CS 744: SPLIT-FS

Shivaram Venkataraman Fall 2021

ADMINISTRIVIA

- Course Project: Check in: Today! --- Instructions on
- Midterm 2 next week!

-> From SQL to TPU paper

1 page - Canvas

Piazza

Serverless Computing



Compute Accelerators



Infiniband Networks



Non-Volatile Memory

PERSISTENT MEMORY

SSDS

	TETANO
ſ	
Inte l	optane
PM	

DRAM

L1, 121PMEM

Tape Prives	D		\searrow
Property	SSD	DRAM	Intel PM
Sequential read latency (ns) Random read latency (ns) Store + flush + fence (ns)	~ 000	81 81 86	169 (2.08×) 305 (3.76×) 91 (1.05×)
Read bandwidth (GB/s) Write bandwidth (GB/s)		120 80	39.4)(0.33×) 13.9 (0.17×)
		proven pc from pc	nigher w Ic 2 4

Prior benchmarking study



BACKGROUND: FILE SYSTEM API

POSIX LPI for files int fd = open(char *path, int flag, mode_t mode) _, open a file, read(int fd, void *buf, size_t nbyte) _____ read a FD write(int fd, void *buf, size_t nbyte) _____ close(int fd)

rename(char *old, char *new) \longrightarrow more a file to a diff breation if breation if the part any state that has been befored

MOTI	VATION: 0	VERHEA)S	Actual appen
		Total time q us	press 8 M	s ~ pu
o virtual	File system	Append Time (<i>ns</i>)	Overhead (ns)	Overhead (%)
	ext4 DAX	9002	8331	(1241%)
	PMFS	4150	3479	518%
	NOVA-Strict	3021	2350	350%
	SplitFS-Strict	1251	580	86%

1160

SplitFS-POSIX

There overheeds are ok for slower storage devices!

488

DAX: mmap pages into virtual memory.

No page caches!

7

73%

SPLIT-FS: GOALS

Low software overhead - - compared to the yest on device

Transparency -> Compatible with existing applications

Minimal data-copy/IO

Flexible semantics

SPLIT FS DESIGN: READ/WRITES



SPLIT FS DESIGN: READ/WRITES



SPLIT FS DESIGN: APPEND



Appends are forwarded to a staging file

flync () - relinked to the original

file mmaps also tracks Staging files used for appends

RFI INK



deta from staging file to the original Avoid file

Actadota La Director La Sizes	y hier Acc	ravchy etc.	SPLI	T-FS M	ODES	Sync -s call returns changes are visible Atomic -> all changes happen simultaneously
	Mode	Sync. Data Ops	Atomic Data Ops	Sync. Metadata Ops	Atomic Metadata Ops	Equivalent to
	POSIX	X	X	X	1	ext4-DAX
	sync	1	×			Nova-Relaxed, PMFS
	strict	1	✓ 	1	✓	NOVA-Strict, Strata
		append	(10),4	append (4)	fsync (fd)	-> Atomicity 10 and 4



SUMMARY

Persistent Memory: New opportunities, new challenges

Split-FS: split Pipelining to use CPU, GPU Partition buffer, BETA ordering Split -FS: Partition FS between User Kernel Date Metadata - Support for efficient appends, atomic Ops usig Staging files

DISCUSSION https://forms.gle/8TwGgqXhVyuiRCpx8

م	ppends be	ar ande	e fut	ter allo
System call	Strict	Sync	POSIX	ext4 DAX
open	2.09	2.08	1.82	1.54
close	0.78	0.69	0.69	0.34
append	3.14	3.09	2.84	11.05
fsync	6.85	6.80	6.80	28.98
read	4.57	4.53	4.53	5.04
unlink	14.60	13.56	14.33	8.60

flyne is more expensive in ext4 -DAX

Ly Split FS la no copy relink

Table 6: SPLITFS system call overheads. The table compares the latency (in us) of different system calls for various modes of SPLITFS and ext4 DAX.

OK! because

data operations are frequent

experises in experises in Spatters need to entries be removed. Data operations are faiter

Metadata sperations are slower

In what ways can SplitFS improve performance of Big Data frameworks like MR/Spark?

NEXT STEPS

Next class:TPU

Project check-ins tonight!

DISCUSSION

Staging files in DRAM?

Page faults are expensive on open()