# Erasure Coding in Distributed Storage Systems: Comprehensive Reference Document

Chiradip Mandal

July 2025

**Abstract**

This document provides a comprehensive guide to erasure coding in distributed storage systems. We examine the fundamental (n,k) parameters, mathematical foundations, performance trade-offs, and practical applications. Crucially, we explain why erasure coding excels for cold storage but fails for primary/hot storage due to severe performance penalties.

# Contents

# 1 Understanding (n,k) Parameters

## 1.1 Core Concept

An **(n,k) erasure code** divides original data into:

- **k data chunks** - pieces of the original file

- **(n-k) parity chunks** - redundancy for fault tolerance

- **n total chunks** stored across the system

**Golden Rule:** Any k chunks out of n total can reconstruct the original data.

## 1.2 Example: (6,4) Configuration

Consider a 1GB file with (6,4) erasure coding:



Figure 1: (6,4) Erasure coding structure

**Key Properties:**

- **Storage Overhead:** $n/k = 6/4 = 1.5\times$ (vs $3\times$ for replication)

- **Fault Tolerance:** Can lose any 2 chunks and still recover

- **Storage Efficiency:** $k/n = 4/6 = 67\%$

## 1.3 Common Configurations

Table 1: Erasure coding configurations in production systems

| Config | Data | Parity | Overhead | Faults | Used By |
|---|---|---|---|---|---|
| (3,2) | 2 | 1 | $1.5\times$ | 1 | Simple systems |
| (6,4) | 4 | 2 | $1.5\times$ | 2 | HDFS, Cassandra |
| (9,6) | 6 | 3 | $1.5\times$ | 3 | Facebook |
| (10,6) | 6 | 4 | $1.67\times$ | 4 | Google Cloud |
| (12,8) | 8 | 4 | $1.5\times$ | 4 | Amazon S3 IA |
| (14,10) | 10 | 4 | $1.4\times$ | 4 | S3 Glacier |
| **3× Replication** | 1 | 2 copies | $3.0\times$ | 2 | Traditional |

# 2 Mathematical Foundations

## 2.1 Reed-Solomon Codes

Reed-Solomon codes represent data as polynomials over finite fields $\mathrm{GF}(2^w)$:

$$f(x) = d_0 + d_1 x + d_2 x^2 + \ldots + d_{k-1} x^{k-1} \tag{1}$$

Where $d_i$ are the data chunks.
**Encoding Process:**

1. Create polynomial $f(x)$ with data chunks as coefficients

2. Evaluate at $n$ distinct points: $c_i = f(\alpha_i)$

3. First $k$ evaluations are data chunks (systematic encoding)

4. Remaining $(n - k)$ evaluations are parity chunks

**Decoding Process:**

1. Given any $k$ chunks, reconstruct polynomial using Lagrange interpolation

2. Extract original data chunks from polynomial coefficients

$$f(x) \quad f(x) = d_0 + d_1 x + d_2 x^2 + d_3 x^3$$

$$D_0 \quad D_1 \quad D_2 \quad D_3 \quad P_0 \quad P_1 \qquad x$$

Figure 2: Reed-Solomon polynomial evaluation

## 2.2 Matrix Representation

Encoding can be expressed as matrix multiplication:

$$\mathbf{c} = \mathbf{d} \cdot \mathbf{G} \tag{2}$$

Where $\mathbf{G}$ is the generator matrix:

$$\mathbf{G} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \\ 1 & \alpha_k & \cdots & \alpha_k^{k-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \alpha_{n-1} & \cdots & \alpha_{n-1}^{k-1} \end{bmatrix} \tag{3}$$

## 2.3 Finite Field Arithmetic

Operations in $\text{GF}(2^w)$:

- **Addition:** XOR operation ($a \oplus b$)

- **Multiplication:** Polynomial multiplication modulo irreducible polynomial

- **Division:** Multiplication by multiplicative inverse

# 3 The Performance Problem

## 3.1 Why NOT Primary Storage?

**Critical Insight:** Erasure coding trades performance for storage efficiency.

### 3.1.1 Read Performance Penalty

Table 2: Read operation comparison

| Approach | Process | Latency |
|---|---|---|
| **Replication** | Read 1 file from 1 node | **5ms** |
| **Erasure Coding** | Contact k nodes + decode | **25ms** |
| **Performance Impact** | **5× slower** | |

**Why so slow?**

- **Network:** Must contact k nodes instead of 1

- **CPU:** Finite field arithmetic for decoding

- **Coordination:** Wait for k responses

- **Memory:** Buffer and process multiple chunks

**Replication:** Request → Read 1 copy → Done 5ms

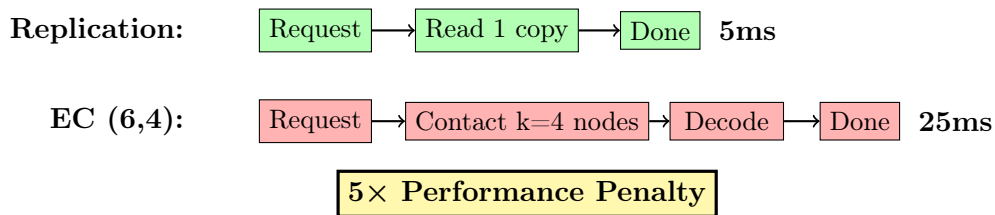**EC (6,4):** Request → Contact k=4 nodes → Decode → Done 25ms

5× Performance Penalty

Figure 3: Read latency comparison: Replication vs Erasure Coding

### 3.1.2 Write Amplification Nightmare

**Worst Case Scenario:** Updating 1 byte in an erasure-coded file

---

**Example: (6,4) EC, change 1 byte in chunk $D_1$**
**Operations Required:**

1. **Read:** All k=4 data chunks (to recalculate parity)

2. **Compute:** New parity chunks (CPU intensive)

3. **Write:** Modified data chunk + all parity chunks

**Total:** 4 reads + 3 writes + encoding CPU
**vs Replication:** 1 read + 3 writes
**Result: 20× slower for small updates!**

---

### 3.1.3 Comprehensive Performance Impact

Table 3: Performance comparison across operation types

| Operation | Replication | EC (6,4) | Impact |
|---|---|---|---|
| Sequential Read | 5ms | 15ms | 3× slower |
| Random Read | 2ms | 25ms | 12× slower |
| Large Sequential Write | 100ms | 150ms | 1.5× slower |
| Small Random Write | 10ms | 200ms | 20× slower |
| Node Failure Recovery | Instant | 5 minutes | 300× slower |
| CPU Usage | Low | High | 10× more |
| Network Bandwidth | 1× | k× | k× amplification |

## 3.2 Real-World Impact Examples

### 3.2.1 Database Query Scenario

- **Replication:** "Get user profile #12345" → 5ms response

- **Erasure Coding:** Same query → 25ms response

- **Impact:** 5× slower for every database query!

### 3.2.2 High-Traffic Application

At 1000 requests/second with 5× read penalty:

- **Additional delay per request:** 20ms

- **Total extra waiting time:** 20 seconds per second!

- **Result:** Completely unacceptable user experience

# 4 Where Erasure Coding Excels

## 4.1 The Tiered Storage Solution

**Access Pattern Rule:** If data is accessed less than once per day, storage savings outweigh performance costs.

Table 4: Strategic storage tier recommendations

| Tier | Data Type | Access | Strategy | Rationale |
|------|-----------|--------|----------|-----------|
| **HOT** | Databases, Active files | Daily | Replication 3× | Performance critical |
| **WARM** | Logs, Cached data | Weekly | Light EC (6,4) | Balanced approach |
| **COLD** | Backups, Analytics | Monthly | Heavy EC (14,10) | Storage efficiency |
| **ARCHIVE** | Compliance, DR | Yearly | Max EC (20,16) | Cost optimization |



Figure 4: Strategic tiered storage architecture

## 4.2 Object Storage Success Stories

### 4.2.1 Amazon S3

- **Standard:** Uses replication for frequently accessed data

- **Standard-IA:** (12,8) erasure coding for infrequent access

- **Glacier:** (14,10) for long-term archival

- **Deep Archive:** Maximum erasure coding for lowest cost

### 4.2.2 Google Cloud Storage

- **Standard:** Replication for hot data

- **Nearline:** (12,8) for monthly access patterns

- **Coldline:** (14,10) for quarterly access

- **Archive:** Heavy erasure coding for yearly access

## 4.3 Database Applications

### 4.3.1 Apache Cassandra

- **Hot data:** Replication for real-time queries

- **Historical data:** (6,4) erasure coding

- **Analytics tables:** (10,6) for data warehouse workloads

### 4.3.2   Hadoop HDFS-EC

- **Active datasets:** Traditional 3× replication

- **Aging data:** Automatic migration to (6,4) EC

- **Archive data:** (10,6) for maximum efficiency

# 5   Implementation Considerations

## 5.1   Hardware Acceleration

Modern erasure coding implementations leverage:

Table 5: Hardware acceleration performance

| Technology | Speedup | Throughput |
|---|---|---|
| Software (single thread) | 1× | 0.1-0.5 GB/s |
| SIMD (AVX-512, NEON) | 8-16× | 1-2 GB/s |
| GPU (CUDA/OpenCL) | 50-100× | 5-15 GB/s |
| FPGA/ASIC | 200× | 20-50 GB/s |

**Intel ISA-L Library:** Production-ready implementation with optimized assembly routines, providing 10-20× speedup over naive implementations.

## 5.2   Network Topology Optimization

**Placement Strategies:**

- **Rack diversity:** Distribute chunks across racks

- **Datacenter placement:** Geographic distribution for disasters

- **Network awareness:** Minimize cross-datacenter traffic

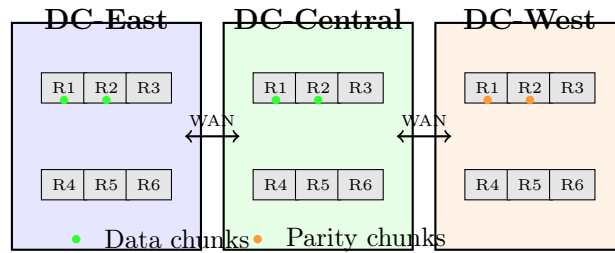- **Hierarchical domains:** Server → Rack → Datacenter fault isolation



Figure 5: Geographic distribution with rack-aware placement

## 5.3   Configuration Selection Guidelines

# 6   Advanced Techniques

## 6.1   Locally Repairable Codes (LRC)

**Problem:** Traditional RS codes require reading k chunks to repair any single failure.

Table 6: Configuration selection matrix

| Access Pattern | Latency Req. | Config | Overhead |
|---|---|---|---|
| Real-time (ms) | <10ms | Replication | 3.0× |
| Interactive (10-100ms) | <100ms | (6,4) EC | 1.5× |
| Batch (seconds) | <10s | (10,6) EC | 1.67× |
| Archive (minutes) | <600s | (20,16) EC | 1.25× |

**LRC Solution:** Add local parity chunks that enable repair from fewer chunks.
**Example LRC(12,8,4):**

- 8 data chunks in 2 groups of 4

- 2 local parity chunks (1 per group)

- 2 global parity chunks

- **Repair bandwidth:** 4 chunks instead of 8 (50% reduction)



Figure 6: LRC vs Traditional Reed-Solomon repair

**Benefits:**

- 2-4× reduction in repair bandwidth

- Faster single-failure recovery

- Lower network traffic during rebuilds

- Used in Azure Storage and Facebook f4

## 6.2 Adaptive Erasure Coding

**Machine Learning Integration:**

- **Workload analysis:** Automatic tier selection based on access patterns

- **Predictive maintenance:** Failure prediction for proactive repair

- **Dynamic tuning:** Parameter adjustment based on system load

- **Cost optimization:** Multi-dimensional optimization algorithms

# 7 Economic Analysis

## 7.1 Total Cost of Ownership

**Key Insights:**

Table 7: TCO breakdown by storage tier (relative costs)

| Cost Component | Hot | Warm | Cold | Archive |
|---|---|---|---|---|
| Storage Hardware | 100 | 60 | 25 | 15 |
| Performance Premium | 50 | 30 | 10 | 5 |
| Management Overhead | 20 | 15 | 35 | 50 |
| Network Bandwidth | 30 | 20 | 10 | 5 |
| **Total Cost/TB** | **200** | **125** | **80** | **75** |

- **Hot tier:** Performance premium justifies 3× replication cost

- **Cold tier:** Storage savings outweigh performance penalty

- **Archive:** Maximum efficiency for rarely accessed data

## 7.2 Break-Even Analysis

**Rule of Thumb:** Erasure coding becomes cost-effective when:

- Access frequency < 1 per day

- Storage cost > 60% of total TCO

- Performance SLA allows >50ms latency

- Data size > 1GB per object

# 8 Implementation Checklist

## 8.1 Phase 1: Planning

☐ Analyze workload patterns and access frequencies

☐ Define performance and durability requirements

☐ Model storage costs vs. performance trade-offs

☐ Assess network topology and bandwidth constraints

☐ Identify appropriate failure domains

☐ Select target (n,k) configurations

## 8.2 Phase 2: Implementation

☐ Choose erasure coding library (Intel ISA-L recommended)

☐ Implement chunk placement algorithms

☐ Set up monitoring and alerting systems

☐ Create automated repair procedures

☐ Test failure scenarios thoroughly

☐ Validate performance benchmarks

## 8.3   Phase 3: Operations

☐ Monitor repair bandwidth and recovery times

☐ Track storage efficiency metrics

☐ Optimize placement policies based on usage

☐ Regular testing of failure recovery procedures

☐ Capacity planning for data growth

☐ Performance tuning and optimization

# 9   Future Directions

## 9.1   Emerging Research Areas

**Next-Generation Codes:**

- **Quantum-resistant codes:** Post-quantum cryptography integration

- **Streaming codes:** Real-time encoding for live data

- **Regenerating codes:** Optimal repair bandwidth minimization

- **Polar codes:** Alternative to Reed-Solomon for specific applications

**Technology Integration:**

- **Edge computing:** Low-bandwidth optimized codes

- **5G networks:** Ultra-low latency requirements

- **Computational storage:** Near-data processing

- **Blockchain integration:** Decentralized storage systems

# 10 Decision Framework

## 10.1 Quick Decision Matrix

**Should I Use Erasure Coding?**
**YES, if:**

- Data accessed less than daily

- Storage cost is primary concern

- Can tolerate 10-100ms latency

- Large objects (>100MB)

- Archival or backup use case

**NO, if:**

- Real-time or interactive applications

- Database primary storage

- Frequent small updates

- Sub-10ms latency requirements

- Small objects (<1MB)

# 11 Conclusion

## 11.1 Key Takeaways

1. **Storage Efficiency Champion:** 50-70% savings over replication

2. **Performance Trade-off:** 3-25$\times$ slower operations

3. **Strategic Application:** Perfect for cold data, problematic for hot data

4. **Tiered Architecture:** Combine replication (hot) + EC (cold) for optimal results

5. **Industry Proven:** Used by all major cloud providers for appropriate workloads

## 11.2 Strategic Recommendations

Table 8: Final recommendations by use case

| Use Case | Strategy | Configuration |
|---|---|---|
| Production databases | Replication | 3$\times$ copies |
| Application logs | Hybrid | Replication $\to$ (6,4) EC |
| Data backups | Heavy EC | (14,10) |
| Compliance archives | Maximum EC | (20,16) |
| Content delivery | Mixed | Hot: replication, Cold: EC |
| Scientific datasets | Tiered EC | (6,4) $\to$ (14,10) |

## 11.3   The Bottom Line

**Erasure coding is transformative technology for the right use cases.**

It enables massive storage savings while maintaining high durability, but only when applied strategically. The key to success is understanding the fundamental trade-off between storage efficiency and performance, then architecting systems that leverage erasure coding's strengths while avoiding its weaknesses.

✓ **Advantages**          ✗ **Disadvantages**

- 50-70% storage savings
- Configurable fault tolerance
- Mathematical guarantees
- Geographic distribution
- Industry proven
- Scales with growth

vs

- 3-25× slower reads
- Write amplification
- High CPU overhead
- Complex recovery
- Network amplification
- Operational complexity

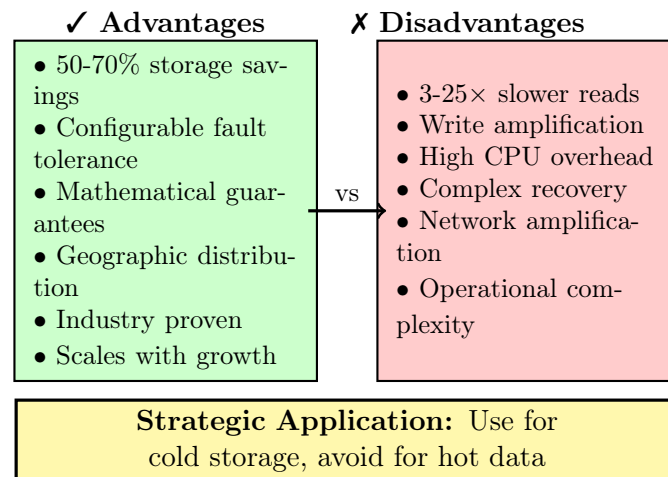**Strategic Application:** Use for cold storage, avoid for hot data

Figure 7: Erasure coding trade-offs summary

For modern distributed storage systems, the optimal approach is **tiered storage** that uses replication for hot data requiring fast access, and erasure coding for cold data where storage efficiency matters more than performance. This hybrid strategy maximizes both cost savings and system performance.

**Remember:** Erasure coding doesn't replace replication—it complements it in a well-designed storage hierarchy.