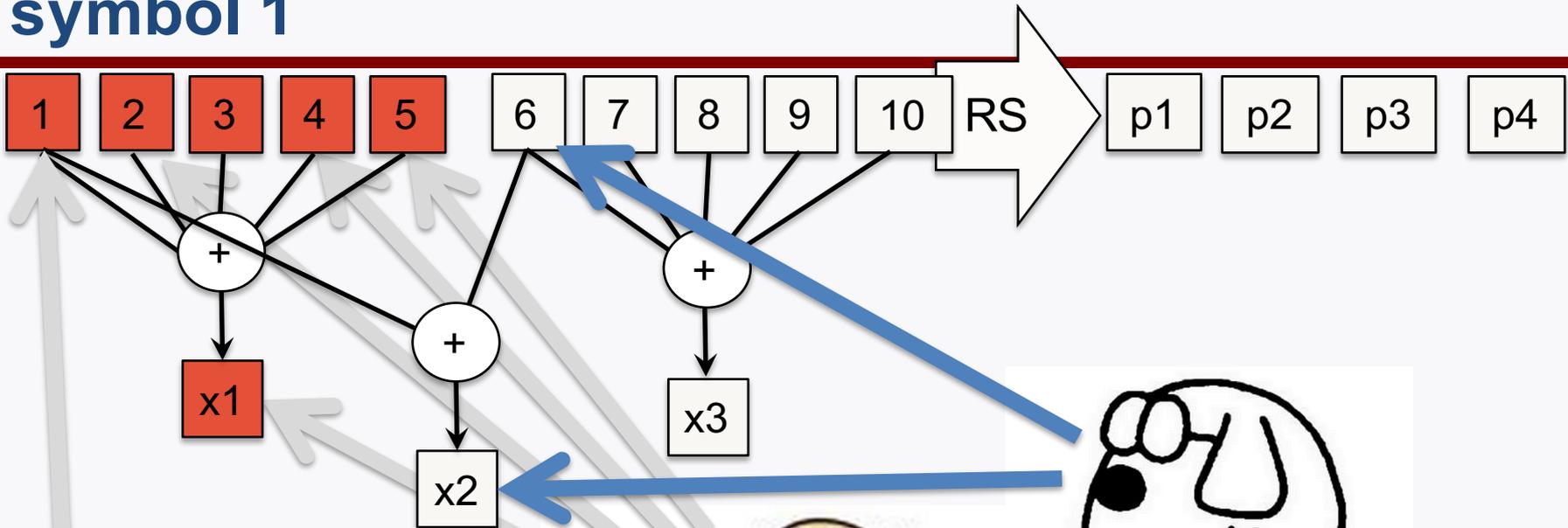


Example of Locality r and availability t for symbol 1

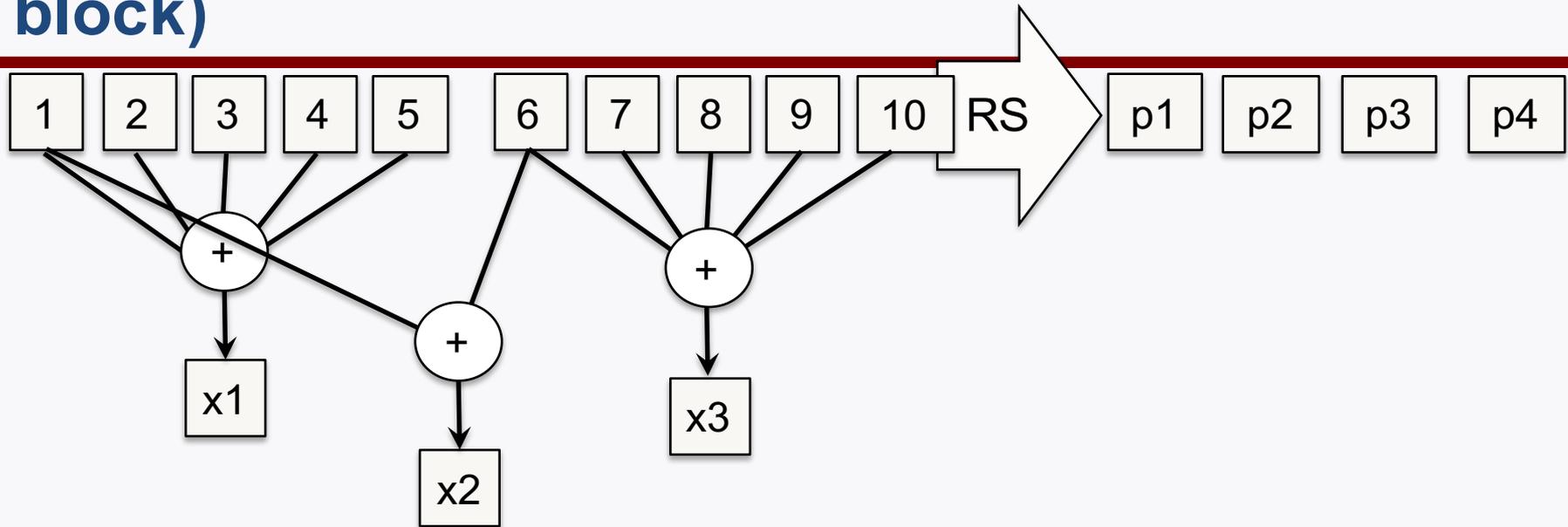


Want to read
Block 1

Example of Locality r and availability t for symbol 1

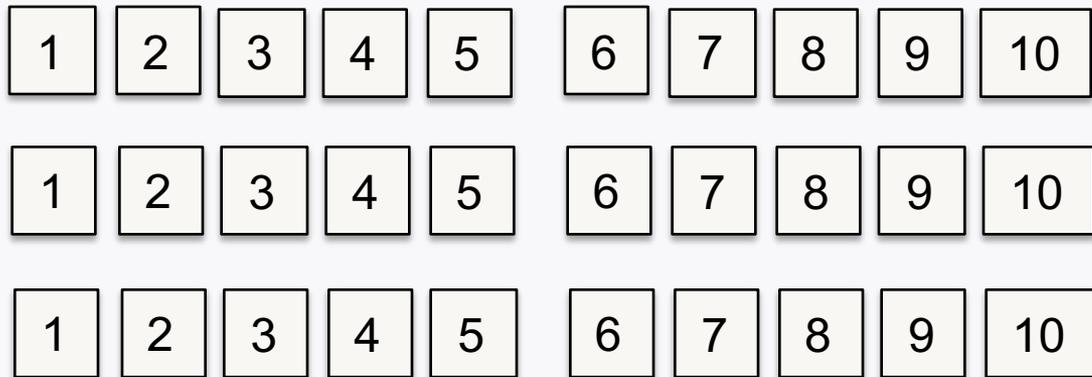


message availability 2 (=2 parallel reads for a block)



- Therefore Block 1 can be read by 1 systematic read + 2 repair reads **simultaneously**
- Block 1 has availability $t=2$ with groups of locality $r_1=5$ and $r_2=2$
- Notice also that the group (2,3,4,5,6,7,8,9,10, p1) of locality $r=10$ can be used to recover 1 (but blocks all others, so not used)

Example: 3 replication



- Each symbol can be read in parallel $t+1 = 3$ times.
- Distance $d=3$. Rate= $1/3$.
- Availability $t=2$. Locality of these reads $r=1$.
- If you want to increase availability, rate goes to zero like $1 / (t+1)$
- Can we get scaling availability with non-vanishing rate?

Our results

We construct codes with scaling availability and small locality.
For any high rate. With near-MDS distance.

- Polynomial Availability (using Combinatorial designs):

$$t = n^{1/3}$$

$$r = n^{1/3 - \epsilon}$$

- Fundamental Bounds: For a given locality r and availability t requirements, what is the best distance possible?
- We obtain some bounds – Sometimes tight.

Related work

- Locally decodable codes
(LDCs imply linear availability, $t = c n$)
- Batch Codes [Ishai, Kushilevitz, Ostrovsky, Sahai STOC'04].

Very similar parallel reads requirement.

Not good distance.

In fact our results imply the first batch codes with near-MDS distance.

Conclusions and Open Problems

Which repair metric to optimize?

- Repair BW, All-Symbol Locality, Message-Locality, Fault tolerance, A combination of all?
- Depends on type of storage cluster
(Cloud, Analytics, Photo Storage, Archival, Hot vs Cold data)

Which repair metric to optimize?

- Repair BW, All-Symbol Locality, Message-Locality, Fault tolerance, A combination of all?
- Depends on type of storage cluster
(Cloud, Analytics, Photo Storage, Archival, Hot vs Cold data)
- Practice involves a spectrum of tradeoffs

Which repair metric to optimize?

- Repair BW, All-Symbol Locality, Message-Locality, Fault tolerance, A combination of all?
- Depends on type of storage cluster
(Cloud, Analytics, Photo Storage, Archival, Hot vs Cold data)
- Practice involves a spectrum of tradeoffs



Which repair metric to optimize?

- Repair BW, All-Symbol Locality, Message-Locality, Fault tolerance, A combination of all?
- Depends on type of storage cluster
(Cloud, Analytics, Photo Storage, Archival, Hot vs Cold data)
- Practice involves a spectrum of tradeoffs



Which repair metric to optimize?

- Repair BW, All-Symbol Locality, Message-Locality, Fault tolerance, A combination of all?
- Depends on type of storage cluster
(Cloud, Analytics, Photo Storage, Archival, Hot vs Cold data)
- Practice involves a spectrum of tradeoffs



Seven open problems

Repair Bandwidth:

- 1. Exact repair region ?
- 2. Practical E-MSR codes for high rates ?
- 3. Better Repair for existing codes (EvenOdd, RDP, Reed-Solomon) ?

Seven open problems

Repair Bandwidth:

- 1. Exact repair region ?
- 2. Practical E-MSR codes for high rates ?
- 3. Better Repair for existing codes (EvenOdd, RDP, Reed-Solomon) ?

Locality:

- 4. Explicit LRCs with Maximum recoverability?

Seven open problems

Repair Bandwidth:

- 1. Exact repair region ?
- 2. Practical E-MSR codes for high rates ?
- 3. Better Repair for existing codes (EvenOdd, RDP, Reed-Solomon) ?

Locality:

- 4. Explicit LRCs with Maximum recoverability?

Availability:

- 5. Distance –availability tradeoff ?
- 6. Practical explicit codes ?
- 7. Approximating GLRC/Index Coding **sum rate** in polylog factor ?

Coding for Storage wiki

[[wiki:definitions:repair_problem]] DISTRIBUTED STORAGE WIKI

Trace: » repair_problem

Show pagesource Old revisions Recent changes Index Login

Navigation

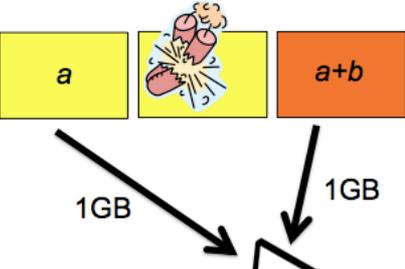
- [Homepage](#)
- [Papers](#)
- [Authors](#)
- [Codes](#)
- [Software](#)
- [Definitions](#)
- [Open Problems](#)
- [Conferences](#)
- [Guidelines](#)
- [Playground](#)
- [Syntax](#)
- [itsociety](#)

The Repair Problem

Consider the very simple ($n=3, k=2$) binary MDS code shown, which is very simply one disk storing the parity of the others



Clearly, this code has the property it can tolerate any single node failure. This means that even after one node fails, a data collector (shown as a laptop) can communicate the information from the remaining two nodes and reconstruct the file. This is shown here:



Distance vs. Locality-Availability trade-off

New Distance Bound

- For (r, t) -Information local codes*:

$$d_{\min} \leq n - k + 1 - \left(\left\lceil \frac{kt}{r} \right\rceil - t \right)$$

Distance vs. Locality-Availability trade-off

New Distance Bound:

- For (r, t) -Information local codes*:

$$d_{\min} \leq n - k + 1 - \left(\left\lceil \frac{kt}{r} \right\rceil - t \right)$$

*The dirty details:

- We can only prove this for scalar linear codes.
- Only one parity symbol per repair group is assumed.
- For some cases we can achieve this using combinatorial designs.

HDFS Xorbas



HOME STARTUPS MOBILE GADGETS EUROPE VIDEO MORE SEARCH

Adobe Marketing Cloud
Marketing works, and we can prove it. [Learn more](#)

HOT TOPICS YAHOO APPLE FACEBOOK TWITTER GOOGLE MICROSOFT NSA

CrunchGov CrunchU Guides Events CrunchBase

Crunch DISRUPT SF 2013 Limited Extra Early Bird Tickets Available Until July 15 [GET TICKETS NOW](#)

Comment 7 Like 173 Tweet 360 Share 42 +1 144

Check Out Facebook's Nerdy Library Of Its Research Papers

JOSH CONSTINE

Thursday, May 16th, 2013 7 Comments

FACEBOOK ACADEMICS

If subjects like "XORing Elephants: Novel Erasure Codes for Big Data" get you all worked up, you'll dig the "Research Publications At Facebook" site, which collects scientific papers written by Facebook employees and researchers. Ranging from hardcore engineering to the sociology of social networks, the library puts Facebook's open-sourced knowledge all in one place.

ZDNet

White Papers Hot Topics Downloads Reviews Newsletters

US Edition NextGen Networks Data Center M2M BYOD CXO Apple Tablets

LEADING EDGE IN CLOUD

MUST READ: *With Windows 8.1, can Microsoft get its mojo*

Topic: Storage Investigate Follow via: RSS Email

How efficient can storage be?

Summary: We want our data protected from failures. After a failure we want our data back quickly. And we want to pay as little as possible. How?

By Robin Harris for Storage Bits | June 19, 2013 -- 07:00 GMT (00:00 PDT)

Comment 1 Vote 1 Like 6 Tweet 33 Share more +

Recent research by hyper-scale system managers - mostly Microsoft and Facebook engineers and scientists - has tried to answer that question. And the answers are way better than what we have today.

In XORing Elephants: Novel Erasure Codes for Big Data, authors Maheswaran Sathiamoorthy, Alexandros G. Dimakis, Megasthenis Asteris, Dhruva Borthakur and Dimitris Papailiopoulos of USC and Ramkumar Vadali and Scott Chen of Facebook delve deeply into the issue. Technically it is related to work that Microsoft presented last year.

High Scalability

Building bigger, faster, more reliable websites.

Home Real Life Architectures Strategies All Posts Advertising Book Store Start Here Contact

All Time Favorites RSS Twitter Facebook G+

« Stuff The Internet Says On Scalability For June 28, 2013 | Main | Leveraging Cloud Computing at Yelp - 102 Million Monthly Visitors and 39 Million Reviews »

Paper: XORing Elephants: Novel Erasure Codes For Big Data

THURSDAY, JUNE 27, 2013 AT 8:45AM

Erasure codes are one of those seemingly magical mathematical creations that with the developments described in the paper [XORing Elephants: Novel Erasure Codes for Big Data](#), are set to replace triple replication as the data storage protection mechanism of choice.

Figure 3: The Markov model used to calculate the MTTDL_{erasure} of (10, 4) RS and (10, 6, 5) LRC.

Scheme	Storage overhead	Repair traffic	MTTDL (days)
Triple replication	2x	1x	2,307,000 - 10

Open Your Eyes

catchpoint

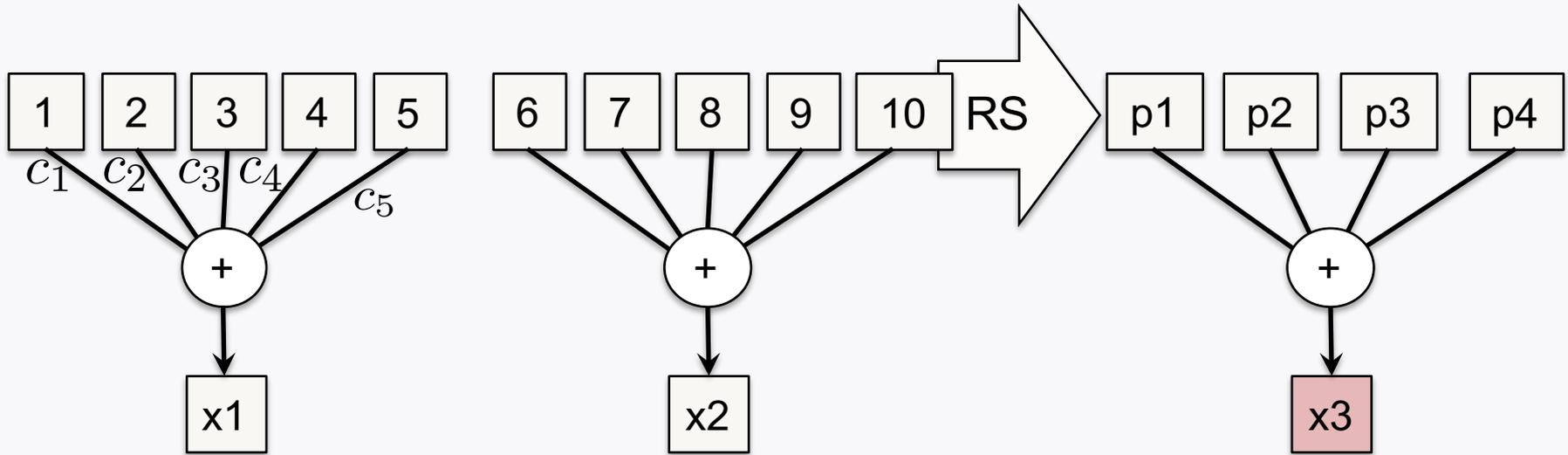
Booking.com

Join our team of engineers in Amsterdam

Now Hiring

NOSQL NOW!

code we implemented in HDFS



Single block failures can be repaired by accessing 5 blocks. (vs 10)

Stores 16 blocks

1.6x Storage overhead vs 1.4x in HDFS RAID.

Implemented this in Hadoop (system available on [github/madiator](https://github.com/madiator))

Java implementation

```
public void encode(int[] message, int[] parity) {  
  
    assert(message.length == stripeSize && parity.length == paritySizeRS+paritySizeSRC);  
  
    for (int i = 0; i < paritySizeRS; i++) {  
        dataBuff[i] = 0;  
    }  
    for (int i = 0; i < stripeSize; i++) {  
        dataBuff[i + paritySizeRS] = message[i];  
    }  
    GF.remainder(dataBuff, generatingPolynomial);  
    for (int i = 0; i < paritySizeRS; i++) {  
        parity[paritySizeSRC+i] = dataBuff[i];  
    }  
    for (int i = 0; i < stripeSize; i++) {  
        dataBuff[i + paritySizeRS] = message[i];  
    }  
    for (int i = 0; i < paritySizeSRC; i++) {  
        for (int f = 0; f < simpleParityDegree; f++) {  
            parity[i] = GF.add(dataBuff[i*simpleParityDegree+f], parity[i]);  
        }  
    }  
}
```

Some experiments

NameNode 'ip-10-168-201-226.us-west-1.compute.internal:54310'

Started: Wed Jan 11 23:49:49 UTC 2012
Version: usc3xor, r
Compiled: Wed Jan 11 20:12:05 UTC 2012 by root
Upgrades: There are no upgrades in progress.

[Browse the filesystem](#)
[Namenode Logs](#)

Cluster Summary

276 files and directories, 1621 blocks = 1897 total. Heap Size is 36.37 MB / 966.69 MB (3%). Committed Heap: 58.69 MB. Init Heap: 16 MB. Non Heap Memory Size is 24.91 MB / 118 MB (21%). Committed Non Heap: 37.5 MB.

WARNING : There are 14 missing blocks. Please check the log or run fsck.

Configured Capacity	:	5.3 TB			
DFS Used	:	100.07 GB			
Non DFS Used	:	284.13 GB			
DFS Remaining	:	4.93 TB			
DFS Used%	:	1.84 %			
DFS Remaining%	:	92.93 %			
DataNodes usages	:	Min %	Median %	Max %	stdev %
	:	1.5 %	1.85 %	2.15 %	0.15 %
Number of Under-Replicated Blocks	:	0			

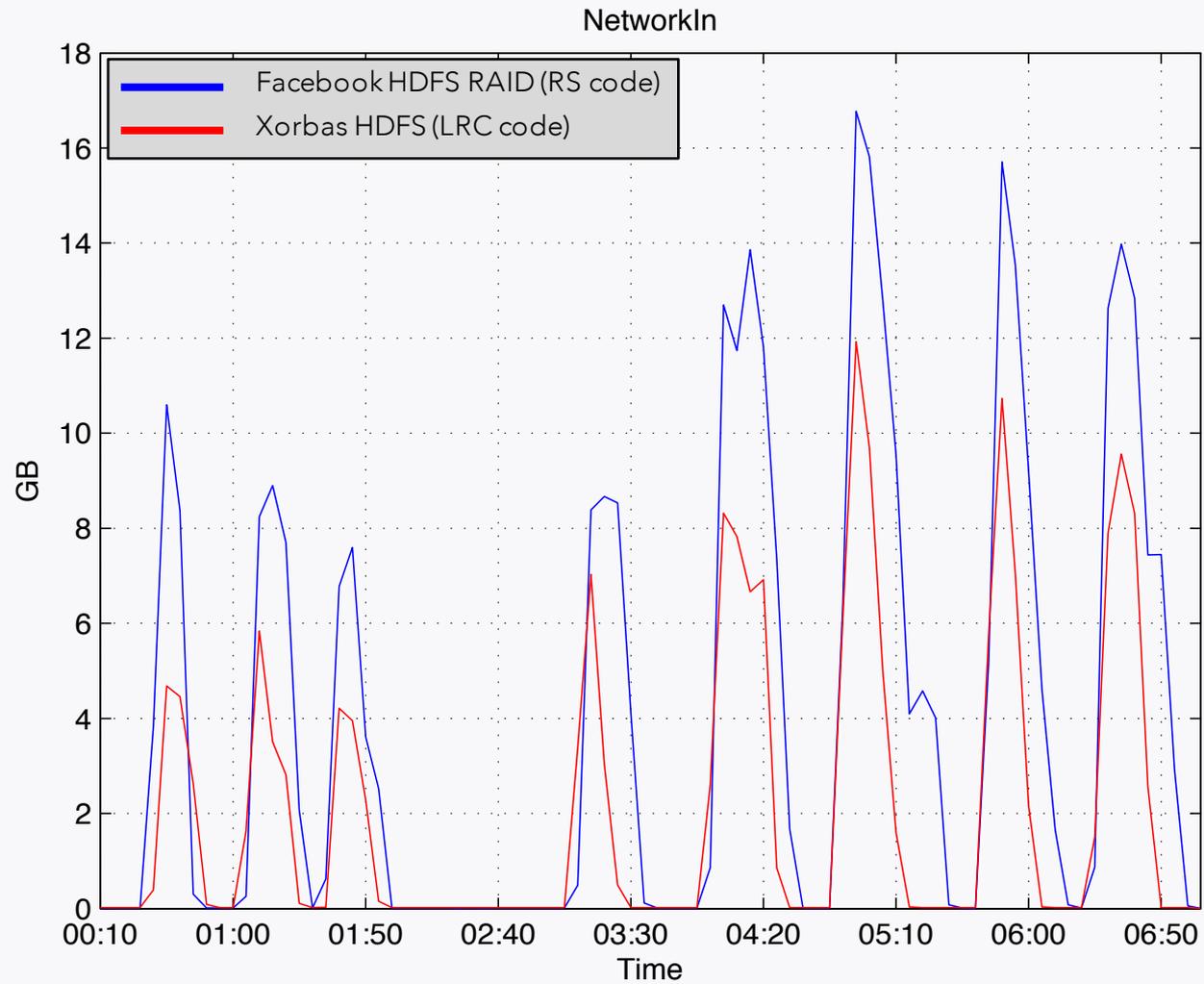
DataNode Health:

Live Nodes	37			
In Service		37		
Decommission: Completed		0		
Decommission: In Progress		0		
Dead Nodes	13			
Excluded		0		
Decommission: Completed			0	
Decommission: Not Completed			0	
Not Excluded		13		

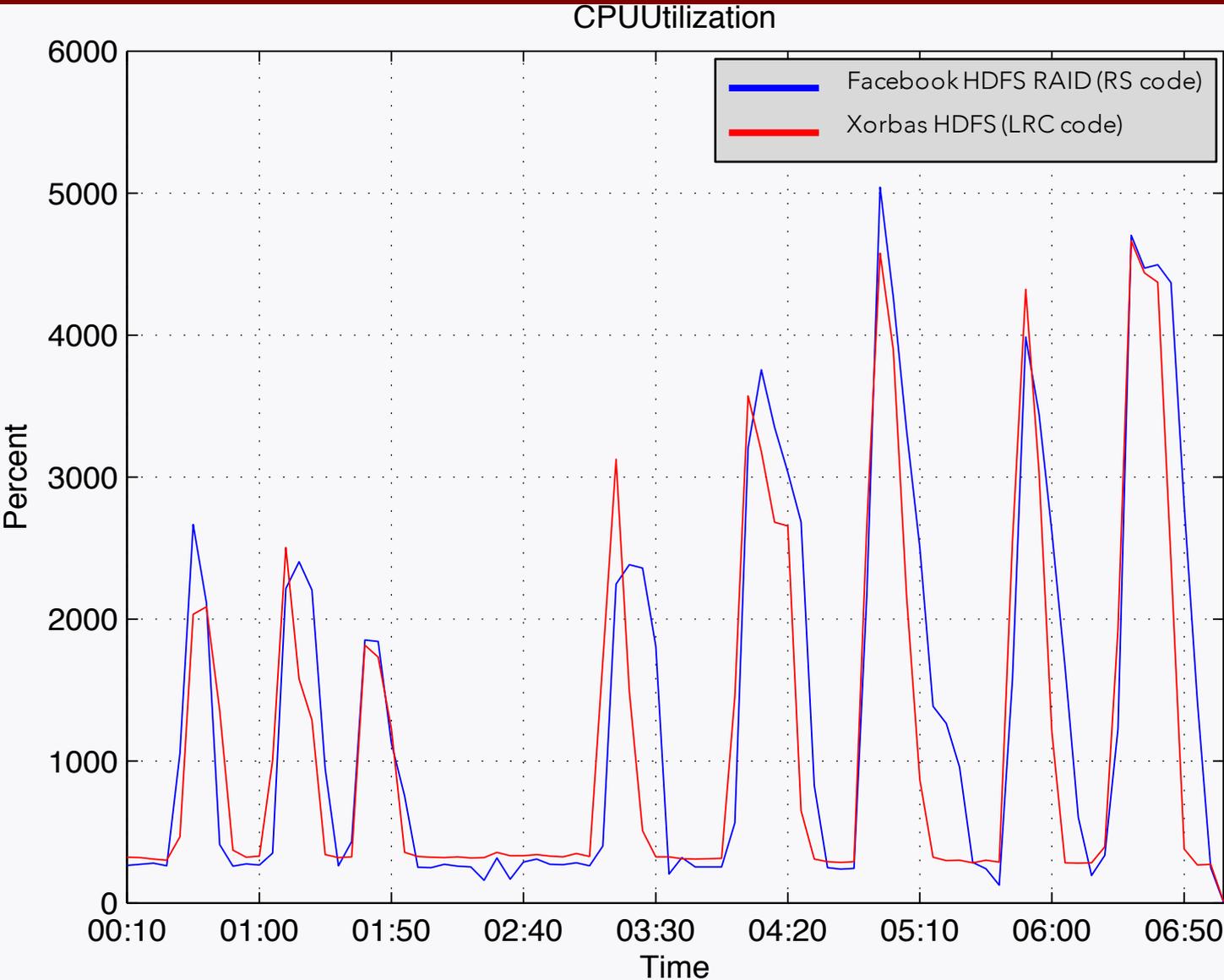
Some experiments

- 100 machines on Amazon ec2
- 50 machines running HDFS RAID (facebook version, (14,10) Reed Solomon code)
- 50 running our LRC code
- 50 files uploaded on system, 640MB per file
- Killing nodes and measuring network traffic, disk IO, CPU, etc during node repairs.

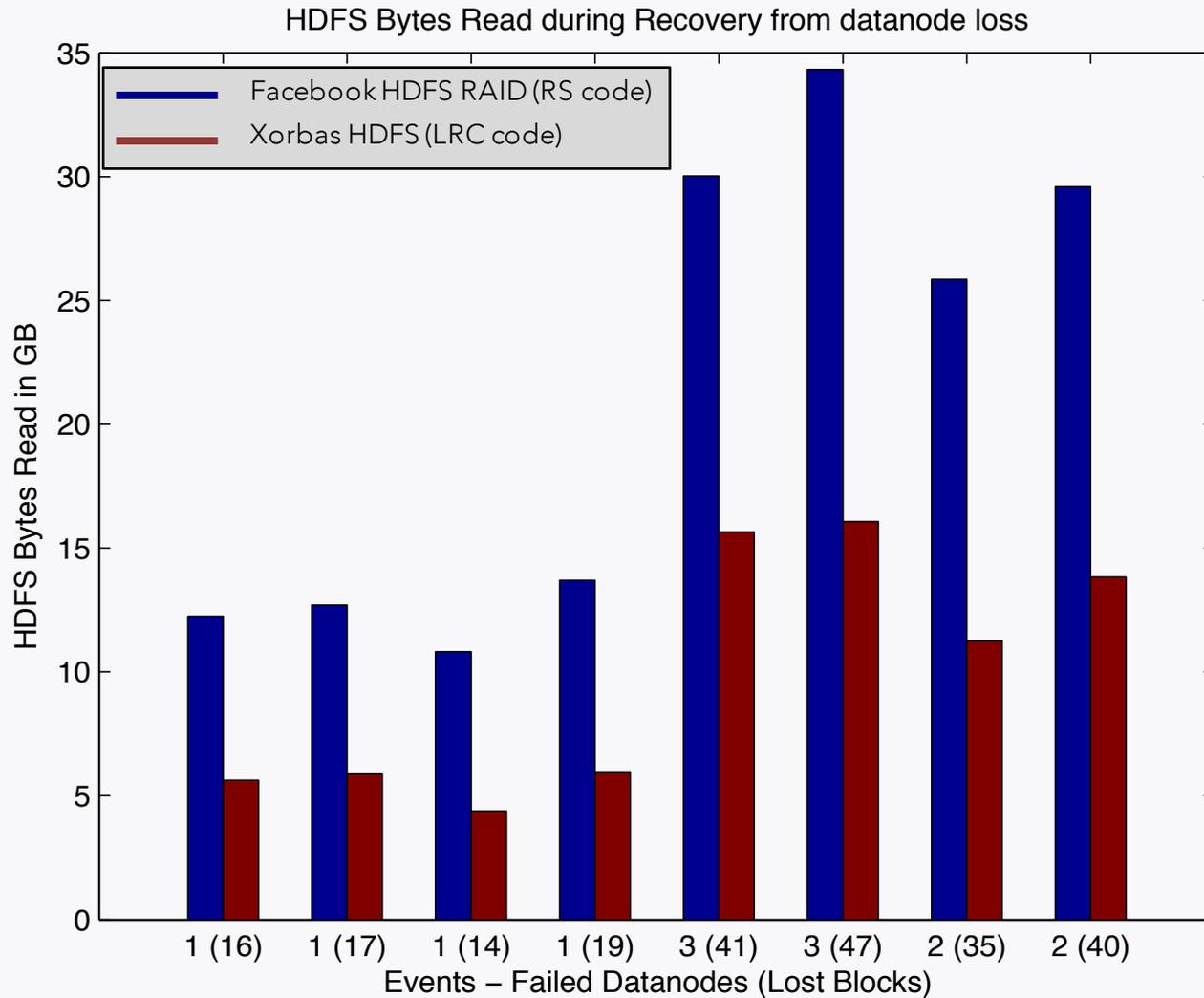
Repair Network traffic



CPU



Disk IO



what we observe

New storage code reduces bytes read by roughly 2.6x

Network bandwidth reduced by approximately 2x

We use 14% more storage. Similar CPU.

In several cases 30-40% faster repairs.

Provides four more zeros of data availability compared to replication

Gains can be much more significant if larger codes are used (i.e. for archival storage systems).

In some cases might be better to save storage, reduce repair bandwidth but lose in locality

(PiggyBack codes: Rashmi et al. USENIX HotStorage 2013)