

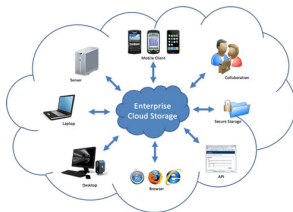
Distributed Storage Systems: Issues, Performance Analysis and Code Design

Vinay A. Vaishampayan

Outline

- Big Picture
 1. Overview
 2. What is repair, why is it important?
 3. Cost benefits.
- Deep Dive
 1. Reliability analysis.
 2. Generalized Singleton Bound:
 3. Code design.
 4. Benefits.
- Open Issues
- References

Cloud Storage



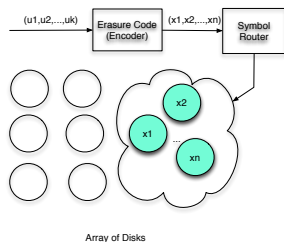
- Trend: 'outsourcing' enterprise and personal computing to the *cloud*. Only 1-2% of the \$25bn software market has moved to cloud. Growth forecasts are 'rosy'.
- **Potential Advantages:** Reduced IT costs. Elasticity. Vendors able to provide low cost storage due to economies of scale. Anywhere availability.
- **Disadvantages/Concerns/Roadblocks:** Security, privacy, reliability, availability, latency, hidden costs, corporate policy.

Economics

- Vendor solutions are expensive.
- Many enterprises consider a 'build our own approach'.
- Use commodity disks.
- but...Annual failure rates for most drives in the range 1% to 10%.
 - Seagate: 1.2 Million Hour MTBF for 1TB ES.2 drive (0.73% AFR)
 - More plausible... $\text{Prob}(\text{Disk will fail in 1 year})$ between 0.01 and 0.1.
- Reliability is a concern. Amazon S3: 11 nines of durability, 4 nines of availability over a given year.
- Translation: $\text{Prob}(\text{Data is lost}) = 10^{-11}$.
 $\text{Prob}(\text{Data is temporarily unavailable}) = 10^{-4}$. (in one year)
- Erasure Codes!!!

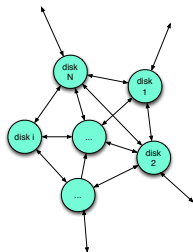
Erasure Codes in a Storage Application

Code Symbol Routing



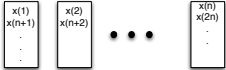
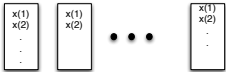
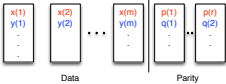
- (n, k) MDS code. Any $n - k$ erasures corrected

Traffic Types



- Ingress/egress traffic.
- Internal traffic for distribution and repair.

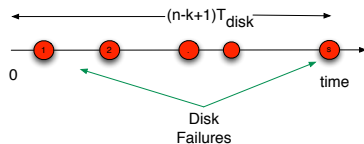
Legacy Solutions: Redundant Array of Independent Disks (RAID)

Type	Configuration	$\rho = 1/\alpha$	Fail Rate
0		1	$1 - (1 - p)^n$
1		$1/n$	p^n
4-6		$\frac{m}{m+r}$	$\sum_{l=r+1}^{r+m} \binom{m+r}{l} p^l (1-p)^{m+r-l}$

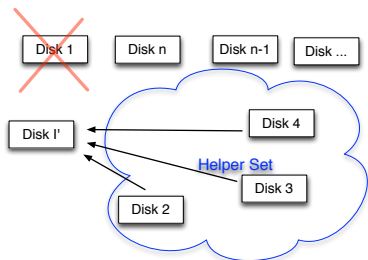
[Wikipedia on RAID Systems](#)

Some Facts and Terminology

Reliability doesn't improve much if single disk failures are not repaired.



Disk Repair



- Disk repair improves reliability.
- Repair time is a crucial parameter.
- Repair time determined by network topology and code characteristics.

Deep Dive

- Reliability analysis: Lower bound on *Reliability Function*

$$R(T) := Pr(\text{No Data Lost in } [0, T])$$

1. **Goal:** to understand dependence on repair duration.
 2. **Methods:** Jensen's inequality + Geometry.
- Code design
 1. Review of tradeoff between storage overhead and repair bandwidth.
 2. Constructions that reduce the repair duration compared to conventional MDS codes.

Error Analysis

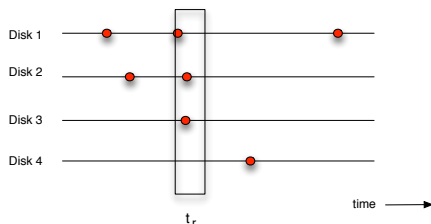
A. Campello and V. A. Vaishampayan, “Reliability of Erasure coded Storage Systems: A Geometric Approach” 2013 IEEE Intl. Conf. on Big Data, Santa Clara, CA, Oct. 2013.

Best and Worst Case Scenarios

When a disk fails repair data is drawn from a set of helper disks.

- **Best Case:** Repair data drawn from a helper disk does not depend on the helper set.
- **Worst Case:** Repair data drawn from a helper disk depends on the helper set.

Best Case Error Region



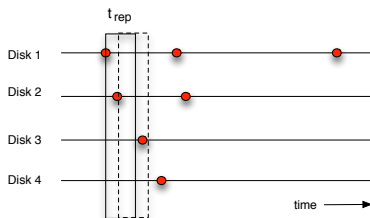
Sort \mathbf{x} in ascending order to get \mathbf{y} .

Theorem

In a best case scenario, a data loss event occurs in $[0, T)$ iff a vector of disk failures $(t_1, t_2, \dots, t_n) \in \mathcal{R}_b$ where

$$\mathcal{R}_b := \{\mathbf{x} = (x_1, x_2, \dots, x_n) \in [0, T)^n : y_{l+n-k} - y_l < t_r \text{ for some } l\}$$

Worst Case Error Region



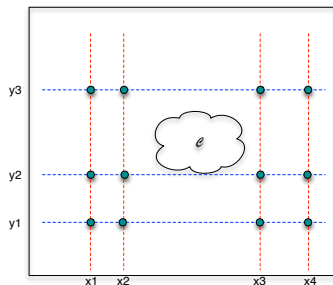
Sort \mathbf{x} in ascending order to get \mathbf{y} .

Theorem

In a worst case scenario, a data loss event occurs in $[0, T)$ iff a vector of disk failures $(t_1, t_2, \dots, t_n) \in \mathcal{R}_w$ where

$$\mathcal{R}_w := \{\mathbf{x} = (x_1, x_2, \dots, x_n) \in [0, T)^n : y_{l+1} - y_l < t_r, l = m, m+1, \dots, m+n-k-1, \text{ for some } m\}$$

Set Avoidance for Cartesian Product of Random Sets



Theorem

Let $\mathcal{X} := \{X_1, X_2, \dots, X_m\}$ and $\mathcal{Y} := \{Y_1, Y_2, \dots, Y_n\}$, where the X_i 's are i.i.d on a set \mathcal{S}_1 and the Y_i 's are i.i.d on a set \mathcal{S}_2 . Let X, Y be generic random variables distributed as X_i and Y_i , resp. Let $\mathcal{C} \subset \mathcal{S}_1 \times \mathcal{S}_2$. Then

$$P(\mathcal{X} \times \mathcal{Y} \cap \mathcal{C} = \emptyset) \geq \left(P(\{X\} \times \{Y\} \cap \mathcal{C} = \emptyset) \right)^{mn} \quad (1)$$

Bounding Reliability: The Set Avoidance Theorem

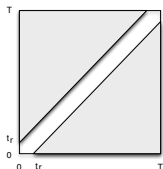
- \mathcal{X}_i : set of random failure instants for disk i .
- $|\mathcal{X}_i| = m_i$.
- $m_i > 0$.
- $Pr[NoDataLoss] = Pr(\mathcal{X}_1 \times \mathcal{X}_2 \times \dots \times \mathcal{X}_n \cap \mathcal{R} = \emptyset)$.

- $$Pr[NoDataLoss] \geq \left(1 - \frac{vol\mathcal{R}}{T^n}\right)^{m_1 m_2 \dots m_n}$$

Example Error Regions

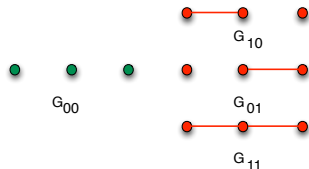
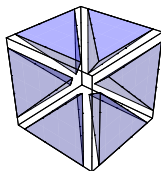
(3, 2) Single Parity Code

(2, 1) Single Parity Code



$$\mathcal{R}_2^C = \{\mathbf{x} : y_2 - y_1 > t_r\}$$

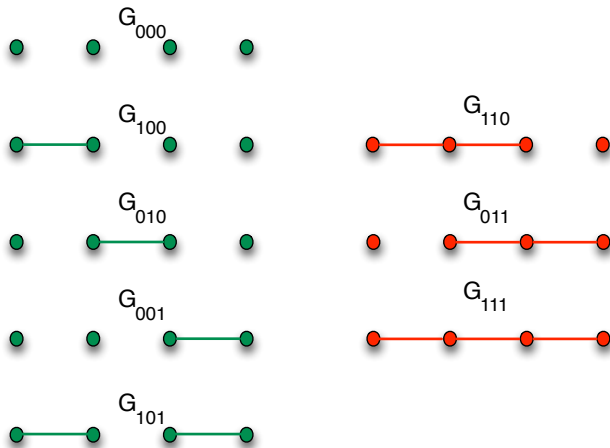
$$\text{vol } \mathcal{R}_2^C = (T - t_r)^2$$



$$\mathcal{R}_3^C = \{\mathbf{x} \in (0, t)^3 : y_2 - y_1 > t_r, y_3 - y_2 > t_r\}$$

$$\text{vol } \mathcal{R}_3^C = (T - 2t_r)^3$$

Graphs for the (4, 2) Code: Worst Case



The Supergraph I

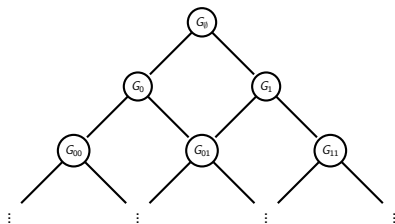


Figure: Representation of the set of graphs $G_{0^i 1^{j-i}}$

$$\text{vol } G_{0^i} = \frac{(T - it_{rep})^n}{n!}.$$

$$\text{vol } G_{0^i 1^{j-i}} = \sum_{l=0}^j \binom{j}{l} (-1)^{j+l} \text{vol } G_{0^{i+j-l}}.$$

Calculations for Non-Trivial Codes

Code	$\text{vol } \mathcal{R}_w^c$
(4, 2)	$-24T^2t_{rep}^2 + 72Tt_{rep}^3 - 64t_{rep}^4 + T^4$
(5, 2)	$-120T^2t_{rep}^3 + 480Tt_{rep}^4 - 540t_{rep}^5 + T^5$
(5, 3)	$-60T^3t_{rep}^2 + 300T^2t_{rep}^3 - 570Tt_{rep}^4 + 390t_{rep}^5 + T^5$
(6, 2)	$-720T^2t_{rep}^4 + 3600Tt_{rep}^5 - 4920t_{rep}^6 + T^6$
(6, 3)	$-360T^3t_{rep}^3 + 2340T^2t_{rep}^4 - 5580Tt_{rep}^5 + 4740t_{rep}^6 + T^6$
(6, 4)	$-120T^4t_{rep}^2 + 840T^3t_{rep}^3 - 2100T^2t_{rep}^4 + 1260Tt_{rep}^5 + 1492t_{rep}^6 + T^6$

$$\Pr(\text{DataLoss}) \leq \frac{n!}{(k-1)!} (\lambda T)^{n-k+1} \left(\frac{t_r}{T}\right)^{n-k}$$

Singleton-Like Bounds

A. Dimakis, P. Godfrey, Y. Wu, M. Wainwright, K. Ramchandran,
“Network Coding for Distributed Storage Systems”, IEEE Trans. In. Th.,
Sept. 2010.

Tradeoff Between Storage and Repair: Setup

Code Definition

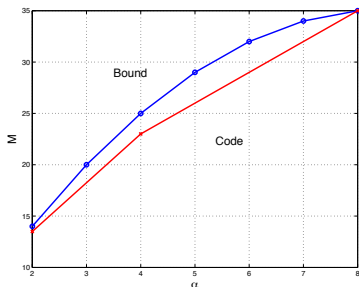
- (i) Code can correct up to $n - k$ erasures in any $n - k$ positions.
- (ii) α code symbols are placed on a single disk.
- (ii) Any single erased code symbol can be reconstructed using β symbols/disk from d helper disks, $d \leq n - 1$.

- M : Number of information symbols.
- Key question: How does M depend on α, β ?
- For a standard (n, k) MDS code
 1. $\beta = k\alpha/d$.
 2. $M = k\alpha$ information symbols can be stored.

Generalized Singleton Bound

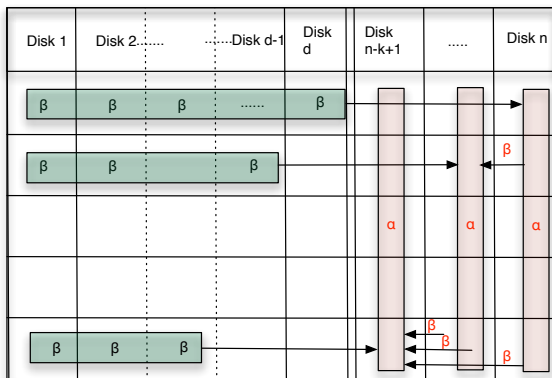
$$\begin{array}{rcccccc}
 d\beta & +(d-1)\beta & +\dots & +\dots & +(d-k+1)\beta & \geq M \\
 \alpha & +(d-1)\beta & +\dots & +\dots & +(d-k+1)\beta & \geq M \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \geq M \\
 \alpha & +\alpha & +\dots & \alpha & +(d-k+1)\beta & \geq M \\
 \alpha & +\alpha & +\dots & \alpha & +\alpha & \geq M
 \end{array}$$

- MSR: $\alpha = (d - k + 1)\beta$ (last two constraints active)
- MBR: $\alpha = d\beta$ (first two constraints active)

Repair/Storage Tradeoff: M vs. α for fixed β 

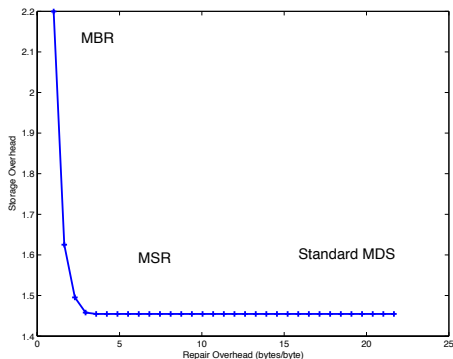
$$n = 9, k = 7, d = 8, \beta = 1$$

Proof: Generalized Singleton Bound



- **Basic Idea:** What is the minimum amount of information required to *fix* a codeword?

Storage Repair-Bandwidth Tradeoff



- $n = 32, k = 22, d = n - 1, \alpha = 48$.
- Repair Overhead: $\beta * d / \alpha$
- Storage Overhead: $(k\alpha) / M$
- Repair time can be reduced by a factor of 7.

Code Design at MSR and MBR points

K. Rashmi, N. Shah, and P. Kumar, Optimal Exact-Regenerating Codes for Distributed Storage at the MSR and MBR Points via a Product-Matrix Construction, IEEE Transactions on Information Theory, vol. 57, no. 8, pp. 5227-5239, Aug. 2011.

Do points superior to the time-sharing bound exist?

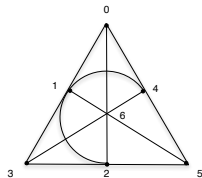
C. Tian, V. Aggarwal and V. A. Vaishampayan, "Exact Repair Regenerating Codes via Layered Erasure Correction and Block Designs, ISIT, 2013.

- Two-layer construction + symbol distribution strategy based on t -designs.

t - Designs

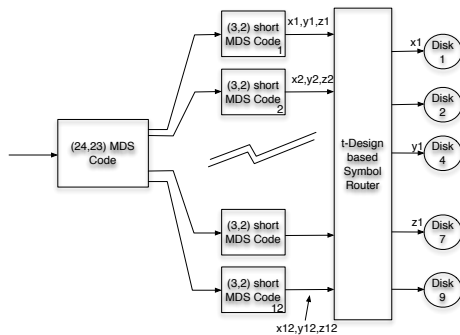
- $\mathcal{X} = \{0, 1, \dots, n - 1\}$.
- Block: r -subset of \mathcal{X} .
- Design is a collection of Blocks.
- $t - (n, r, \lambda)$ Design: Any t -subset of \mathcal{X} is contained in exactly λ blocks.
- Block Designs: $t = 2$
- Incomplete Design: $r < n$
- Steiner Systems: $\lambda = 1$.

$2 - (7, 3, 1)$ Block Design or Steiner System



- $\mathcal{X} = \{0, 1, \dots, 6\}$
- $r = 3$
- $\{\{0, 1, 3\}, \{2, 3, 5\}, \{0, 4, 5\}, \{1, 5, 6\}, \{3, 4, 6\}, \{0, 2, 6\}, \{1, 2, 4\}\}$
- $t = 2, \lambda = 1$. Any 2-subset is contained in exactly 1 block.

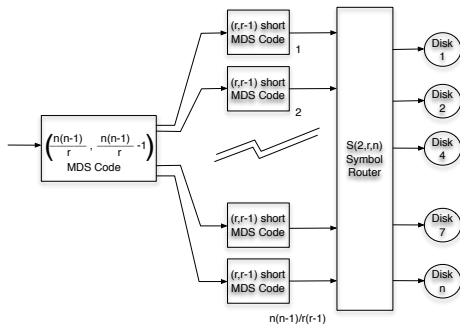
A Multi-Level Repair Efficient Code



t -design: 2 – (9, 3, 1) has 12 blocks.

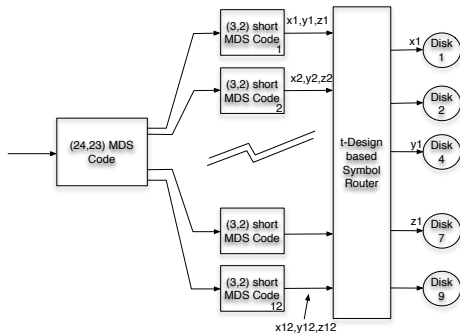
$\{\{1, 4, 7\}, \{2, 5, 8\}, \{3, 6, 9\}, \{1, 6, 8\}, \{2, 4, 9\}, \{3, 5, 7\},$
 $\{1, 5, 9\}, \{2, 6, 7\}, \{3, 4, 8\}, \{1, 2, 3\}, \{4, 5, 6\}, \{7, 8, 9\}\}$

General Construction



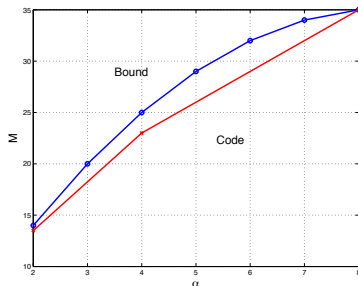
$$\alpha = \frac{n-1}{r-1}; \beta = \frac{\alpha(r-1)}{n-1} = 1$$

Why Use a Steiner System?



If two disks fail at most one short parity group will lose two symbols.

Repair/Storage Tradeoff: Vary r



$$n = 9, k = 7, d = 8, \beta = 1$$

Code achieves $M = 13.4$ and $\alpha = 2$ corresp. to $r = 5$

Code achieves $M = 35$ and $\alpha = 8$ corresp. to $r = 2$

Summary

- Erasure Coding to improve reliability of distributed storage systems.
- Reliability depends on error correction capability of the code and also the repair bandwidth.
- Generalized Singleton bound: storage/ repair bandwidth tradeoff.
- Exhibited low-repair-bandwidth code based on two-layer code construction and block designs.

Thank You!

References

- [1] A. Campello and V. A. Vaishampayan, “Reliability of Erasure coded Storage Systems: A Geometric Approach” 2013 IEEE Intl. Conf. on Big Data, Santa Clara, CA, Oct. 2013.
- [2] C. Tian, V. Aggarwal and V. A. Vaishampayan, “Exact Repair Regenerating Codes via Layered Erasure Correction and Block Designs, under review.
- [3] A. Dimakis, P. Godfrey, Y. Wu, M. Wainwright, K. Ramchandran, “Network Coding for Distributed Storage Systems”, IEEE Trans. In. Th., Sept. 2010.
- [4] K. Rashmi, N. Shah, and P. Kumar, Optimal Exact-Regenerating Codes for Distributed Storage at the MSR and MBR Points via a Product-Matrix Construction, IEEE Transactions on Information Theory, vol. 57, no. 8, pp. 5227-5239, Aug. 2011.