# Discriminating Hierarchical Storage (DHIS)

**Chaitanya Yalamanchili,**
**Kiron Vijayasankar,**
**Erez Zadok**
*Stony Brook University*

**Gopalan Sivathanu**
*Google Inc.*

http://www.fsl.cs.sunysb.edu/

---

## Outline

- **Introduction**
- Background
- DPROTO Framework
- DHIS Design
- Evaluation
- Related Work
- Conclusions

---

## Large-scale Storage Systems

**EMC® Symmetrix 8830**
80 333Mhz processors, 64GB RAM

**NetApp FAS6080**
1170 Disks, 1170TB Storage

- Storage interfaces unchanged
- Processing power and memory available

---

## Motivation

- Storage System Tradeoffs (RAID)
  - Availability
    - Replication, Multipathing, Failover
  - Cost
    - Power, Backup, Capacity
  - Performance
    - Striping, Load balancing
- Storage Management difficult
- Make better use of processing power at the storage system level to balance tradeoffs

---

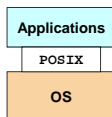## Information Gap

- Layered System Design
  - Layers: fundamental to modern systems
    - Modularity
    - Independent innovation
  - But layers hide information
    - Information gap
    - Age-old problem in computer systems
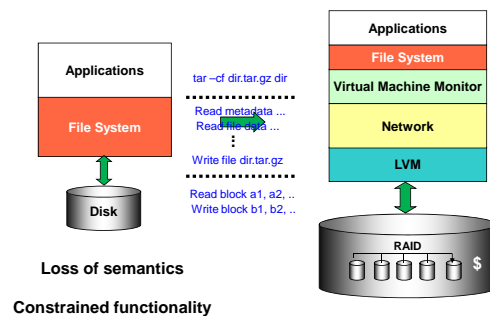    - Constrained functionality within layers

**Applications**

`POSIX`

**OS**

---

## Information Gap (cont.)

**Applications**

**File System**

**Disk**

tar –cf dir.tar.gz dir

Read metadata ...
Read file data ...

Write file dir.tar.gz

Read block a1, a2, ..
Write block b1, b2, ..

**Loss of semantics**

**Constrained functionality**

**Applications**

**File System**

**Virtual Machine Monitor**

**Network**

**LVM**

**RAID**    $

## DHIS Overview

- Bridges the information gap with an extended storage system interface
- Uses hints from higher-level software (applications, file systems)
- Storage management done inside the storage system firmware
- RAID levels and NVRAM comprise the hierarchy
- Ordering in the hierarchy depends on the hints (configurable)

STONY BROOK

## Outline

- Introduction
- **Background**
- DPROTO Framework
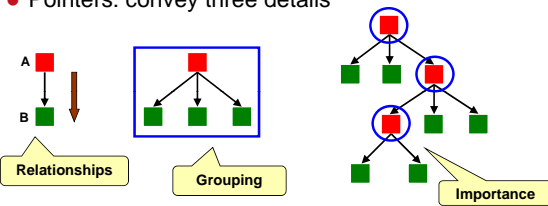- DHIS Design
- Evaluation
- Related Work
- Conclusions

STONY BROOK

## Block-Based Storage

- Two basic entities: data and pointers
- Pointers: convey three details



A

B

Relationships

Grouping

Importance

- Today's disks are unaware of pointers

STONY BROOK

## Type-Aware Storage

- Bridge information gap through pointers
  - Disks aware of pointers
- Higher level software communicate pointers
  - File systems or user applications
  - Explicit disk interface extension
- Type-Safe Disks (TSDs)
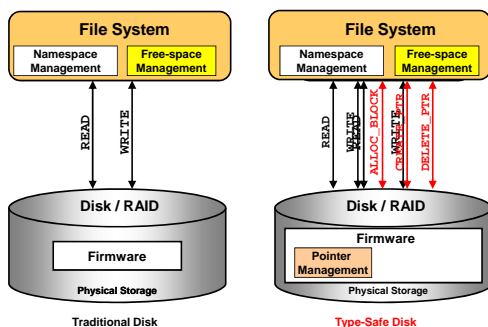  - Track pointers and enforce constraints
  - Manage free-space

STONY BROOK

## TSD Infrastructure



Traditional Disk    Type-Safe Disk

STONY BROOK

## TSD Interface

- **READ(Block)**
- **WRITE(Block)**
- **ALLOC(Ref, Count, Hint)**
- **CREATE_PTR(Src, Dest)**
- **DELETE_PTR(Src, Dest)**

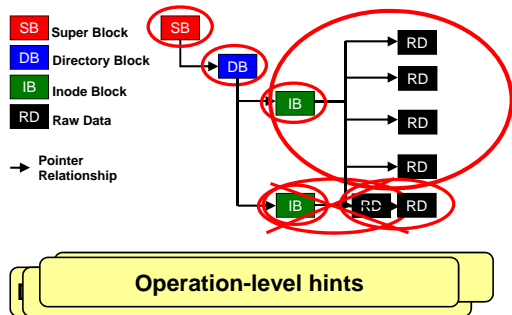STONY BROOK

2

## TSDs and File Systems

- **SB** Super Block
- **DB** Directory Block
- **IB** Inode Block
- **RD** Raw Data
- → Pointer Relationship

**Operation-level hints**

---

## TSD Project

- Type-Safe Disks [OSDI 2006]
- Exploiting Type-Awareness in a Self-Recovering Disk [ACM StorageSS 2007]
- Selective Versioning in a Secure Disk System [Usenix Security 2008]

---

## Outline

- Introduction
- Background
- **DPROTO Framework**
- DHIS Design
- Evaluation
- Related Work
- Conclusions

---

## DPROTO Framework
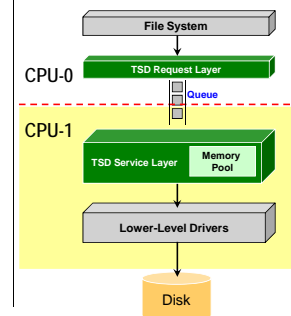
- **CPU**
  - Dual CPU machine
  - Isolated CPU for disk
  - Linux "`cpusets`" interface
- **Memory**
  - Isolated memory pool
  - Pre-allocated
  - Leveraged "`mempool`"
- **Broader uses**
  - Prototype other disk-level features

File System

CPU-0 　 TSD Request Layer

Queue

CPU-1

TSD Service Layer 　 Memory Pool

Lower-Level Drivers

Disk

---

## Outline

- Introduction
- Background
- DPROTO Framework
- **DHIS Design**
- Evaluation
- Related Work
- Conclusions

---

## DHIS
### A **D**iscriminating **HI**erarchical **S**torage system

- Communicate data properties to disk
  - Granularity: groups of data conveyed through pointers
  - Well-defined *attributes*
- Placement decisions based on attributes
  - Reliability
  - Performance
  - Cost
- Series of configurable hierarchies to place data in

3

## Attributes

- Importance: High / Low
- Access Pattern: Random / Sequential
- Read-most / Write-most
- Hot / Cold
- Life-time: Temporary / Long-lived

Can add new attributes

STONY BROOK

---

## DHIS Interface

- Adds **SETATTR(Block, attr)**
  - ◆ Used by higher-level software to pass hints
- Supports the TSD interface
  - ◆ Relates blocks using logical pointers
  - ◆ Attributes inherited from parents

STONY BROOK

---

## Ext2DHIS File System

- Unlike Ext2, Ext2DHIS does not perform free-space management
- Block allocation via the **ALLOC** primitive
- Issues **CREATE_PTR** or **DELETE_PTR** whenever a new pointer is added or removed for a meta-data block
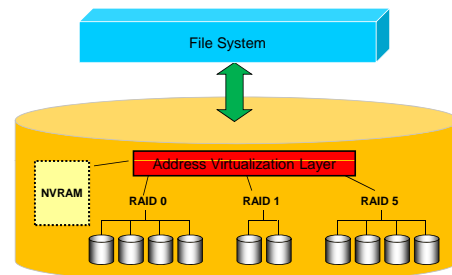- Uses **SETATTR** to set attributes

STONY BROOK

---

## DHIS Architecture

STONY BROOK

---

## DHIS Detailed Design

STONY BROOK

---

## RAID Placement Policy

| IMPORTANT | ACCESS PATTERN | READ/WRITE MOST | HOT/COLD | Preferred RAID Levels |
|-----------|----------------|-----------------|----------|----------------------|
| Low | Any | Any | Any | 0, 5, 1 |
| High | Any | Any | Cold | 5, 1, 0 |
| High | Not-Set | Not-Set | Not-Set / Hot | 5, 1, 0 |
| High | Random | Not-Set / Write | Not-Set / Hot | 1, 5, 0 |
| High | Random | Read | Not-Set / Hot | 5, 1, 0 |
| High | Sequential | Any | Not-Set / Hot | 5, 1, 0 |

Configurable, can add new RAID levels

STONY BROOK

4

## Policies

- Temporary files
  - ◆ RAID-0
  - ◆ Placed in an isolated portion of disk
    - ▪ Reduce disk fragmentation
- Meta-data blocks
  - ◆ Identified as those having outgoing pointers
  - ◆ Placed in the RAID level of highest reliability and best random access performance (RAID 1 in our setup)

STONY BROOK

## Policies (cont.)

- NVRAM caching
  - ◆ Caching candidates
    - ▪ All meta-data
    - ▪ "Hot" and "write-most" data
  - ◆ Absorbs write latency
  - ◆ Sequential workloads are not chosen
  - ◆ Use block liveness info to remove freed blocks

STONY BROOK

## Outline

- Introduction
- Background
- DPROTO Framework
- DHIS Design
- **Evaluation**
- Related Work
- Conclusions

STONY BROOK

## Postmark

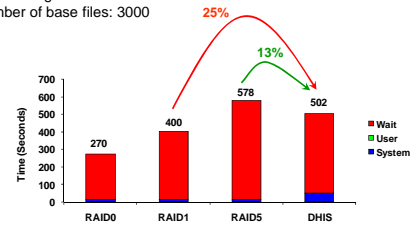File size range: 400 – 600 KB
Number of base files: 3000



25%
13%

| | RAID0 | RAID1 | RAID5 | DHIS |
|---|---|---|---|---|
| Time (Seconds) | 270 | 400 | 578 | 502 |

Legend: Wait, User, System

Attributes: IMPORTANT | RANDOM ➔ RAID1 policy

STONY BROOK

## Micro-Benchmarks

**Sequential Read 7.5 GB**

18%
10%

| RAID1 | RAID5 | DHIS |
|---|---|---|
| 153 | 114 | 126 |

**Sequential Write 7.5 GB**

33%
16%

| RAID1 | RAID5 | DHIS |
|---|---|---|
| 198 | 158 | 132 |

Legend: Wait, User, System

Attributes: IMPORTANT | SEQUENTIAL | READ-MOST

Attributes: IMPORTANT | SEQUENTIAL | WRITE-MOST

STONY BROOK

## Micro-Benchmarks (cont.)

**Random Read**
20,000 4K reads

4%
18%

| RAID1 | RAID5 | DHIS |
|---|---|---|
| 254 | 207 | 245 |

**Random Write**
150,000 4K writes

2%
53%

| RAID1 | RAID5 | DHIS |
|---|---|---|
| 116 | 252 | 118 |

Legend: Wait, User, System

Attributes: IMPORTANT | RANDOM | READ-MOST

Attributes: IMPORTANT | RANDOM | WRITE-MOST

STONY BROOK

## OLTP Workload



Latency and Throughput vs. no. of reader
threads (FileBench OLTP)

Attributes: IMPORTANT | RANDOM | READ-MOST ➔ RAID5 policy

STONY BROOK

## Kernel Compile Workload



Compiled the linux
kernel 2.6.28
sources

Attribute:
TEMPORARY ➔
RAID0 policy

Similar (+7%)

STONY BROOK

## Outline

STONY BROOK

## Related Work

- **Object-Based Storage Devices**
  - Objects versus Blocks
  - Objects support attributes
  - **Problem:** Require fundamental changes to higher-level software
- **Self-\* Storage**
  - Automated administration
  - Notion of supervisors, workers, and routers
  - Works better when workers are intelligent (like DHIS)
- **HP AutoRAID**
  - Newly written data placed in RAID 1
  - Data migrated to RAID 5 as it gets cold
  - **Problem:** Limited to cold/hot attribute; data migration can be costly

STONY BROOK

## Related Work (cont.)

- **ExRAID**
  - Expose fault boundaries and redundancy information to the file system
  - **Problem:** Managing redundancy within the file system can be difficult, requiring the careful placement of inodes and data blocks to ensure efficient operation under failure.
- **RAIF**
  - Stackable fan-out file system
  - **Problem:** Rule management done at the file system level; extra layer adds overhead.
- **Semantically smart disks**
  - Automatically infer higher-level operations and data structures
  - **Problem:** Inference is not always accurate

STONY BROOK

## Conclusions

- Enables easy storage management
  - Fine-grained policies
- Attribute association can be automated
  - Ext2DHIS file system
  - Attributes based on file extension
- Online attributes
  - Obviates need for data migration

- Future: _someone_ should offer a more intelligent storage system...

STONY BROOK