Ming Hsieh Department of Electrical Engineering

Novel Codes for Cloud Storage

Maheswaran Sathiamoorthy*, Megasthenis Asteris*, Dimitris Papailiopoulos*, Alexandros G. Dimakis*, Ramkumar Vadali**, Scott Chen**, Dhruba Borthakur** * - USC, ** - Facebook



School of Engineering

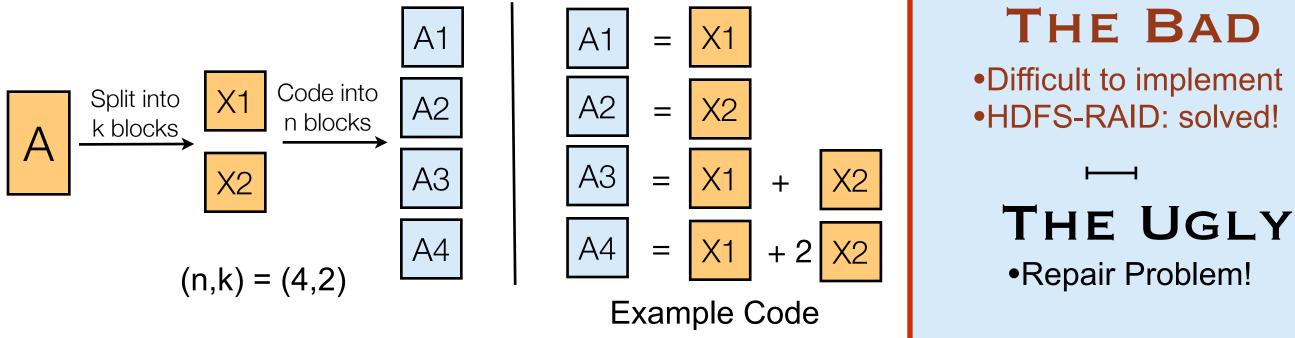
Cloud Storage

- Thousands of machines in data centers
- Failure is the norm rather than the exception
- To provide reliability, replication is used.
 - each file is stored multiple times (generally thrice).
 - replication costly, so could coding be an alternative?
- Hadoop: Distributed storage system, widely used.

Erasure Coding

(n,k) coding:

- •Split file in k blocks and then code into n blocks
- •Any k of the n blocks enough to recover the file
- •e.g.: Reed Solomon codes



THE GOOD

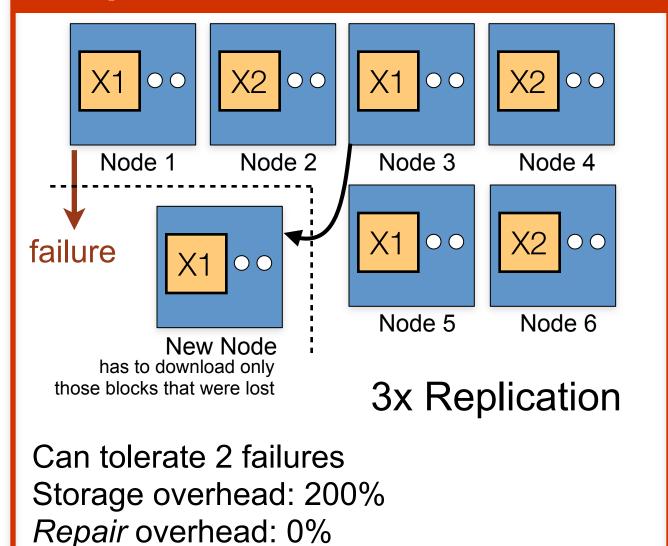
•Higher reliability

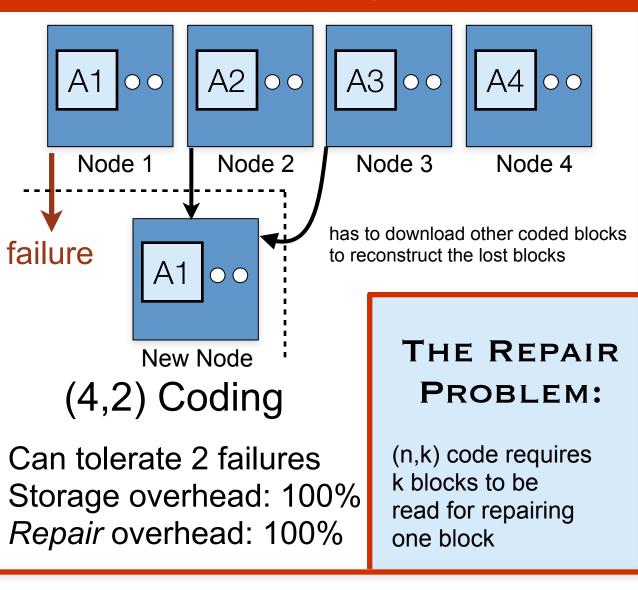
•Lower storage

Replication

Coded Storage

Facebook clusters

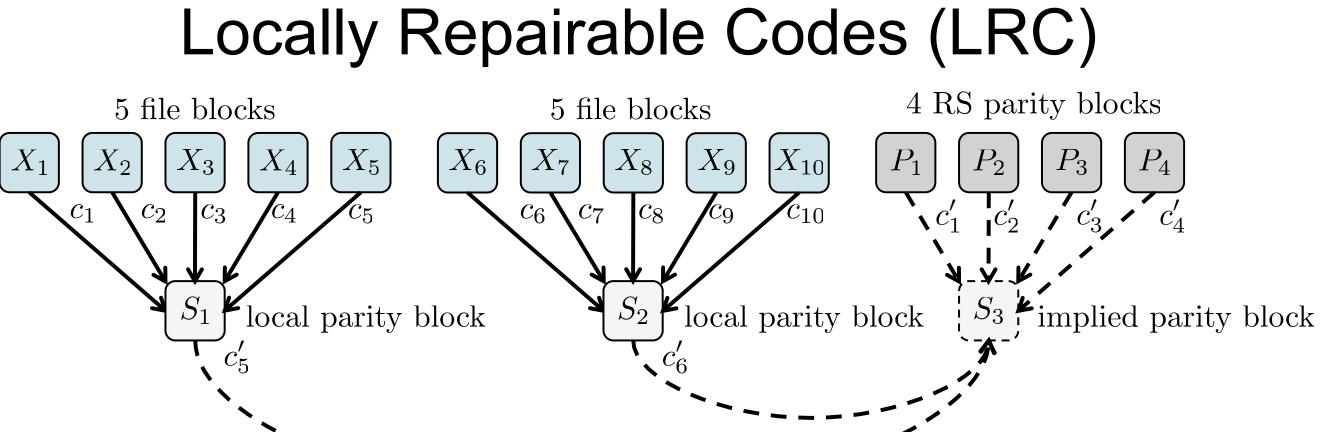




- Facebook has one of the largest Hadoop clusters
 - >3000 machines
 - >30 PB logical data
- 3x replication costly
- Facebook uses HDFS-RAID with (14, 10) Reed Solomon Code which we will call HDFS-RS.
- Every lost block needs 10 more blocks for repair!
- Network is bottlenecked
- So only 8% of cold data is encoded

Xorbas

- LRC Codes to the rescue!
 Facebook's HDFS-RAID is open source.
- We used it to develop our own version of Hadoop called HDFS-Xorbas*, which implements LRC codes



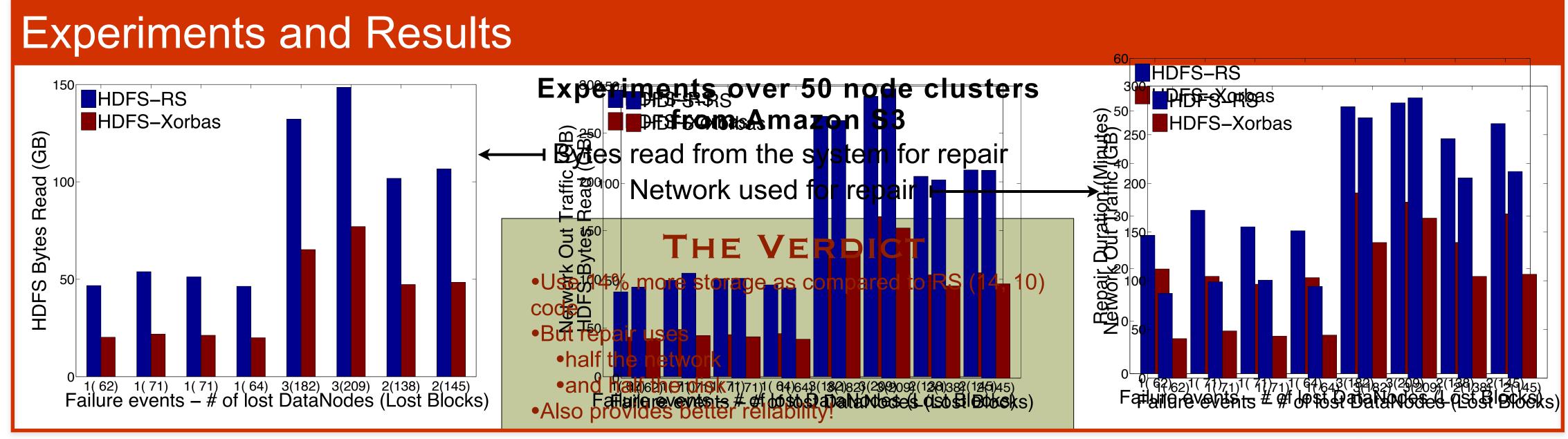
We use a (14, 10, 2) LRC Code
2/14 = 14% extra storage
but mitigates the repair problem

LRC (14,10,2) Recipe:

•Split file into 10 blocks

•Then create 4 extra parity blocks using Reed Solomon encoding

- •Use the first 5 and second 5 file blocks to create 2 "local parity blocks" S1 & S2
- •One lost block requires only 5 others blocks to be read.



Ming Hsieh Department of Electrical Engineering

<u>msathiam@usc.edu</u> *-code opensourced at <u>http://github.com/madiator/hadoop-20</u>