# GCTrees: Garbage Collecting Snapshots

Chris Dragga

Doug Santry

# GCTrees

Support snapshots for enterprise workloads

Minimize write overhead

Prototype implemented in ext4:  gcext4

Reduce overhead up to 68x from state of the art

NetApp®

# Outline

1) Overview

2) Theory

3) Implementation

4) Evaluation

**NetApp**

# Snapshot Overview

# Copy-on-write

# Copy-on-Write:  Appending

# Hierarchical Refcounts

# Hierarchical Refcounts: Deletion

# Problem: Fan-out



Even worse for writes

# Potential Solutions

**NetApp®**

# Delay

Write refcount changes to log, then checkpoint

Hope that they'll cancel

No guarantee this will actually happen

Update storm

**NetApp**®

# GCTrees

Don't count references at all

Track lineage

Use garbage collection (GC) to determine what can be freed

Keep deletion, write overhead minimal

# Basic Metadata Structure

# Basic Metadata Structure

# Basic Metadata Structure

# Deletion: Child

# Deletion:  Child



Metadata

Data 1

Data 2

1 extra write

# Deletion: Parent

# Deletion:  Parent

1 extra write



Metadata-COW

Data 1

Data 2

Data 3

# Deletion in a Snapshot Chain

# Deletion in a Snapshot Chain



Metadata → Head → Metadata-COW2

Metadata ← Source ← Metadata-COW2

**2 extra writes**

# What about B+ trees?

Tree ops can move pointers

Need additional pointers:  next and previous

Can increase deletion overhead

**NetApp**

# GCTree summary

One extra write per metadata block when writing

One to two extra writes for most deletions

Avoid deletion overhead with background scans

**NetApp**

# Implementing gcext4

Ext4:  Extent-based file system

Two metadata types:  inodes, extent blocks

# GCTree Metadata Layout

| Source 48 b | Head 48 b | Next 48 b | Prev 48 b | Borrowed bitmap 1 B / 42 B |
|---|---|---|---|---|

# Adding GCTree metadata

ext4 inode (256 B):

| Inode Header | Inode Data | Inode Tail | X-attrs |
|---|---|---|---|

gcext4 inode (256 B):

| Inode Header | Inode Data | Inode Tail | GCTree | X-attrs |
|---|---|---|---|---|

NetApp

# Adding GCTree metadata

ext4 extent (4096 KB):

| Extent Header | Extent Pointers (340) | Extent Tail |
|---|---|---|

gcext4 extent (4096 KB):

| Extent Header | Extent Pointers (334) | GCTree | Extent Tail |
|---|---|---|---|

**NetApp**

# Adapting ext4 to COW:

Straightforward, but fiddly

COW once per snapshot

Inodes proved problematic

**NetApp**

# Fixed Location Inodes



Directory Tree

# Solution: ifile



Directory Tree

# Implementing Deletion

Separate kernel threads act as scanners

Deletion enqueues a message

Scanners process message, do actual deletion

Removing a snapshot deletes an ifile

NetApp®

# Evaluation

Do they work?

How do they compare?

NetApp®

# Experimental Set-up

3GHz Core 2 Duo, 6GB RAM

7200 RPM, 160GB hard drive

NetApp

# Basic Benchmarks

Fileserver

OLTP

VM:  fileserver in VirtualBox

NetApp®

# Benchmark Configuration

Storage footprint: 2x memory

5 repetitions, 3 hours each

Snapshot per hour for Fileserver, VM

Snapshot per five minutes for OLTP

NetApp®

# Comparison to ext4

# Hierarchical Refcount Comparison

Direct performance comparison infeasible

Look at block-write overhead

**NetApp**

# File Systems for Comparison

Btrfs:  refcounting file system

Simulation in gcext4:  accounts for differences in btrfs

# Simulation Methodology

Assumes refcounts stored in a contiguous map

Use a 15MB durable log

**NetApp**

# Results:  Traditional Workloads

# Results:  Enterprise Workloads



OLTP

VM

# Conclusion

GCTrees:  snapshots for enterprise workloads

Substantial gulf between refcounts and GCTrees

Optimal choice depends on workload

**NetApp®**

# Thank you!