EMC²

CHARACTERISTICS OF BACKUP WORKLOADS IN PRODUCTION SYSTEMS
Grant Wallace, Fred Douglis, Hangwei Qian*, Philip Shilane, Stephen Smaldone, Mark Chamness, Windsor Hsu
EMC Corporation, *Case Western Reserve University

**B**ackup
**R**ecovery
**S**ystems

Purpose Built Backup Appliances will protect 8 EB of data by 2015

Data Protection (backup) systems must ingest entire primary stores at regular intervals

**Data Domain Backup Appliance**

Deduplication drives higher backup throughput

Weekly full and incremental backups deduplicate against a base copy

Average of 10x deduplication allows backup storage to scale

Data is growing rapidly:
IDC study estimates 1 ZB in 2010 growing to 35 ZB in 2020

Analyze Auto-Support Statistics

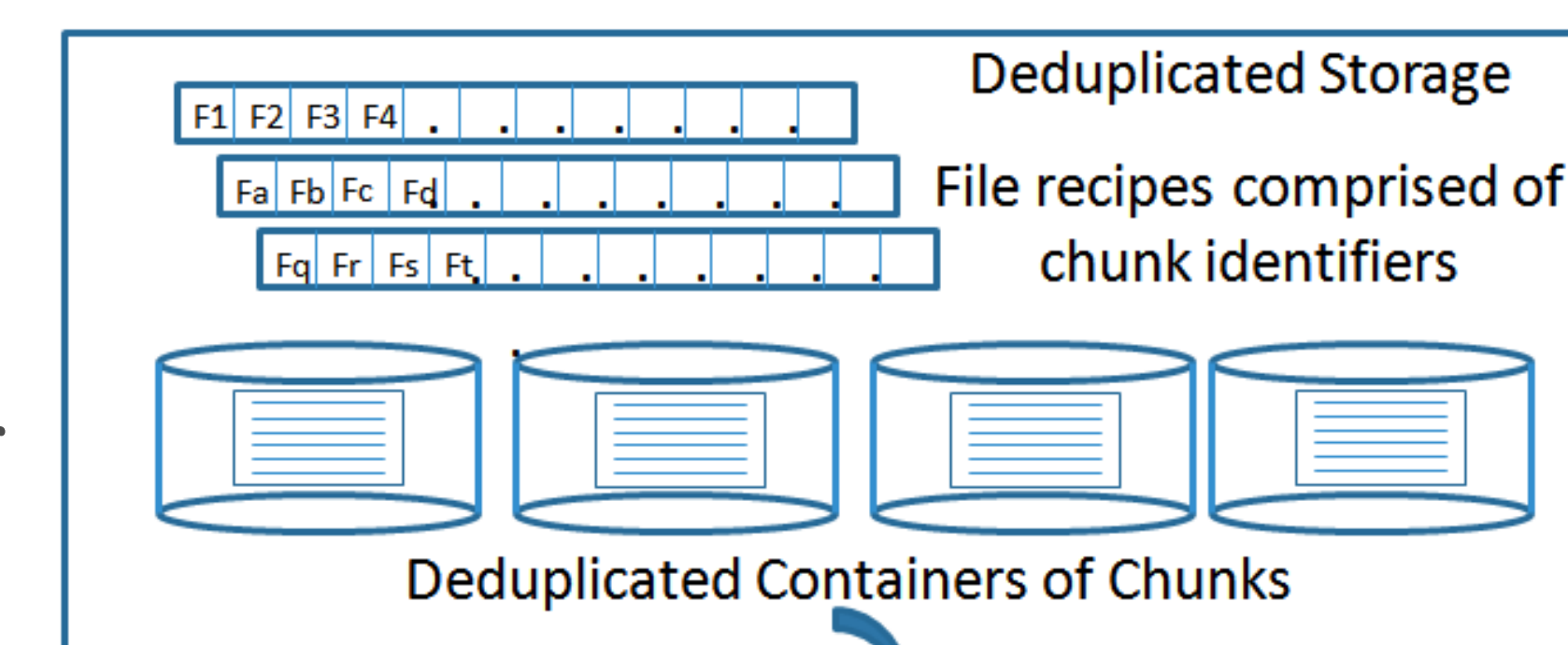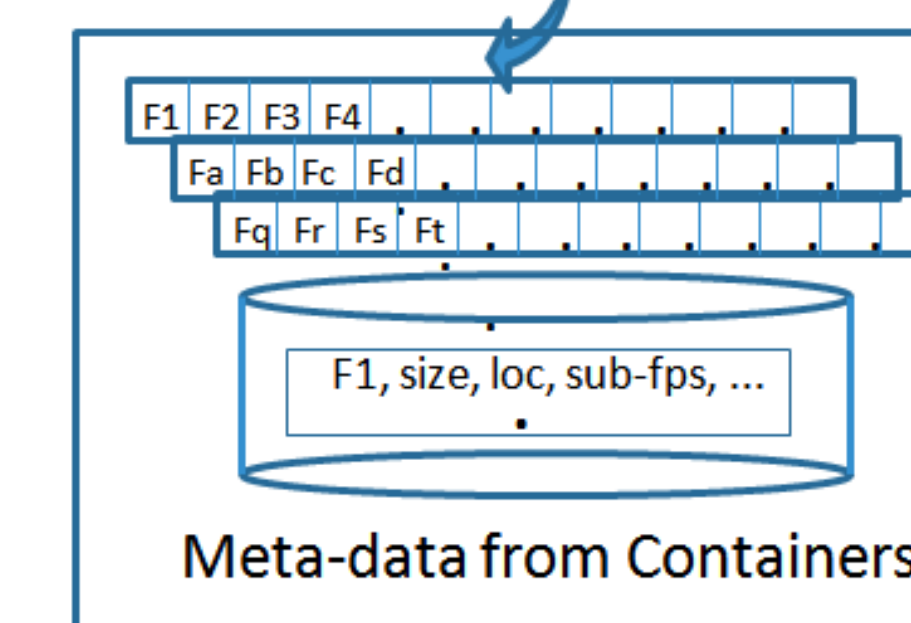1 week, 10,000+ production systems

Collect metadata traces for analysis

Deduplicated Storage
File recipes comprised of chunk identifiers
Deduplicated Containers of Chunks

Perform data collection

Meta-data from Containers
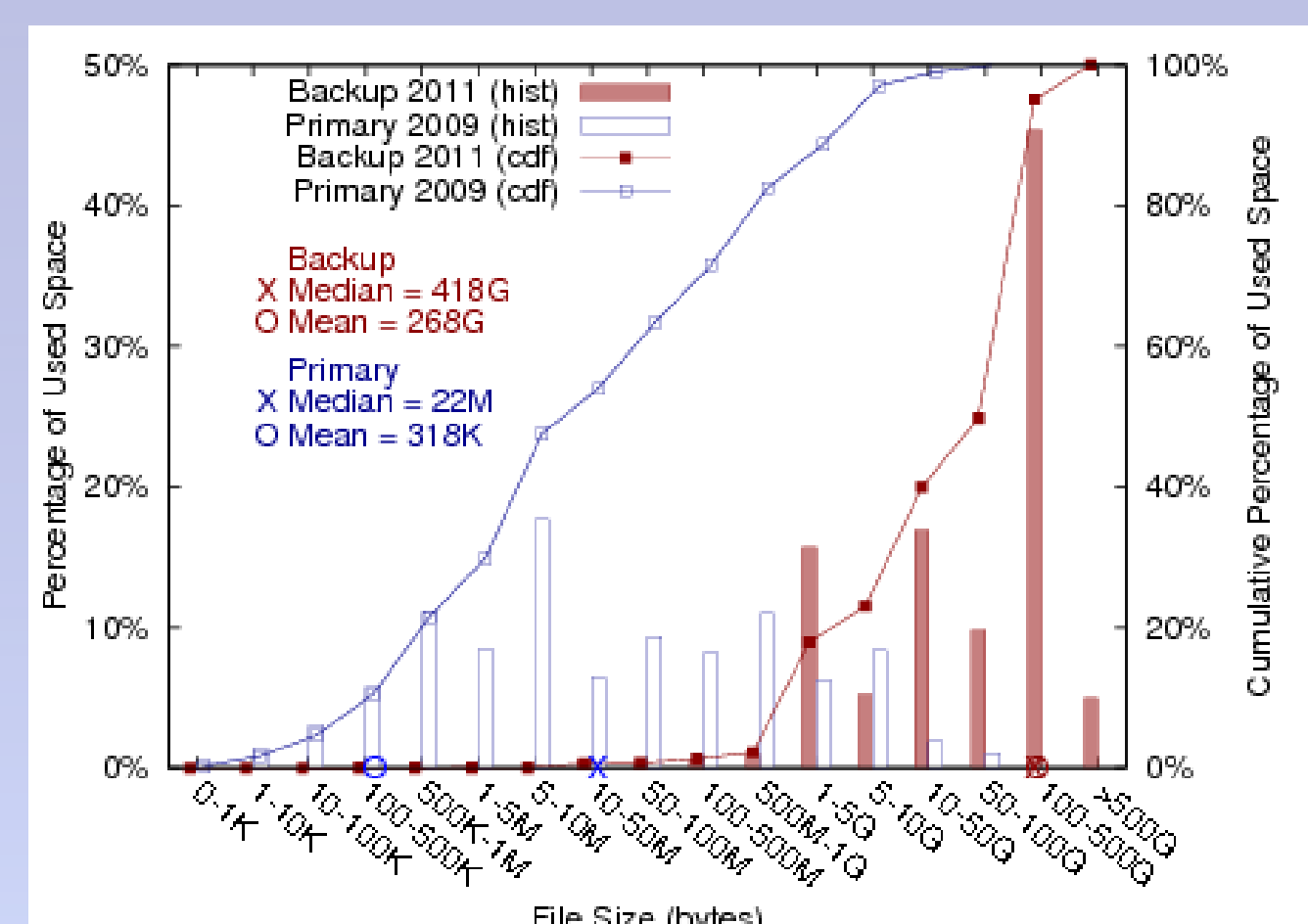
Metadata collected includes anonymized per-chunk fingerprints, sizes, and locations
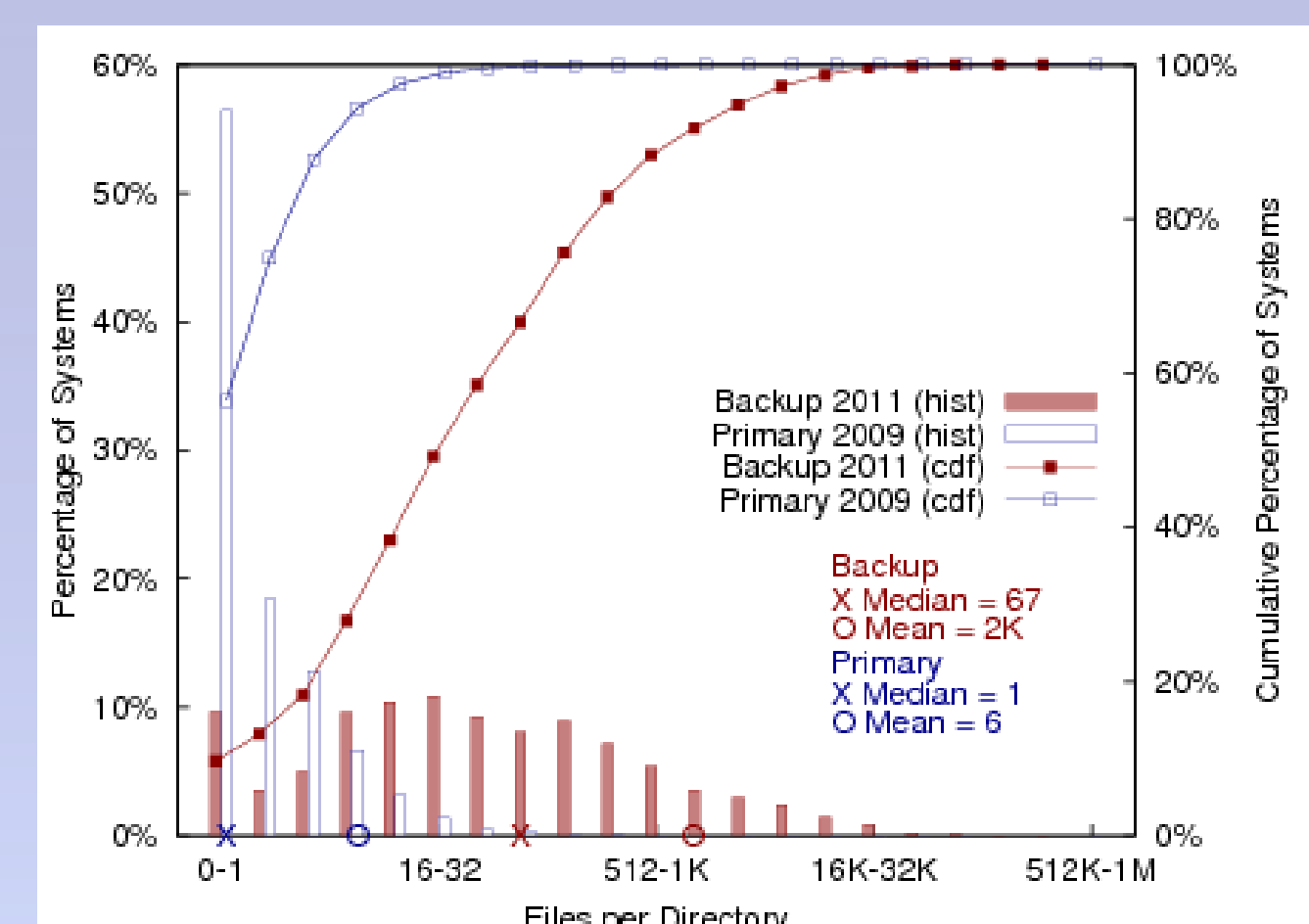
Resulting trace is a time-ordered list of chunk references, with size and location
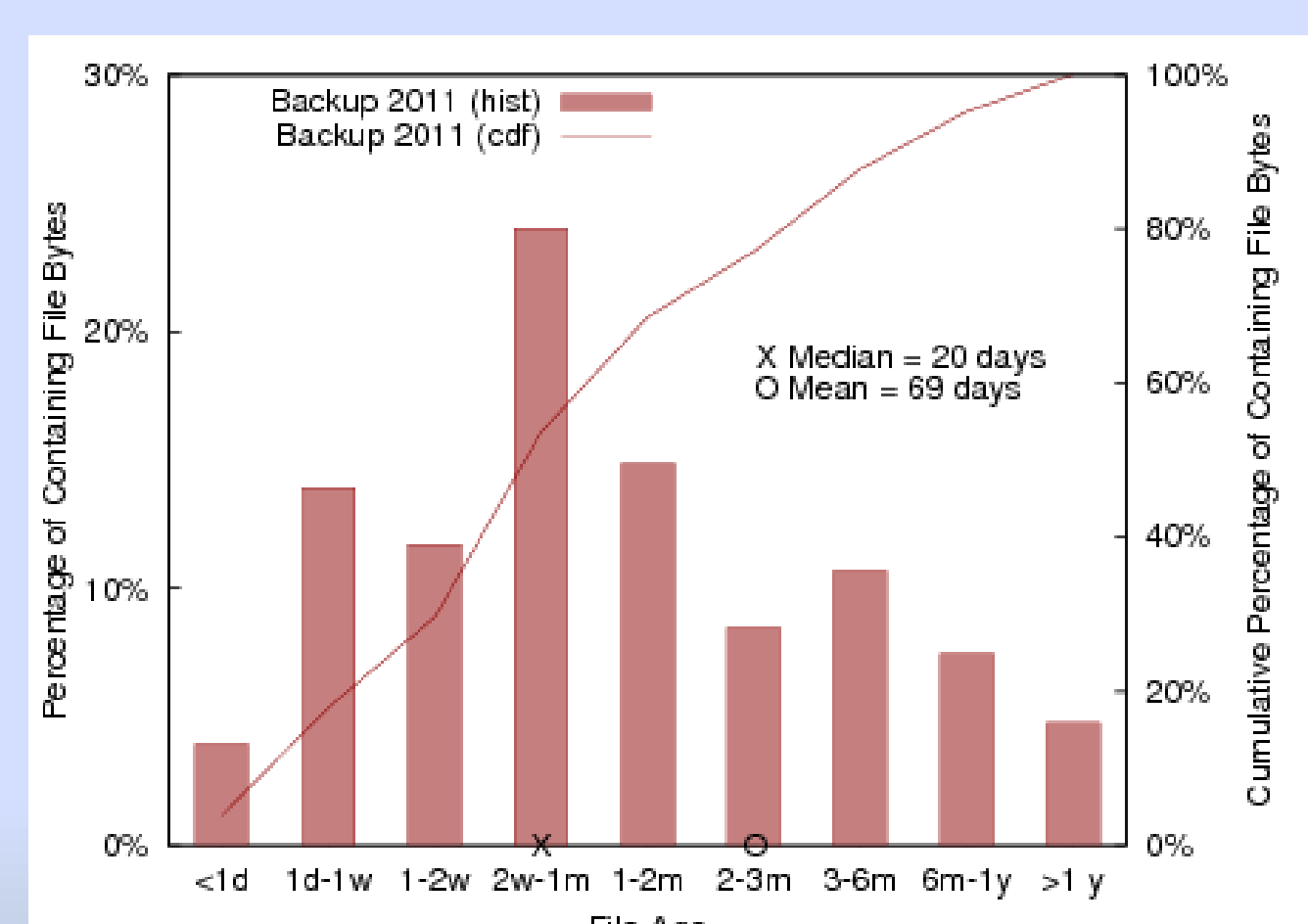
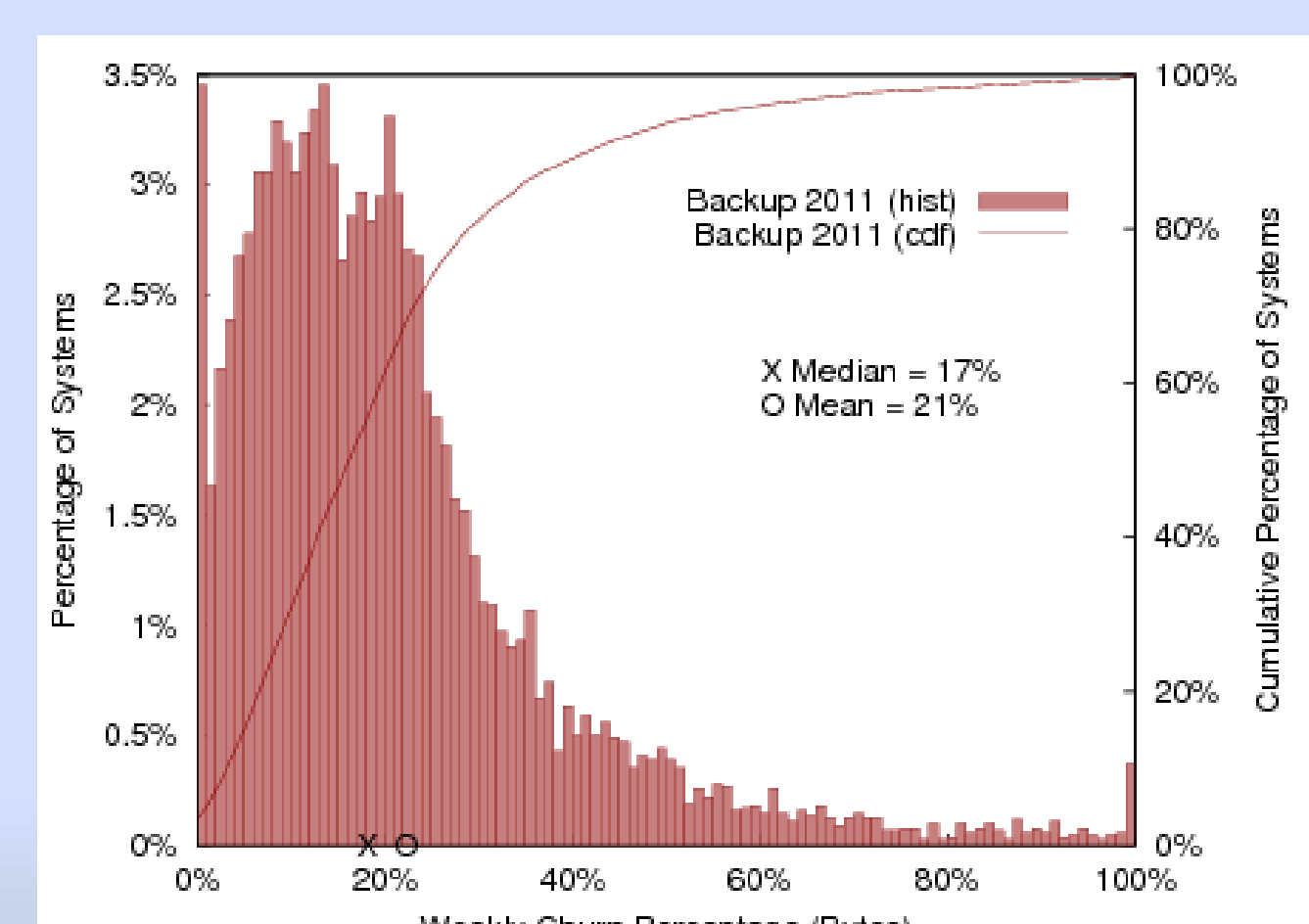**Backup** workloads have **unique properties** compared to **Primary** workloads

Much larger file size

Many more files per directory
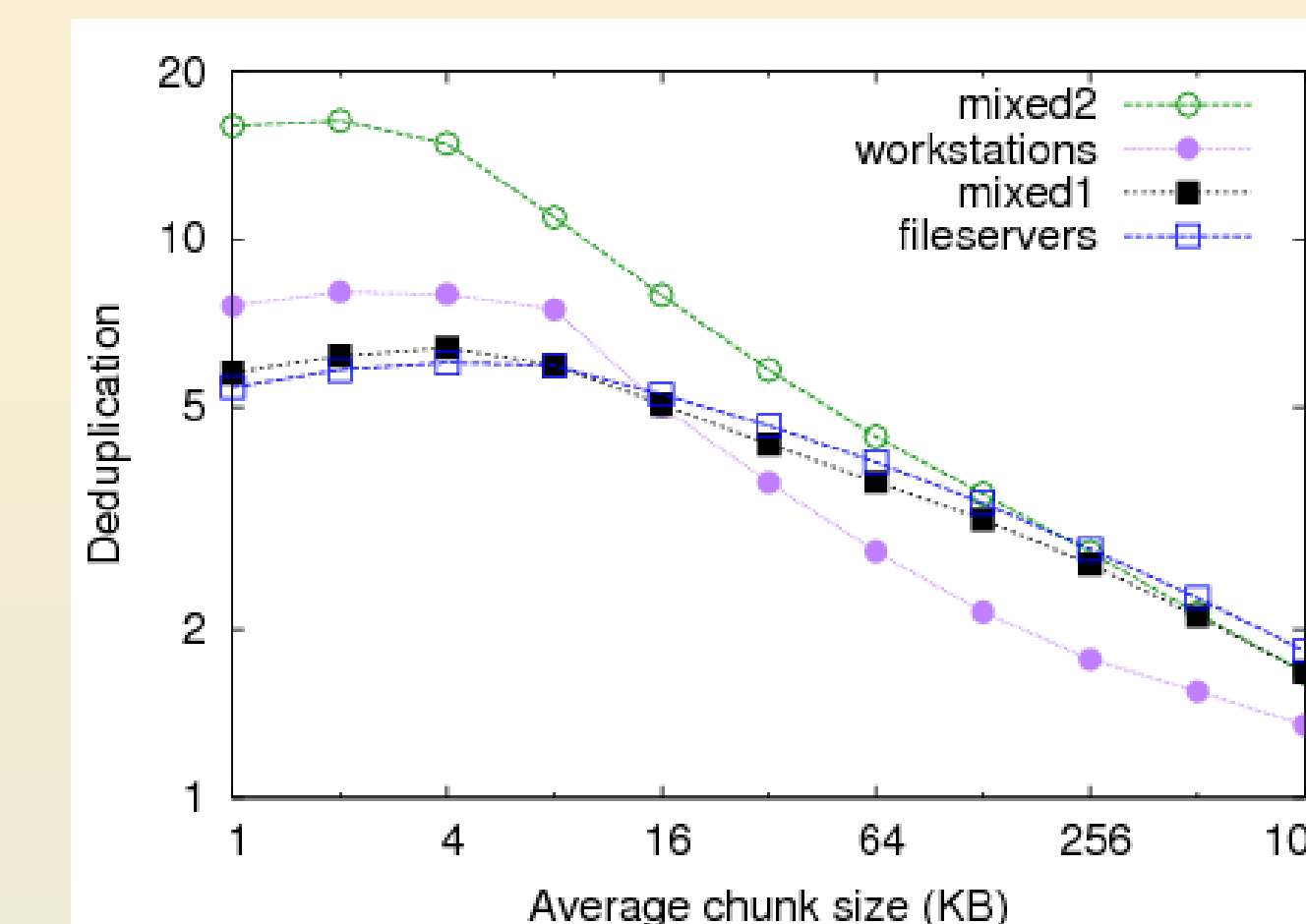
Shorter lived files

Higher churn (20% of data freed and rewritten every week)

**Deduplication** and **effective index caching** are critical to meet large capacity, high throughput demands
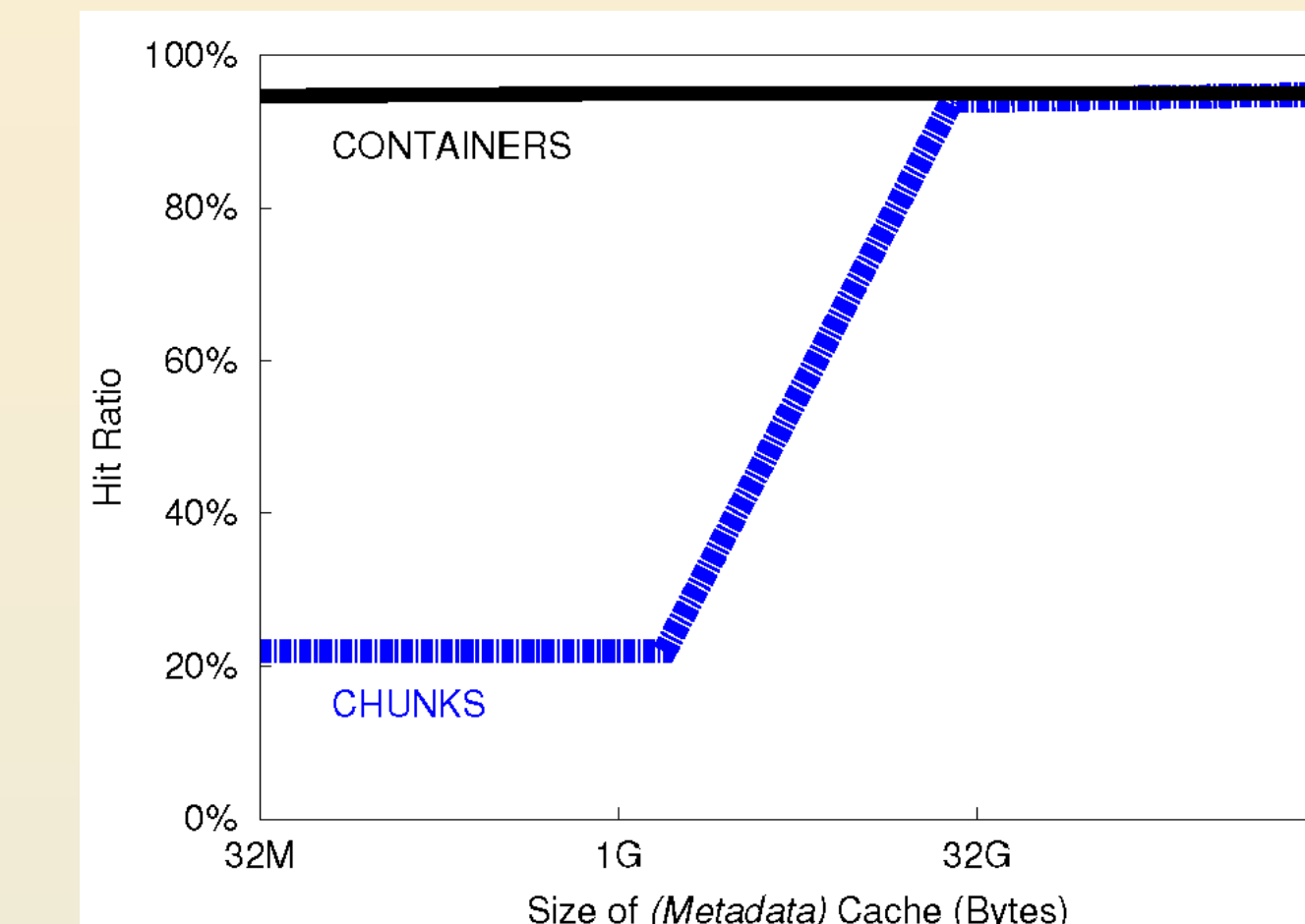
Data Set Characteristics

| Dataset | Size (TB) | Deduplication | Median Age (weeks) |
|---|---|---|---|
| Homedirs | 201 | 14x | 3.5 |
| Database1 | 177 | 5x | 2.2 |
| Email | 146 | 10x | 1.4 |
| Fileservers (Exchange, DB) | 60 | 6x | 5.8 |
| Mixed1 (NAS) | 47 | 6x | 3.2 |
| Mixed2 (Workstations & Servers) | 43 | 11x | 9.4 |
| Workstations | 5 | 8x | 13.6 |
| Database2 | 4 | 2x | 0.2 |

Caching Results (Simplified)

Metadata traces collected from over 700 TB of backup data representing various dataset types and sizes

Best deduplication with 4KB chunks, but 8KB is often sweet spot for data-structure sizes and cleaning

➤ Temporal LRU caching of chunk references needs large cache
➤ Stream-informed LRU caching of regions of chunk references is effective with a small cache

**Backup != Primary**