

REVIEW, SUMMARY

Shivaram Venkataraman CS 537, Spring 2020

ADMINISTRIVIA

Project 5 grades on the way. \longrightarrow Regrade request by early Final Exam: \longrightarrow Monday !! Everything up to including NFS. \leftarrow up to & including the previous May 4th, 10:05am-12:05pm becture Lipiazza ion today! Ly Study / Practice for final No discussion today!

AGENDA / LEARNING OUTCOMES

What are some alternate designs for networked filesystems? NFS -> wid for

What is the role of OS in context of new trends like cloud computing?

RECAP

hent Serv	DISTR	IBUTE	ED SYSTE	MS	
	× C	2			+
		Go	odle	Gmail Images	Sign in
			gic	Ļ	
		Google Search	I'm Feeling Lucky		



NFS SUMMARY

NFS handles client and server crashes very well; robust APIs that are: `

- stateless: servers don't remember clients _____ File handle & dient tracks offsets
- idempotent: doing things twice never hurts

Client retry requests to bouble faibure of server Caching and write buffering is harder, especially with crashes

- Problems: stale cartes when do writes become visible Oroblems: stale cartes when do cartes get invelidated
 - Consistency model is odd (client may not see updates until 3s after file closed)
 - Scalability limitations as more clients call stat() on server

attribute ache duration

Server

re,

CS D

ALTERNATE DESIGN: ANDREW FILE SYSTEM (AFS) Project at CMU our lab machines !!

Open - Close semantics

WHOLE-FILE CACHING

Upon open, AFS client fetches whole file (even if huge), storing in local memory or disk \rightarrow Upon close, client flushes file to server (if file was written) $\frac{1}{2}$ real (helloic, 4) server C1helloc Convenient and intuitive semantics: entre AFS needs to do work only for open/close Reads/writes are local 3 read belloic Use same version of file entire time between open and close Prob > reads / writes incur no network traffic Cons > might need to fetch whole file even if you update just a few bytes.





AFS solution: Tell clients when data is overwritten

- Server must remember which clients have this file open right now Trade-off's

When clients cache data, ask for "callback" from server if changes

 Clients can use data without checking all the time reduces menter

Server no longer stateless!



Mary G. Baker, John H. Hartman, Michael D. Kupfer, Ken W. Shirriff, and John K. Ousterhout

WORKLOAD PATTERNS (1991) UIJY university

	N	cs derign			Camp
	, r	7			
File Usage	Type of Transfer	Acce	esses (%)	В	ytes (%)
	Whole-file	(78)	(64-91)	89	(46-96)
Read-only	Other sequential	19	(7-33)	5	(2-29)
	Random	3	(1-5)	7	(2-37)
	Whole-file	67	(50-79)	69	(56-76)
Write-only	Other sequential	29	(18-47)	19	(4-27)
	Random	4	(2-8)	11	(4-41)
Read/write	Whole-file	0	(0-0)	0	(0-0)
	Other sequential	0	(0-0)	0	(0-0)
	(Random)	100	(100-100)	100	(100-100)



OS/FILESYSTEMS FOR THE CLOUD?

FROM MID 2006

Rent virtual computers in the "Cloud" ~ 11

On-demand machines, spot pricing



Or	on demand AMAZON EC2 (2018)						
Ŭ	Machine	Memory (GB)	Compute Units (ECU)	Local Storage (GB)	Cost / hour		
	t2.nano	0.5		0	\$0.0058		
_	<mark>r5d.24</mark> xlarge	244 768	104 -96	4x900 NVMe	\$6.912		
-	x1.32xlarge	2 TB	4 * Xeon E7	3.4 TB (SSD)	\$13.338		
-	<mark>p3</mark> .16xlarge	488 GB	8 Nvidia Tesla V100 GPUs	0	\$24.48		





Google data centers in The Dulles, Oregon

DATACENTER EVOLUTION

Capacity:

~10000 machines



Bandwidth: 12-24 disks per node



Latency: 256GB RAM cache OY 27B RAM





Some of our customers have been d Worth Data Center. Others of you minimerruption like this is not up to our function like this from occurring in the function of the function



Worth Data Center. Others of you mi More on today's Gmail issue

Posted: Tu		Sign U
Posted by		
Gmail's w		Entire Site 🝷
people rel problem v	Amazon EC2 and Amazon BDS Service Disru	Intion
a list of th		puon

The Joys of Real Hardware

Typical first year for a new cluster:

- ~0.5 overheating (power down most machines in <5 mins, ~1-2 days to recover)
- ~1 PDU failure (~500-1000 machines suddenly disappear, ~6 hours to come back)
- ~1 rack-move (plenty of warning, ~500-1000 machines powered down, ~6 hours)
- ~1 network rewiring (rolling ~5% of machines down over 2-day span)
- ~20 rack failures (40-80 machines instantly disappear, 1-6 hours to get back)
- ~5 racks go wonky (40-80 machines see 50% packetloss)
- ~8 network maintenances (4 might cause ~30-minute random connectivity losses)
- ~12 router reloads (takes out DNS and external vips for a couple minutes)
- ~3 router failures (have to immediately pull traffic for an hour)
- ~dozens of minor 30-second blips for dns
- ~1000 individual machine failures

~thousands of hard drive failures

slow disks, bad memory, misconfigured machines, flaky machines, etc.

Long distance links: wild dogs, sharks, dead horses, drunken hunters, etc.



FEEDBACK! $\sim 41\%$ https://aefis.wisc.edu/

I.What was one idea or concept that you learnt in this course that you appreciated the most?

2. What are some future opportunities that you look forward to based on content from 537?

ALTERNATE DESIGN: GOOGLE FILE SYSTEM (GFS)

~ 2003 open source ~ 2006

GFS: WORKLOAD ASSUMPTIONS

Building a wieb search index Mapheduce workloads

-> "Modest" number of large files (vs. large number of small files)

Two kinds of reads: Large Streaming and small random

Writes: Many large, sequential writes. No random

High bandwidth more important than low latency



The Datacenter Needs an Operating System

Matei Zaharia, Benjamin Hindman, Andy Konwinski, Ali Ghodsi, Anthony D. Joseph, Randy Katz, Scott Shenker, Ion Stoica University of California, Berkeley

1 Introduction

Clusters of commodity servers have become a major computing platform, powering not only some of today's most popular consumer applications—Internet services such as search and social networks—but also a growing number of scientific and enterprise workloads [2]. This rise in cluster computing has even led some to declare that "the datacenter is the new computer" [16, 24]. However, the tools for managing and programming this new computer are still immature. This paper argues that, due to the growing diversity of cluster applications and users, the datacenter increasingly needs an operating system.¹ and Pregel steps). However, this is currently difficult because applications are written independently, with no common interfaces for accessing resources and data.

~ 20 0

In addition, clusters are serving increasing numbers of concurrent users, which require responsive time-sharing. For example, while MapReduce was initially used for a small set of batch jobs, organizations like Facebook are now using it to build data warehouses where hundreds of users run near-interactive ad-hoc queries [29].

Finally, programming and debugging cluster applications remains difficult even for experts, and is even more challenging for the growing number of non-expert users

DATACENTER OPERATING SYSTEMS



COURSE SUMMARY

OPERATING SYSTEMS: THREE EASY PIECES

Three conceptual pieces

I.Virtualization

2. Concurrency

3. Persistence

VIRTUALIZATION





CONCURRENCY

Events occur simultaneously and may interact with one another

Need to

Manage concurrency with interacting processes 2 process Mare data de abstractions (locks some '

Provide abstractions (locks, semaphores, condition variables etc.)

Correctness: mutual exclusion, ordering

Performance: scaling data structures, fairness

Common Bugs!

L> Deadlocks

PERSISTENCE

Managing devices: key role of OS!

Hard disk drives ---- Device

Rotational, Seek, Transfer time

Disk scheduling: FIFO, SSTF, SCAN

Filesystems API

File descriptors, Inodes _____ metadata

Directories

Hardlinks, softlinks

PERSISTENCE

Very simple FS

Inodes, Bitmaps, Superblock, Data blocks

FFS

Placement in groups, Allocation policy

LFS

Write optimized, Garbage collection

Journaling, FSCK

NFS: Partial failures retry, cache consistency

LOOKING BACK. LOOKING FORWARD future opportunities? I dea or concept you learnt -7 Concurrency is important! CS grad courses (744) Security as a correr Eubernetes Is OS research? dore

NEXT COURSES



THANK YOU!