# MarFS and Multi-Tier Erasure
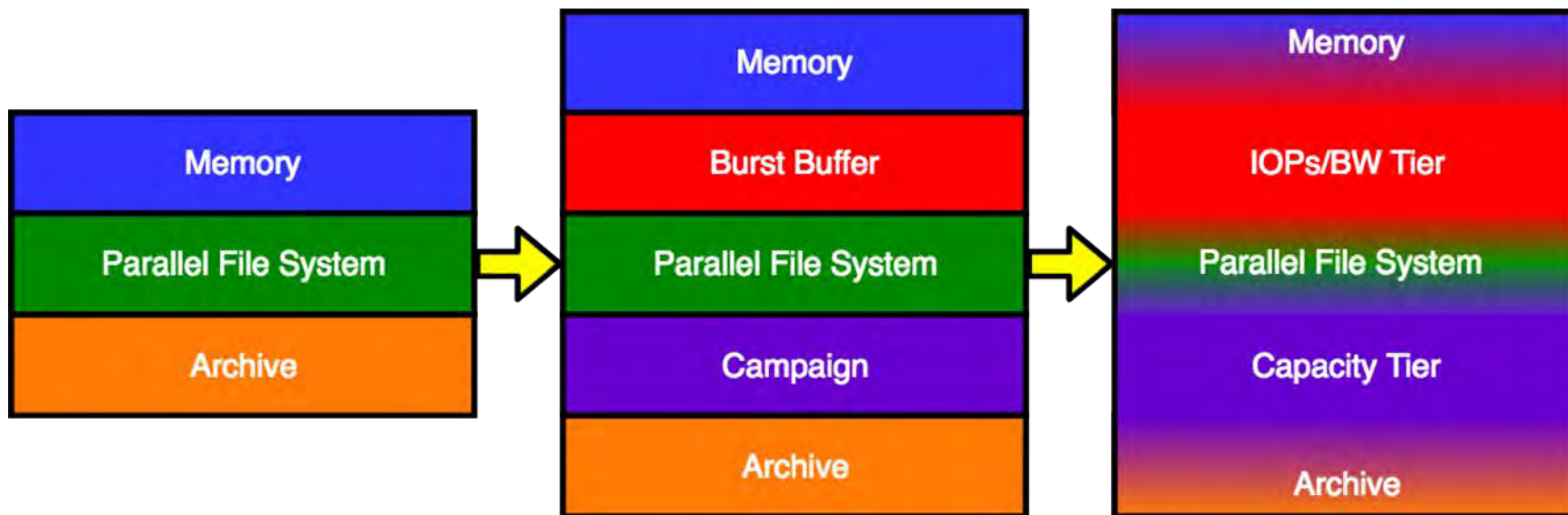
**Garrett Ransom (HPC-SYS)**

05/15/18

# HPC Storage Trends at LANL

# HPC Storage Trends at LANL – Capacity Tier

- **Data sets are growing faster than long-term storage can support**
  - Trinity: 2PB of memory / 4PB of flash
  - Crossroads: (maybe) 4PB mem / 10-100PB flash
  - HPSS Archive ~60PB Total, near capacity
- **Bandwidth of archive is a limiting factor**
  - Usable bandwidth of traditional archive (HPSS) ingest is roughly 3 GB/sec
  - HPSS ingests data much faster than it recalls
  - Storing petabytes of job checkpoints is infeasible
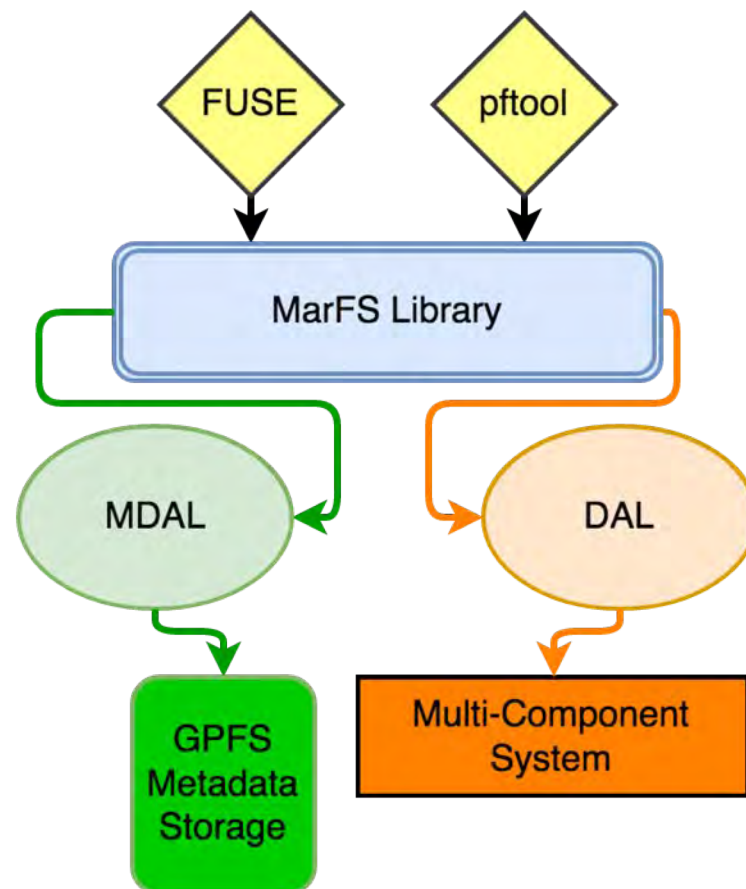
# Implementing a Capacity Tier

- **Tape is likely not the approach to take**
  - Tape is effective for truly cold data, not data sets that require periodic recall
  - Designing tape storage solutions is complex
- **Object Storage seems promising**
  - Flat namespace allows easy scalability
  - Erasure coding allows for cheaper disk media
- **Object Storage has limitations**
  - Machines love object-IDs, people generally don't
  - Potentially billions $ in applications expecting 'POSIX-like' file trees

# What is MarFS?

- **A near-POSIX interface layered over distinct metadata and data implementations**

  - Data stored as erasure coded objects

  - Metadata mirrored within a parallel file system (GPFS)

  - Object IDs stored as extended attributes of metadata files

- **Familiar semantics, fast metadata, stable objects**

  - Storing metadata to a real PLFS gives us POSIX-style directory trees and permissions almost for free

  - Storing data as objects simplifies implementation and data protection

- **With tradeoffs, of course**

  - no update in place or file locking

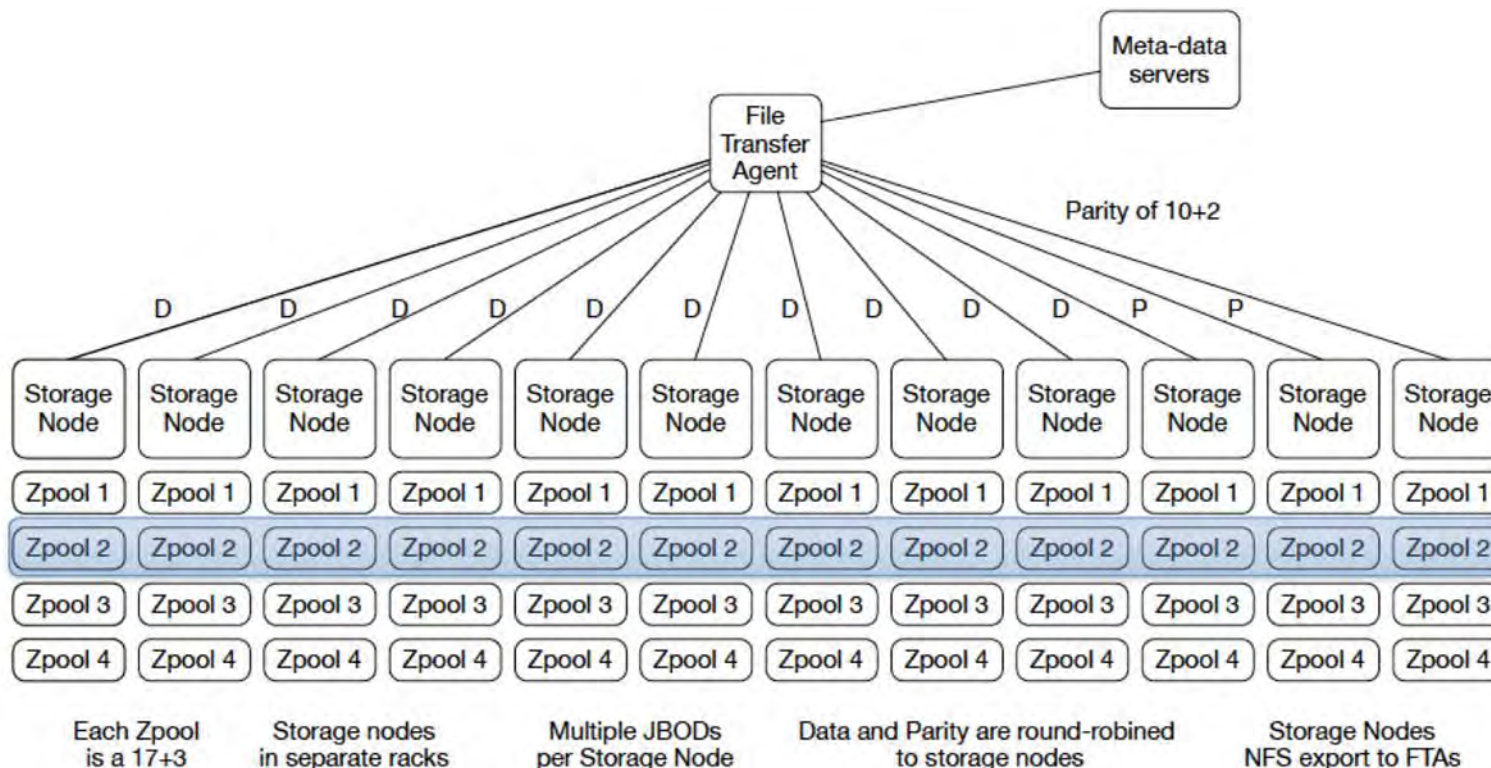  - restricted interactive use

# The Structure of MarFS

- **Pftool**
  - Parallel MPI file transfer utility
- **FUSE mount**
  - Provides user view of metadata
- **MarFS Library**
  - Heart of the software infrastructure
- **DAL/MDAL**
  - Abstraction layers atop data and metadata respectively
  - Allow easy swapping of underlying storage

# Multi-Component

- **Current data storage solution for Campaign**

  - Integrated via the MarFS DAL

  - Stores data as pseudo-objects

- **Cross-server erasure coding atop ZFS pools**

  - Allows failure tolerance at both the disk and server level

  - More reliability allows the use of cheaper disk

  - Erasure coding performed through Intel's Storage Acceleration Library (isa-l)

- **Performance through parallelism**

  - Threaded I/O to multiple servers

# Multi-Component



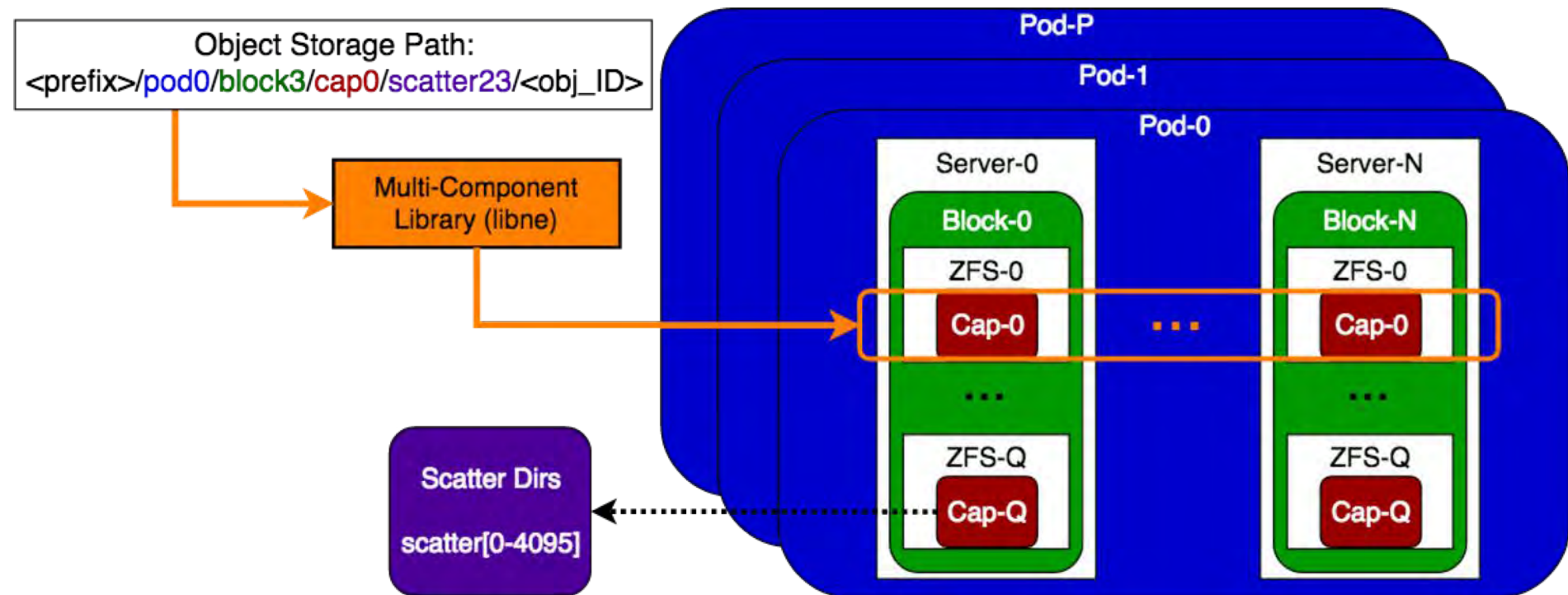Thanks to Kyle Lamb, Dave Bonnie, and Jeff Inman

# Multi-Component

- **Parallelism**
  - Objects hashed across all available servers
- **Transparency**
  - Object stripes simply stored as files in ZFS
- **Resiliency**
  - Data recoverable even after losing entire server racks

# Multi-Component: Parallelism

- **GB Objects are hashed across a multi-level tree**

  - …/pod#/block#/cap#/scatter#/…

    - Pod: Corresponds to a server set (12, for 10+2)

    - Block: Corresponds to an individual starting server (object will still be striped across all servers in a given pod)

    - Capacity-Unit (Cap): Corresponds to a ZFS Pool within each storage server

    - Scatter Directory: Purely logical subdivision to avoid overloaded directories

- **Provides automatic load-balancing**

  - Utilizes all available disk bandwidth, provided enough I/O is issued in parallel
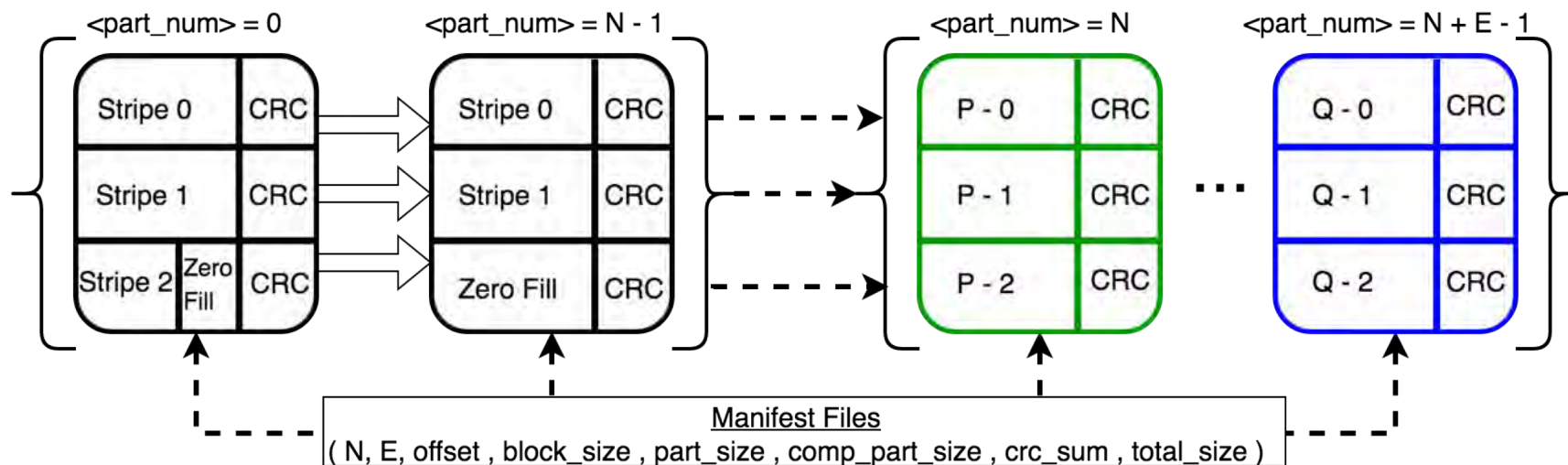
# Multi-Component: Parallelism

# Multi-Component: Transparency

- **Storing to ZFS systems means that objects are plainly visible to administrators**
  - Each object 'part' is paired with a manifest file, providing data and erasure structure info
  - Admins can literally 'ls' object parts and 'cat' manifest info
- **Utilities exist for interacting directly with objects**
  - Read/write data independent of the entire MarFS stack
  - Object integrity checks
  - Erasure rebuild of damaged objects
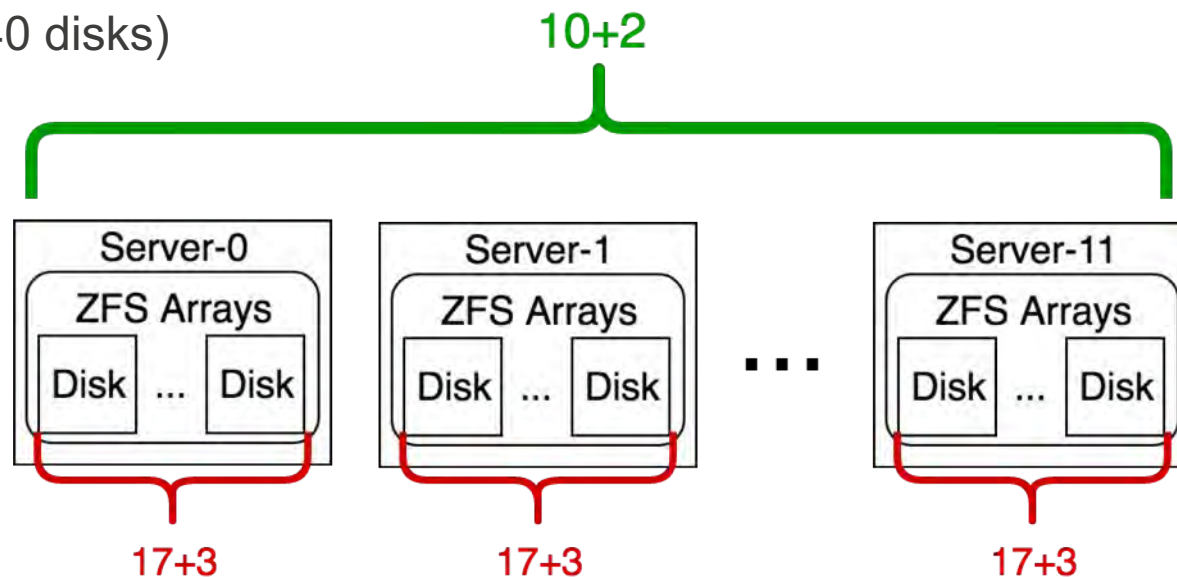
# Multi-Component: Transparency



Storage Path: <prefix>/repo<N>/pod<P>/block<Q>/cap<X>/scatter<Y>/<obj_ID>.<part_num>
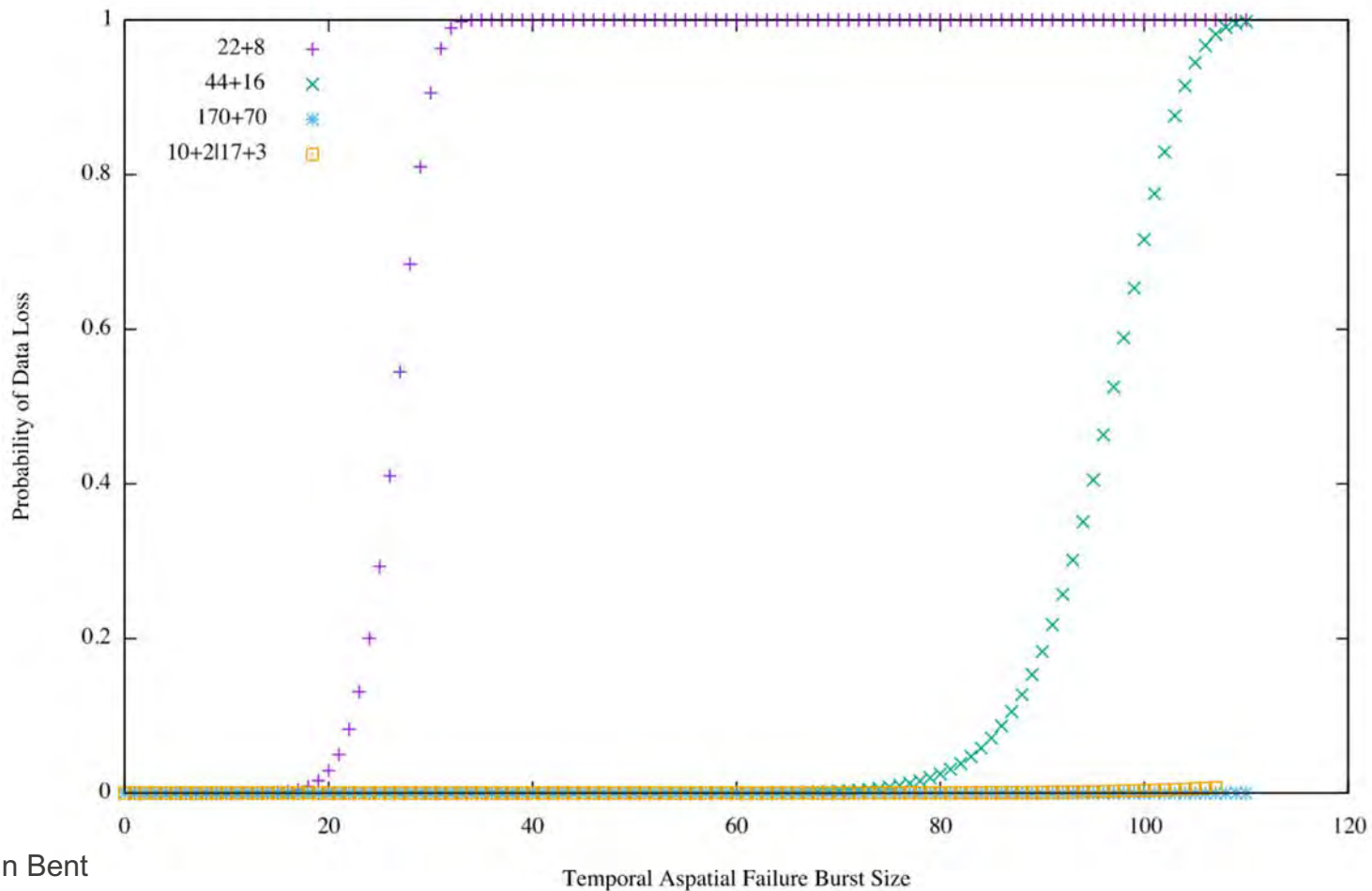
# Multi-Component: Resiliency

- **Multi-tier erasure**
  - Multi-Component: **10+2** across servers
  - ZFS: **17+3** across disks
  - Tolerates min 11 and max 70 disk failures per stripe (set of 240 disks)

# Multi-Component: Resiliency



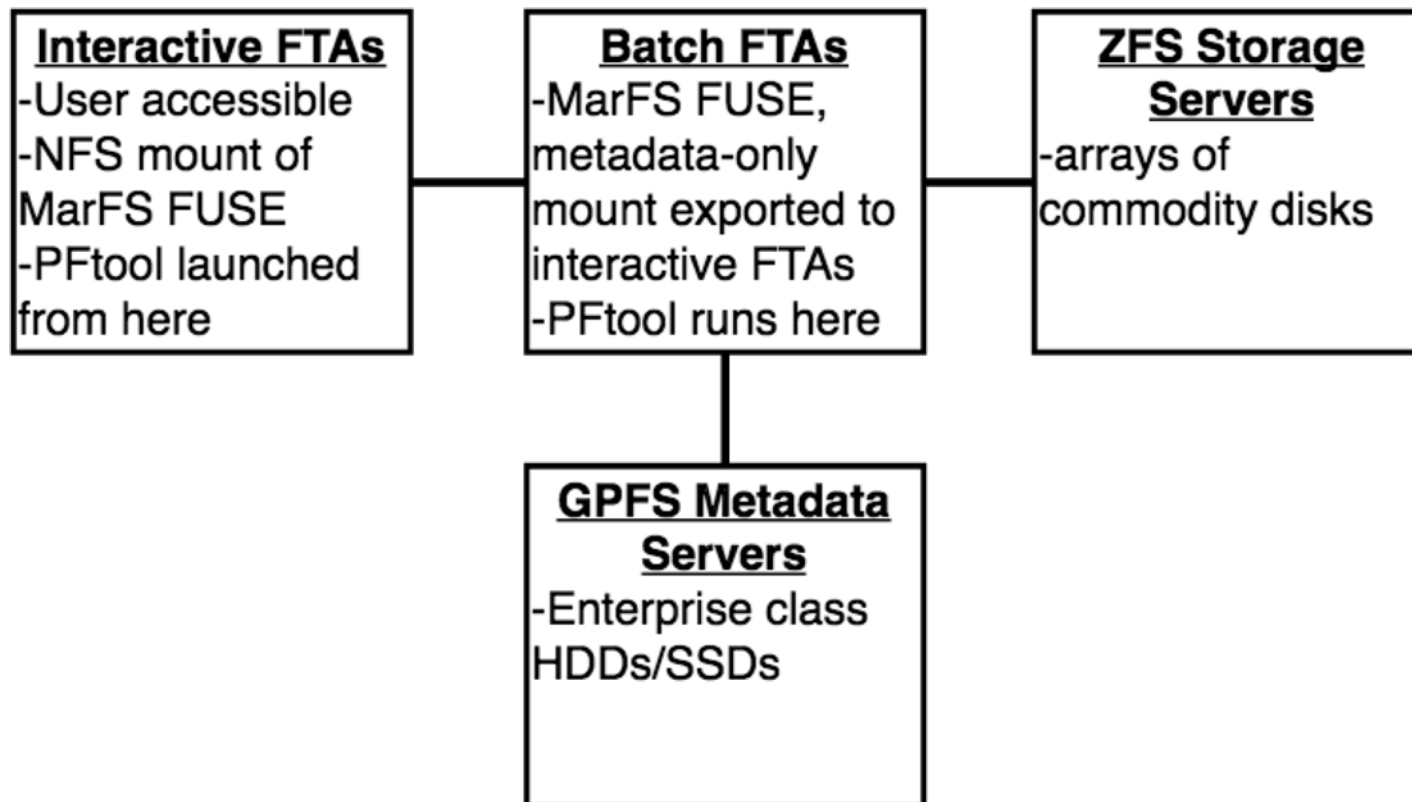Data Loss Probabilities with One Trillion Objects

Thanks to John Bent
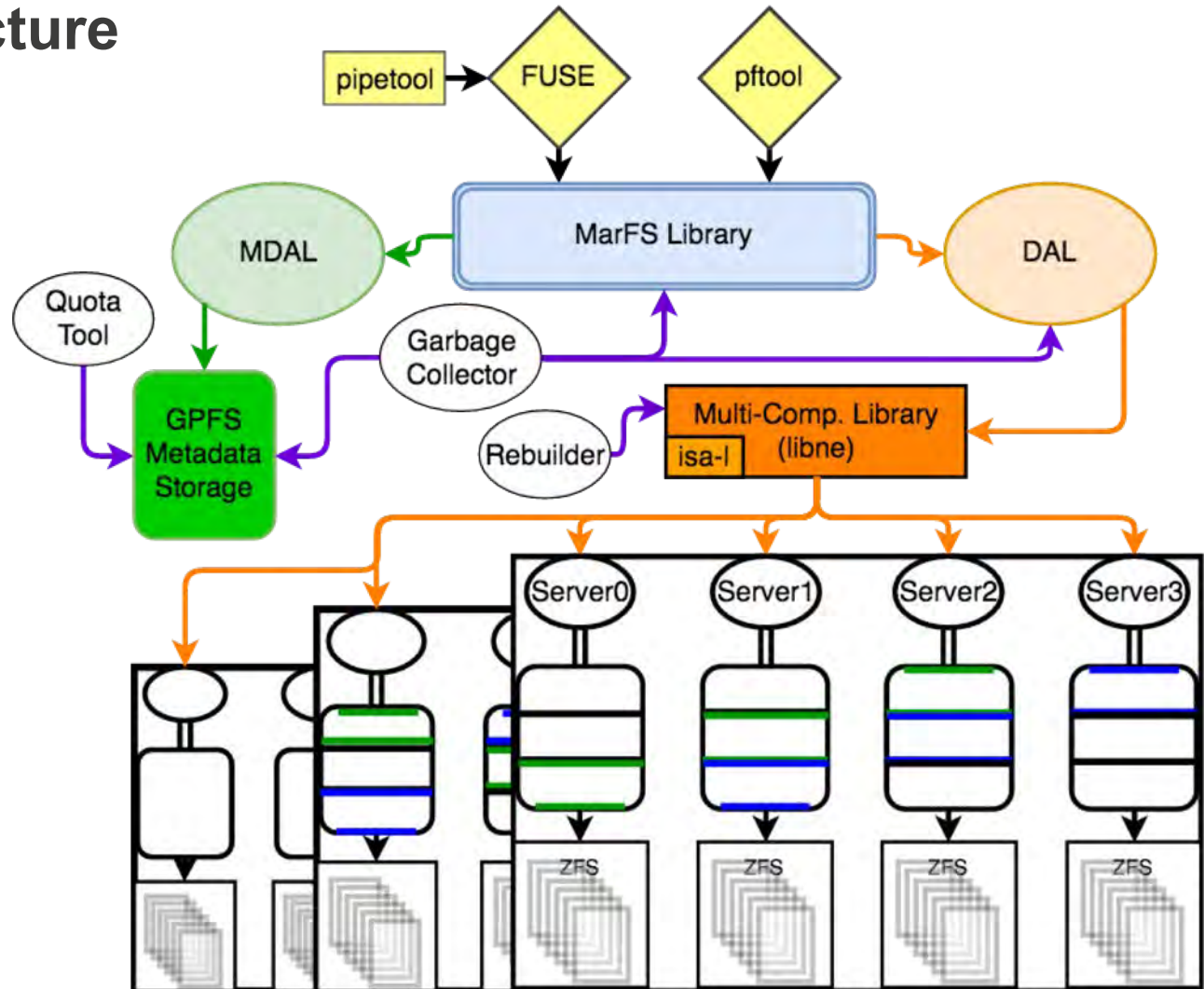
# Multi-Component: Resiliency

- **The Multi-Component Library (libne) automatically responds to read failures and makes use of erasure code if necessary**

- **'Degraded' writes are permitted up to a configurable safety threshold**

  - 10+2 can be written as 10+1 or 9+2 and later rebuilt to full width

- **Any degraded objects detected during read/write are logged**

  - A utility crawls these logs and calls into the Multi-Component Library to rebuild objects

  - The rebuilder utility can also be run against entire capacity-units to assess overall health or recover from catastrophic failures
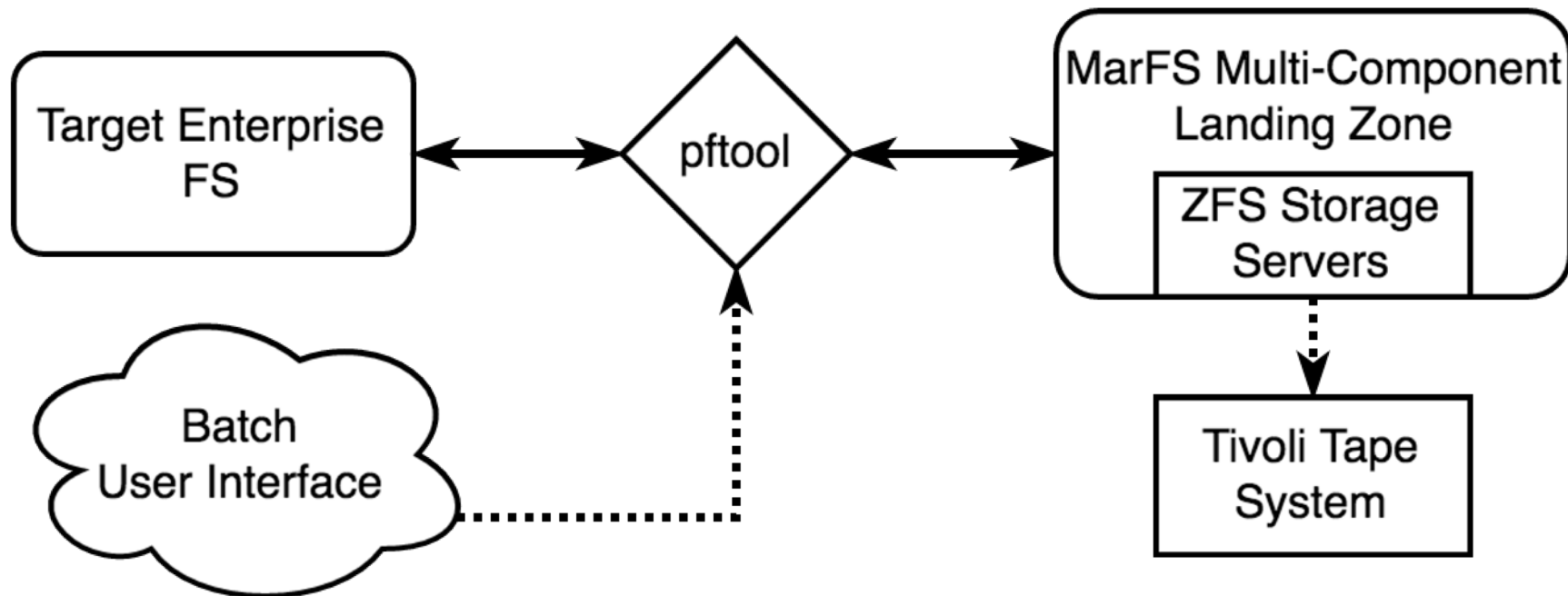
# Deployment

# The Big Picture
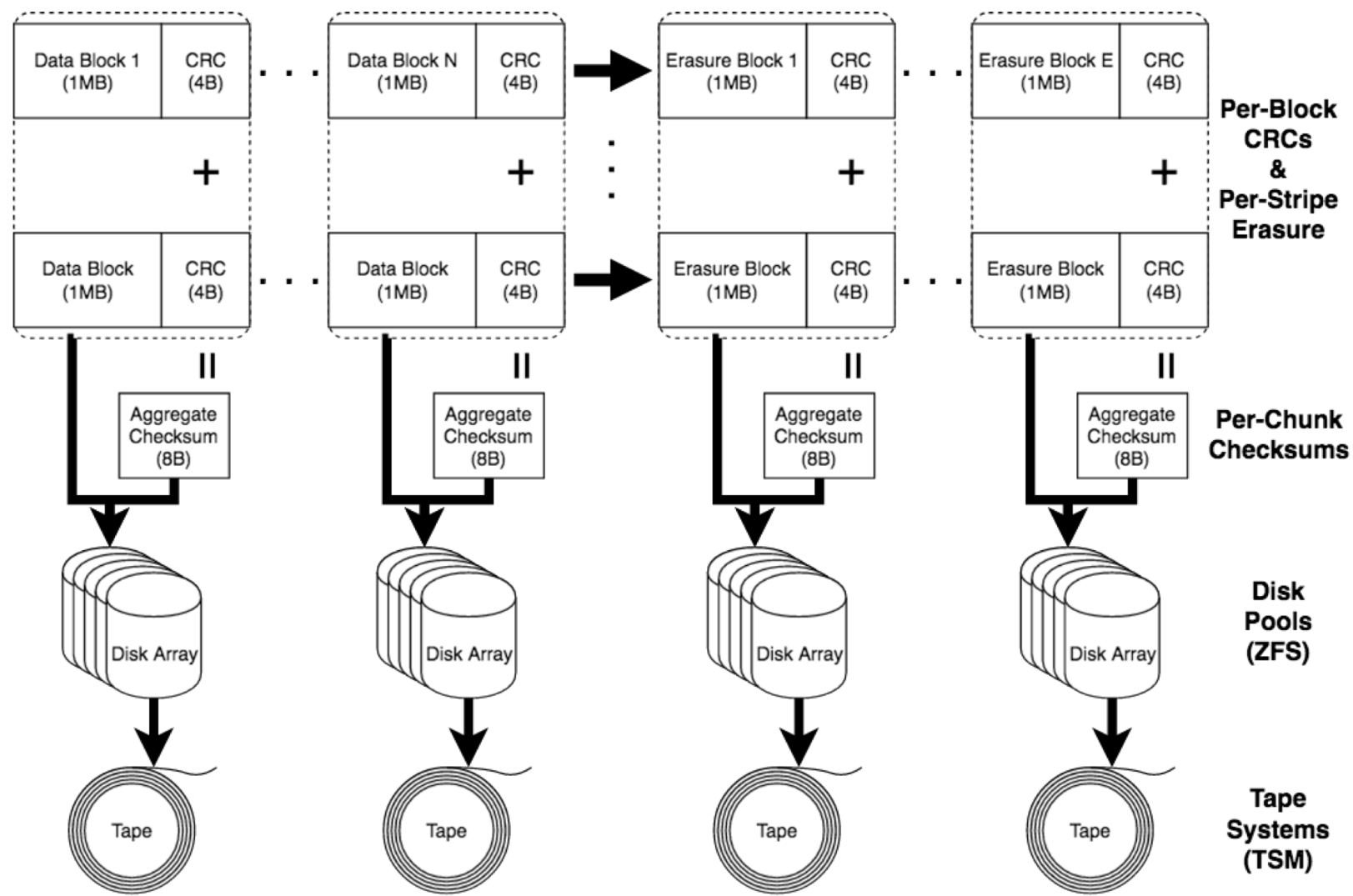
# Current Status

- **MarFS Multi-Component has been in production use for two years now**

  - 5 interactive and 25 batch FTAs

  - 60PB of total storage

  - Roughly 25 GB/sec aggregate bandwidth for both read and write

  - NFS seems to be the bottleneck

- **Additional deployments in progress for this year**

- **Infiniband RDMA based Multi-Component currently in development**

  - Initial testing yielded 35 GB/sec with a fraction of the mpi ranks

# Future Work

- **Better documentation and administrator controls**

- **Capacity-unit migration**

  - Migration of objects to new storage with no downtime

- **Job scheduling for transfers?**

  - Current lack of scheduling means simultaneous transfers always compete for resources

- **Tape-backed MarFS Multi-Component (Marchive)**

# Marchive Conceptual Implementation

**Github Organization – https://github.com/mar-file-system**

- **MarFS – https://github.com/mar-file-system/marfs**

- **LibNE – https://github.com/mar-file-system/erasureUtils**

**Pftool – https://github.com/pftool/pftool**