Facebook F4 & Google Search & **Virtual Machines**

Hung-Wei Tseng





Recap: GFS v.s. WAS v.s. Facebook

	GFS (OSDI 2003)	Facebook Haystack (OSDI 2010)
File organizations	file chunk block	volume needle
System architecture	master chunkserver	directory haystack store
Data updates		
Consistency models	relaxed consistency	
Data formats	files	photo/needle
Replications	intra-cluster replication	RIAD-6 & geo-replication
Usage of nodes	chunk server can perform both	



WAS (SOSP 2011)				
stream extent record				
stream manager extent nodes				
append only updates				
strong consistency				
multiple types of objects				
geo-replication				

separate computation and storage

Recap: What kind of data is f4 optimized for?

- Regarding the type of data that f4 aims at, please identify how many the following statements is/are correct.
 - ① f4 is optimized for most frequently requested data in Facebook services
 - ② f4 is optimized for frequently created, deleted data
 - ③ f4 is optimized for reducing the access latency of long-term storage
 - If 4 is optimized for read-only data
 - If 4 is optimized for data that are not accessed very frequently.
 - A. 1
 - B. 2
 - C. 3
 - D. 4
 - E. 5



"Temperature" of data





Recap: Facebook storage architecture



Recap: Storage efficiency

- Reed-Solomon erasure coding
 - Strips: 10GB data + 4GB parity 1.4x space efficiency
 - One volume contains 10 strips
- XOR Geo-replication
 - Use XOR to reduce overhead further (e.g., Azure makes full copies)
 - Block A in DC1 + block B in DC2 -> parity block P in DC3
 - Any two blocks can be used to generate the third
 - 1.5x space efficiency
- 1.4*1.5 = 2.1x space efficiency in total





data center 3

Current scoreboard





- Facebook's F4 (cont.)
- Google Search
- Virtual machines

Fault tolerance

- 1%-2% HDD fail in a year
 - replicate data across multiple disks
 - Use erasure coding for storage efficiency
 - n blocks -> n + k blocks, can tolerate k simultaneous failures
 - higher cost for recovering data when there is a failure
- Host failures (periodically)
 - replicate coded blocks on different hosts
- Rack failures (multiple times/year)
 - replicate coded blocks on different racks
- Datacenter failures (rare, but catastrophic)
 - replicate blocks across data centers
 - use XOR to reduce overhead further (e.g., Azure makes full copies)
 - block A in DC1 + block B in DC2 -> parity block P in DC3
 - any two blocks can be used to generate the third
- Index files
 - use normal triple replication (tiny, little benefit in coding them)

What happens if fault occurs?

- Drive fails
 - Reconstruct blocks on another drive
 - Heavy disk, Network, CPU operation
 - one in background
- During failure, may need to reconstruct data online
 - rebuilder node reads BLOB from data + parity, reconstructs
 - only reads + reconstructs the BLOB (40KB), not the entire block (1GB)



Performance of f4





- Each cell contains 14 racks of 15 hosts, each host contains 30 4TB H.D.Ds.
- A unit of acquisition, deployment
- Storage for a set of volumes
- Similar to the idea of stamps

Recap: GFS/WAS



Partition layer					
Stream layer					
Stream Manager					
ktent ode	Extent node	Extent node	Extent node		
ktent Iode	Extent node	Extent node	Extent node		

Web search for a planet: The Google cluster architecture Luiz Andre Barroso, Jeffery Dean ; Urs Holzle Google

Google search architecture

- How many of the following fulfill the design agenda of the Google search architecture described in this paper?
 - ① Reduce the hardware cost by using commodity-class and unreliable PCs
 - ② Use RAID to provide efficiency and reliability
 - ③ Use replication for better request throughput and availability
 - ④ Optimize for the peak performance
 - A. 0
 - B. 1
 - C. 2

D. 3

E. 4



Google search architecture

- How many of the following fulfill the design agenda of the Google search architecture described in this paper?
 - ① Reduce the hardware cost by using commodity-class and unreliable PCs
 - ② Use RAID to provide efficiency and reliability
 - ③ Use replication for better request throughput and availability
 - ④ Optimize for the peak performance
 - A. 0
 - B. 1
 - C. 2

D. 3

E. 4



Google search architecture

- How many of the following fulfill the design agenda of the Google search architecture described in this paper?
 - Reduce the hardware cost by using commodity-class and unreliable PCs
 - ² Use RAID to provide efficiency and reliability
 - Use replication for better request throughput and availability replica, replica, replica
 - ④ Optimize for the peak performance for performance per dollar

A. 0 **B**. 1 C. 2 D. 3 E. 4

- Price/performance beats peak performance. We purchase the CPU generation that currently gives the best performance per unit price, not the CPUs that give the best absolute performance.
- Using commodity PCs reduces the cost of computation. As a result, we can afford to use more computational resources per query, employ more expensive techniques in our ranking algorithm, or search a larger index of documents.

- Software reliability. We eschew fault-tolerant hardware features such as redundant power supplies, a redundant array of inexpensive disks (RAID), and highquality components, instead focusing on tolerating failures in software.
- Use replication for better request throughput and availability. Because machines are inherently unreliable, we replicate each of our internal services across many machines. Because we already replicate services across multiple machines to obtain sufficient capacity, this type of fault tolerance almost comes for free.



Also reliability and fault-tolerance

Why Google Search Architecture?

- The demand of performing search queries efficiently
 - Each query reads hundreds of MBs of data
 - Support the peak traffic would require expensive supercomputers or high-end servers
- We need a cost-effective approach to address this demand
 - Google search is compare against "AltaVista" search engine that uses DEC's high-performance alpha-based multiprocessor systems
 - AltaVista is later acquired by Yahoo! and you know the later story...



What Google proposes?

- Using commodity-class / unreliable PCs
- Provide reliability in software rather in hardware
- Target the best aggregate request throughput, not peak server response time



Google query-serving architecture





Replication is the key

- Scalability: simply add more replicas, the service capacity can improve
- Availability: even though one machine fails, another replica to take over



Poll close in 1:30

What kind of processors Google search needs

- If we are designing a processor just for Google search or similar type of applications, how many of the following targets/features would fulfill the demand?
 - ① Can execute many instructions from the same process/thread simultaneously
 - ② Can execute many processes/threads simultaneously
 - ③ Can predict branch outcome accurately
 - ④ Have very large cache capacity
 - A. 0
 - B. 1
 - C. 2
 - D. 3

E. 4

Poll close in 1:30

What kind of processors Google search needs

- If we are designing a processor just for Google search or similar type of applications, how many of the following targets/features would fulfill the demand?
 - ① Can execute many instructions from the same process/thread simultaneously
 - ② Can execute many processes/threads simultaneously
 - ③ Can predict branch outcome accurately
 - ④ Have very large cache capacity
 - A. 0
 - B. 1
 - C. 2
 - D. 3

E. 4

What kind of processors Google search needs

- If we are designing a processor just for Google search or similar type of applications, how many of the following targets/features would fulfill the demand?
 - Can execute many instructions from the same process/thread simultaneously (1)
- Can execute many processes/threads simultaneously
- Can predict branch outcome accurately
- ④ Have very large cache capacity
- A. 0
- **B**. 1
- C. 2
- D. 3

E. 4

given how little ILP our application yields, and shorter pipelines would reduce or eliminate branch mispredict penalties. The avail-

For such workloads, a memory system with a relatively modest sized L2 cache, short L2 cache and memory latencies, and longer (perhaps 128 byte) cache lines is likely to be the most effective.

end ones. Exploiting such abundant threadlevel parallelism at the microarchitecture level appears equally promising. Both simultaneous multithreading (SMT) and chip multiprocessor (CMP) architectures target thread-level parallelism and should improve the performance of many of our servers. Some early

none concurrencity and mas superior branch prediction logic. In essence, there isn't that much exploitable instruction-level parallelism (ILP) in the workload. Our measurements suggest that the level of aggressive out-oforder, speculative execution present in modern processors is already beyond the point of diminishing performance returns for such programs.

Hardware

- Processor
 - Index search has little ILPs doesn't need complex cores
 - Index search can be highly parallelized processors with threadlevel parallelism would be a good fit (e.g. Simultaneous) Multithreading, SMT and Chip Multicprocessor, CMP)
 - Branch predictor matters
- Memory: Good spatial locality. Moderate cache size will suffice
- Storage: No SCSI, No RAID not worth it
- Power: is an issue, but only \$1,500/mo operating bill vs \$7,700 capital expense

Will their architecture work for other things?

As mentioned earlier, our infrastructure consists of a massively large cluster of inexpensive desktop-class machines, as opposed to a smaller number of large-scale sharedmemory machines. Large shared-memory machines are most useful when the computation-to-communication ratio is low; communication patterns or data partitioning are dynamic or hard to predict; or when total cost of ownership dwarfs hardware costs (due to management overhead and software licensing prices). In those situations they justify their high price tags.

At Google's scale, some limits of massive server parallelism do become apparent, such as the limited cooling capacity of commercial data centers and the less-than-optimal fit of current CPUs for throughput-oriented applications. Nevertheless, using inexpensive PCs to handle Google's large-scale computations has drastically increased the amount of computation we can afford to spend per query, thus helping to improve the Internet search experience of tens of millions of users. A t first sight, it might appear that there are few applications that share Google's characteristics, because there are few services that require many thousands of servers and petabytes of storage. However, many applications share the essential traits that allow for a PC-based cluster architecture. As long as an application orientation focuses on the price/performance and can run on servers that have no private state (so servers can be replicated), it might benefit from using a similar architecture. Common examples include highvolume Web servers or application servers that are computationally intensive but essentially stateless. All of these applications have plenty of request-level parallelism, a characteristic exploitable by running individual requests on separate servers. In fact, larger Web sites already commonly use such architectures.

Metrics we care about data center design

- Costs machine architecture, distributed system architecture, replication strategies
- Power machine architecture
- Energy machine architecture
- Space-efficiency erasure coding, replication, distributed
- Throughput replication, distributed
- Reliability replication

Virtual Machines



Taxonomy of virtualization





Virtual machine architecture

Applications

Guest OS

Virtual Machine Monitor

The Machine





Three Laws of Robotics

- A robot may not injure a human being or, through inaction, allow a human being to come to harm.
- A robot must obey orders given it by human beings except where such orders would conflict with the First Law.
- A robot must protect its own existence as long as such protection does not conflict with the First or Second Law.





Back to 1974...

Formal Requirements for Virtualizable Third Generation Architectures

Gerald J. Popek University of California, Los Angeles and Robert P. Goldberg Honeywell Information Systems and Harvard University

A virtual machine is taken to be an efficient, isolated duplicate of the real machine. We explain these notions through the idea of a virtual machine monitor (VMM). See Figure 1. As a piece of software a VMM has three essential characteristics. First, the VMM provides an environment for programs which is essentially iden-Fidelity tical with the original machine; second, programs run in this environment show at worst only minor decreases **Performance** in speed; and last, the VMM is in complete control of Safety and isolation system resources.

. . . . AL a Grat

Recap: virtualization However, we don't want everything to pass through this API!









Recap: privileged instructions

- The processor provides normal instructions and privileged instructions
 - Normal instructions: ADD, SUB, MUL, and etc ...
 - Privileged instructions: HLT, CLTS, LIDT, LMSW, SIDT, ARPL, and etc...
- The processor provides different modes
 - User processes can use normal instructions
 - Privileged instruction can only be used if the processor is in proper mode

Least privileged



Ring 3

Ring 2

Ring 1

Ring 0

Kernel

Device Drivers

Device Drivers

Recap: How applications can use privileged operations?

- Through the API: System calls
- Implemented in "trap" instructions
 - Raise an exception in the processor
 - The processor saves the exception PC and jumps to the corresponding exception handler in the OS kernel



user mode

kernel/privileged mode

Hosted virtual machine





Applications







Virtual machine monitors on bare machines





Applications



Three main ideas to classical VMs

- De-privileging
- Primary and shadow structures
- Tracing



CPU Virtualization: Trap-and-emulate





Announcement

- Group photo next lecture!
- Project revision
 - Allows you to revise your project with 20% of penalty on the unsatisfactory parts/test cases after the first-round of grading (firm deadline 3/11)
 - Say you got only 60% in the first-round, and you fixed everything before 3/11 - you can still get 60% + 80% * 40% = 92%
- iEVAL count as an extra, full-credit reading quiz, due 3/11
- Final contains two parts (each account for 50%)
 - Part 1: unlimited time between 3/12-3/17, open-ended questions
 - Part 2: 80 minute multiple choices/answers questions + two problem sets of comprehensive exam questions

Computer Science & Engineering





