

FUSE-Based File System Monitoring and I/O Classification with IOFS

Marcus Boden, Lars Quentin, Julian Kunkel



2022-06-02

Classification

Outline



2 Monitoring



4 Outlook

Motivation

- Many (new) users have limited knowledge about their I/O patterns
- I/O hard to analyze for users
 - no simple tools for analysis
 - analysis tools may need advanced skills or elevated access
- often treated as a black-box, if at all
- Different file systems have different characteristics
 - bad I/O patterns can slow down global file systems

Goals



easy to use

- no changes in code necessary
- no extra libraries/linking
- easy to test and set up
 - dockerfiles for database and Grafana
 - analysis dashboards available as json

FUSE

FUSE (Filesystem in USErspace)

- Userspace interface to create file system
- kernel module and libfuse library
- rewrite file system calls
- examples: sshfs, ntfs-3g

Security considerations:

- setuid binary
- only the user can access the mount
 - can be relaxed with allow_root or allow_other

Introduction	Monitoring	Classification	Outlook
00	0000000		00

Architecture

- FUSE-mount a directory to monitor it
- Intercept all operations on mount
 - start timer
 - change path
 - do operation
 - end timer
 - write timing and size to global struct
- global struct with counter for all operations
- second thread for reporting
 - read counters
 - send to database
 - clear the counters











Introduction	Monitoring	Classification	Outlook
00	○○○○○○●		00
Performance impact			

IO500 tests with SCC rules

10 Nodes, 80 tasks

Result: 50%-60% of native performance

I/O Classification

- Idea: Create models for different I/O Operations
- Label accesses by mapping access sizes to characteristics
- Workflow:
 - 1 Create Benchmarks isolating different characteristics
 - 2 Analyze each Benchmark; find possible clusters
 - 3 Create regression models from analyzed benchmarks
 - 4 Classify new I/O operations via regression models

Creating Benchmarks

- Goal: Find expected durations for different types of I/O requests
 - Different benchmark variations to isolate access features
 - File access patterns
 - Access size of I/O operation
 - bypassing cache via 0_DIRECT
 - Prereading every block
 - Dropping cache after every benchmark
 - Deleting the benchmarked file afterwards

Introduction	Monitoring	Classification	Outlook
00	0000000	00000	00
			6

Analyzing a Single Benchmark





Misleading histogram

Introduction	Monitoring	Classification	Outlook
00	0000000	00000	00

Kernel Density Estimation

We approximate the PDF f with the KDE

$$\hat{f}(x) = \frac{1}{N} \sum_{i=1}^{N} \mathcal{K}(x - x_i)$$

with *K* being the gaussian normal distribution.

- **a** $x x_i$ to map x_i onto the kernel K
- **1**/*N* normalizes to $\int \hat{f}(x) dx = 1$

Visually:

 We put a normal distribution at each observation



Analyzed Benchmark, Clusters Coloured

Classification

Creating Models from Benchmarks

- Durations increase linearly in access size ⇒ linear interpolation!
- Interpolate over all analyzed benchmarks
 - Points defined as (access_size, access_time)
- Compare each new observation with computed models
 - Classify new operations into our benchmark categories

ToDo

- per file accounting
- Slurm integration
- DB-agnostic reporting
- take file system load into account
- integrate modeling and classification to IOFS
- make recommendations based on classifications

Questions

Thanks for your attention

Questions? Feedback? Ideas? The source code can be obtained at: https://github.com/gwdg/iofs/ https://github.com/lquenti/blackheap/