

Distributed Operating Systems

Prof. Chi-Sheng Shih
Graduate Institute of Networking and Multimedia
Department of Computer Science and Information
Engineering
National Taiwan University

High Performance Computing

Bit-coin

Cloud Computing

大量運算、儲存

Are they the same?

Cloud Computing = High Performance Computing?

- Shared properties:
 - Large amount of computation resources are interconnected to provide coherence services.
 - Located in a server room/data center and connected via networks.

機房

Workloads

- Real-time Weather forecast
- Nuclear fusion research
- Stock trading
- Facebook
- Online Gaming

Workloads

Both need large amount of computation

HPC

- Real-time Weather forecast:
 - Large amount of data, short latency (< 10s)
- Nuclear fusion research
 - Low latency, generating large amount of data during and after the work load.
- Stock trading
 - Real-time response (< 10⁻²s), 10⁶ requests per second, guaranteed ordering.
- Facebook
 - 10³ of participants per message, long latency, guaranteed ordering, and number of messages increase over time.
- Online Games
 - Short latency (< 1s), 10³ players per game, number of games increase over time.

CC

Scalability

要很容易擴充

HPC vs. Cloud Computing

- High performance computing: **Performance**
 - The majority of the workloads are computation intensive and can only tolerate short latency among sub-workloads.
 - The systems are built with high performance processors, and high bandwidth bus. However, it is not easy to deploy additional computation resources.
- Cloud Computing: **Scalability (可大可小) Both** < **scale up** **scale down**
 - The majority of the workloads can be partitioned and conducted independently.
 - The systems are built with low cost processors, and computer networks. However, it is designed to add/remove computation resources at any time.
 - The performance are improved by adding more computation resources.
- HPC and Cloud Computing are distributed computing in general.

Why distributed computing systems?

Network 高速

- Personal computers are cheap and powerful.
- Why bother to use distributed computing systems?
 - (Do you use peer-to-peer file sharing/streaming?)
 - Broadband connection is becoming popular
 - Inherently distributed applications
 - Communication and resource sharing possible
 - Economics – price-performance ratio
 - Higher reliability
 - Scalability
 - Potential for incremental growth
- What should be done to make it possible/better?
 - Distribution-aware platforms, operating systems and applications.
 - Security and privacy.

concern

成本低

要重新思考

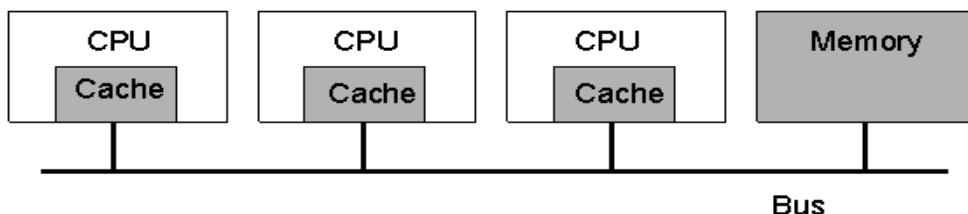
Distributed Computing System Models

- A distributed system:
 - Multiple connected processors/computing devices working together.
chip chip + memory + storage
 - A collection of independent computers that appear to its users as a single coherent system *對 user 來說像一台*
- Examples of distributed computing system models:
 - Minicomputer Model
 - Workstation Model
 - Workstation-server Model
 - Processor-pool Model
 - Hybrid Model



Hardware Concepts: Multiprocessors (1)

- Multiprocessor dimensions
 - Memory: could be shared or be private to each CPU
 - Interconnect: could be shared (bus-based) or switched
- A bus-based multiprocessor.



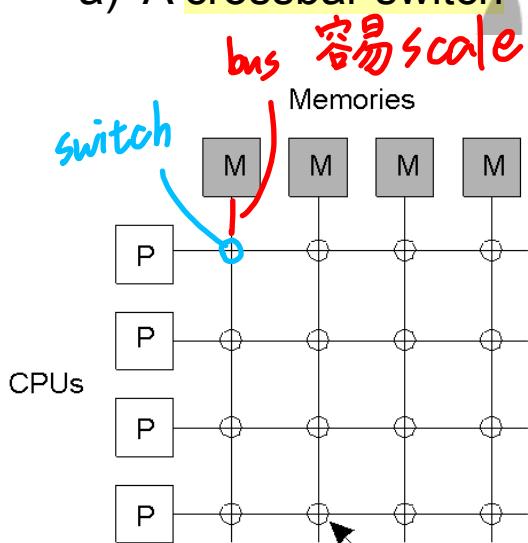
- Q: What are the potential problems for bus-based multiprocessor?

pros: 容易 share memory (∴ public memory)

cons: 修改(先後順序), 不易 scale up.
Bandwidth 固定

Multiprocessors (2)

a) A crossbar switch



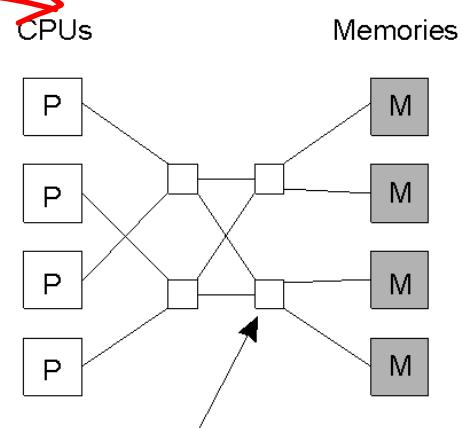
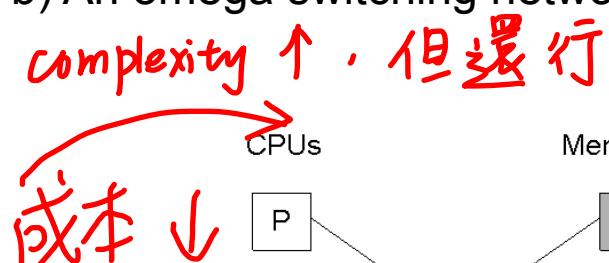
成本高:
switch latency 要很低

(a)

一個 instruction $\approx 10^{-9} - 10^{-12}$ s

代表 switch 要更快，且小

b) An omega switching network



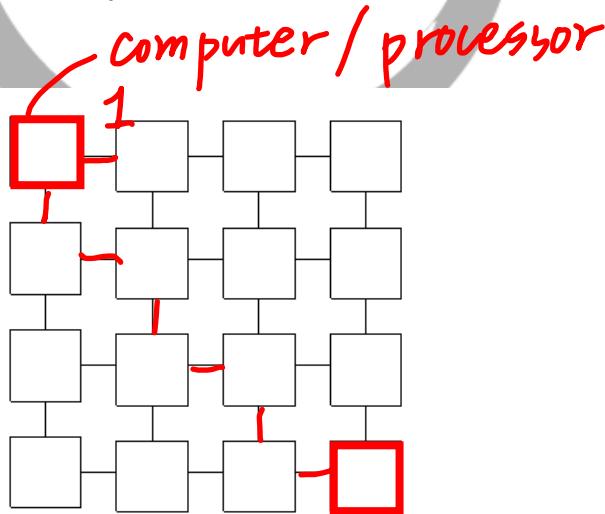
(b)

Data center 常見

Homogeneous Multicomputer Systems

最短: 1
最长: 6

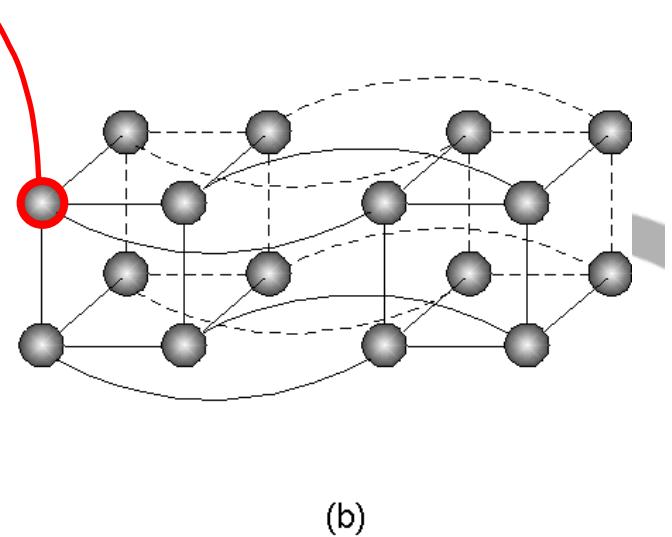
a) Grid



latency
不固定 (大家一起慢)

最短: 1
最长: 4

b) Hypercube



(b)

Distributed Systems Models

- Minicomputer model (e.g., early networks)
 - Each user has local machine **CPU + memory + storage**
 - Local processing but can fetch remote data (files, databases) e.g. NAS
- Workstation model (e.g., Sprite) **remote端口 computing**
 - Processing can also migrate **local沒做什麼計算**
- Client (workstation)- server Model (e.g., V system, world wide web) **HTML, JavaScript**
 - Each user has local workstation
 - Powerful workstations serve as servers (file, print, DB servers)
- Processor pool model (e.g., Amoeba, Plan 9) **一般user少用**
 - Terminals are Xterms or diskless terminals
 - Pool of backend processors handle processing
- Cloud Computing
 - Relationship between client and servers changes over time.
 - Computation capacity on servers are dynamically adjustable.

有差別

哪個程式跑在 server, client { 是事先定義好的

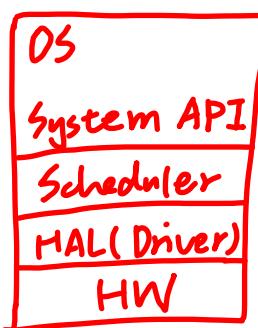
{ 不是 : (scale) }

Distributed Operating Systems

- What's an operating system?
 - To present users with a virtual environment that is easier to program than the underlying hardware.
 - To manage the various resources of the system.

Uniprocessor Operating Systems

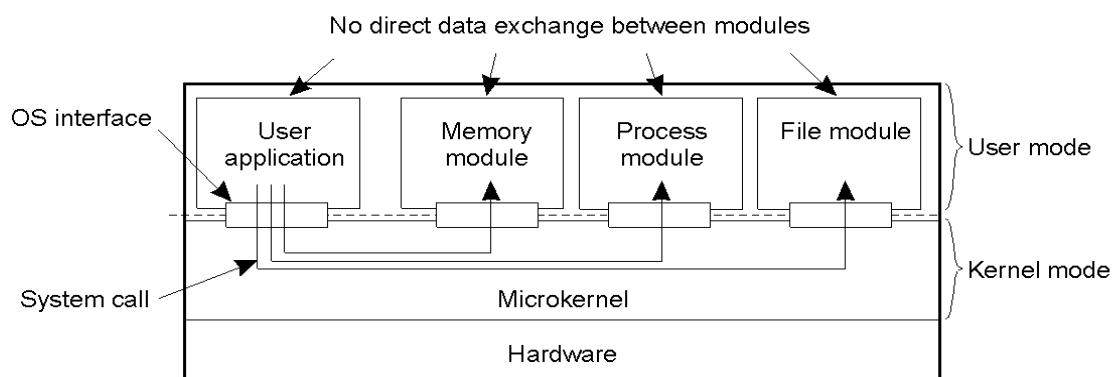
- An OS acts as a resource manager or an arbitrator
 - Manages CPU, I/O devices, memory *: multi-proc*
- OS provides a virtual interface that is easier to use than hardware
- Structure of uniprocessor operating systems
 - **Monolithic** (e.g., MS-DOS, early UNIX)
 - One large kernel that handles everything
 - Layered design
 - Functionality is decomposed into N layers
 - Each layer uses services of layer N-1 and implements new service(s) for layer N+1



Uniprocessor Operating Systems

Microkernel architecture

- Small kernel
- user-level servers implement additional functionality



Distributed Operating Systems

- The operating system for distributed computing systems:
 - Network operating systems
 - Distributed operating systems

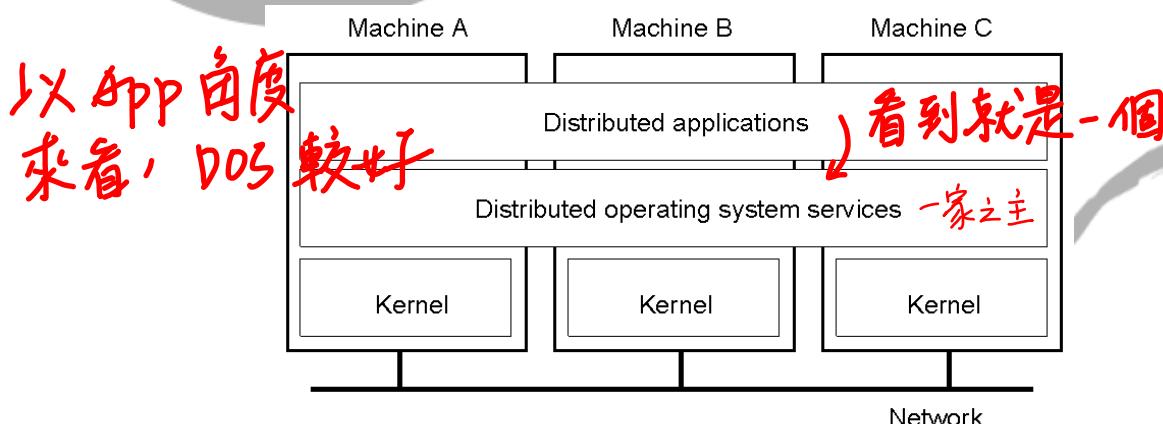
Types of OSs for distributed systems

System	Description	Main Goal
DOS Distributed	\exists consensus <u>Tightly-coupled</u> operating system for multi-processors and <u>homogeneous</u> multicomputers	Hide and manage hardware resources
NOS Network	每個人可以做自己的決定 <u>Loosely-coupled</u> operating system for heterogeneous multicomputers (LAN and WAN)	Offer local services to remote clients
Middleware	Additional layer atop of NOS implementing general-purpose services	Provide distribution transparency

Application
middleware
NOS

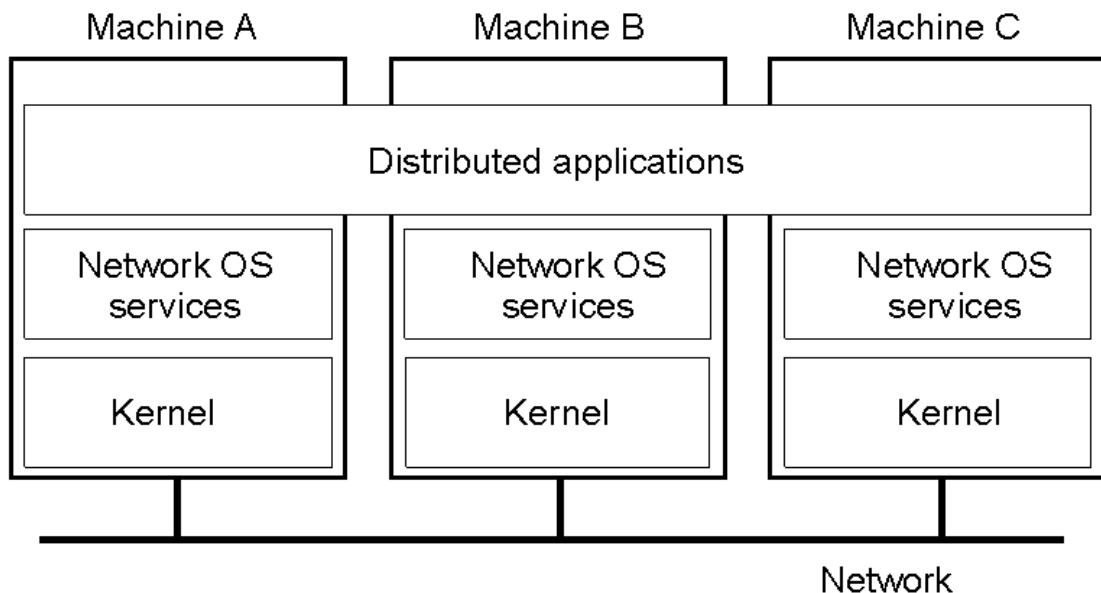
Distributed Operating Systems

- Each computing unit has its own hardware (including memory, CPU, and IO devices)
- A distributed operating system
 - looks like a uni-processor operating system but operates on multiple independent computing units
 - manages multiple computing units transparently to the user



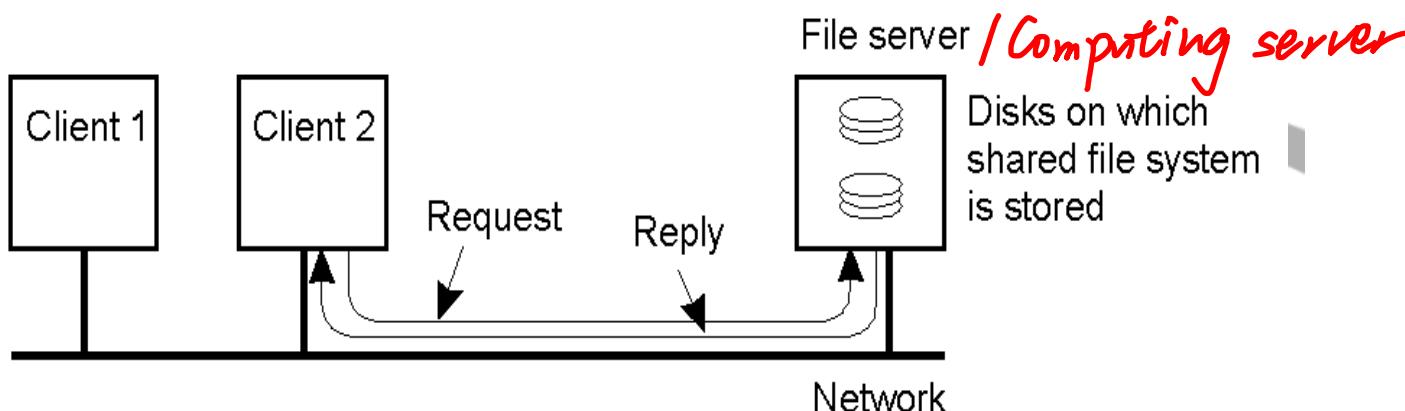
Network Operating System

以實作來講，NOS
“只要在座各位，用的 Ubuntu, Windows 等 NOS”



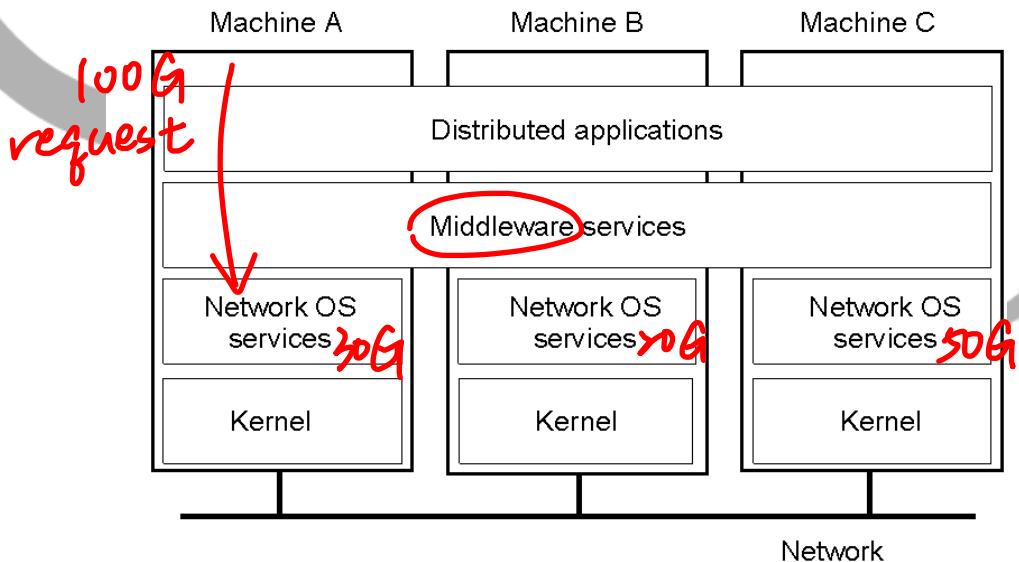
Network Operating System

- Employs a client-server model
 - Minimal OS kernel
 - Additional functionality as user processes



Middleware-based Systems

- General structure of a distributed system as middleware.



Network OS vs. Distributed OS

	Network OSs	Distributed OSs
System Image	A collection of distinct machines connected by a communication subsystem.	A virtual uniprocessor system.
Autonomy 自主性	Each computer functions independently of other computers.	Processes and resources are managed globally.
Fault tolerance capability	Little or no fault tolerance.	High fault tolerance.

Goals of developing distributed systems

- Connecting Users and Resources
- Transparency
- Openness
- Scalability
- Flexibility
- Reliability
- Performance

Transparency in Distributed Systems

Transparency	Description
Access	
Location	
Migration	
Relocation	
Replication	
Concurrency	
Failure	
Persistence	

Transparency in Distributed Systems

Transparency	Description
Access	Hide differences in data representation and how a resource is accessed user 不需要知道 data 是怎麼存的才可以 access
Location	Hide where a resource is located
Migration	Hide that a resource may move to another location 這東西這次存取跟下次存取可能不一樣。
Relocation	Hide that a resource may be moved to another location while in use
Replication	Hide that a resource is replicated. 自己要寫 n 份。
Concurrency	Hide that a resource may be shared by several competitive users OS 決定誰先，誰後，確保最後結果是正確的。
Failure	Hide the failure and recovery of a resource RAID
Persistence	Hide whether a (software) resource is in memory or on disk 寫進去後不會因為系統 shutdown 就消失了。

Openness

- An open distributed system offers services according to standard rules.
- Services are generally specified through interfaces in Interface Definition Language.
 - Completeness
 - Neutrality
- Distributed services should have
 - Interoperability and 容易和他人對接
 - Portability VMware, Docker 之間無法互通，所以沒有 portability。

SLC → MLC Reliability

- A fault in a system may cause system failure. Multiple resources may not be able to increase the system reliability.
 - Fail-stop failure Hardware failure 之後，就不會再往下動了。
 - Byzantine failure Software failure : 1. 報了假的情報回來，導致 system crash 掉。
- Fault-handling mechanisms:
 - Fault Avoidance 以硬碟來說，盡量平均 workload，不要把所有事放在同一個碟上。
e.g. wear leveling
 - Fault Tolerance
 - Redundancy technique e.g. RAID (把資料 duplicate 好幾份)
 - Distributed control 多個 algorithms
- Fault Detection and Recovery
 - Atomic transaction 一次完成，要嘛全做，要嘛全不，要記錄一開始什麼狀態，若 detect failure 要能 roll back。e.g. 習大大：中國、台灣不可分割
e.g. Gmap <-> stateful (keep track)
 - Stateless servers e.g. Gmap <-> stateful (keep track)
 - Acknowledgements and timeout-based retransmission of messages 若網路斷線，要能夠重傳。e.g. TCP/IP <-> UDP (沒收到就謝謝再連絡)

Byzantine failure

- A Byzantine Fault is an incorrect operation (algorithm) that occurs in a distributed system that can be classified as:
 - Omission Failure – a failure of not being present such as failing to respond to a request or not receiving a request.
 - Execution Failure or Lying – a failure due to sending incorrect or inconsistent data, corrupting the local state or responding to a request incorrectly.
- Examples
 - Round off errors passed from one function to another and then another, etc.
 - Corrupted system databases where the error is not detected
Compiler errors
 - An undetected bit flip producing a bad message.
- This is a worse case model since the Byzantine Fault can generate misleading information causing a maximum of confusion.

Flexibility

- Why flexibility is an important feature of a distributed operating system?
 - Ease of modification
 - Ease of enhancement 效能上、功能上
- A flexible distributed system should be organized as a collection of relatively small and easily replaceable or adaptable components.
- Separating policy from mechanism
 - Web caching
 - Policy=?
 - Mechanism=? FIFO, stack (FILO)

Performance

- Distributed systems vs. centralize systems
 - Google search

我可以完全掌握我的 OS，也沒有 communication 的 overhead

Scalability Problems

Concept	Example
Centralized services	A single server for all users
Centralized data	A single on-line telephone book
<u>Centralized algorithms</u>	Doing routing based on complete information

Examples of scalability limitations.

我要得到所有資訊才有辦法算

How much you can sort within one minute?

- Sort is a fundamental function for data analysis.
- Minute Sort:
 - Amount of data that can be sorted in 60.00 seconds or less.
 - One of the sort benchmark on sortbenchmark.org

Latest Results for Soft Benchmark

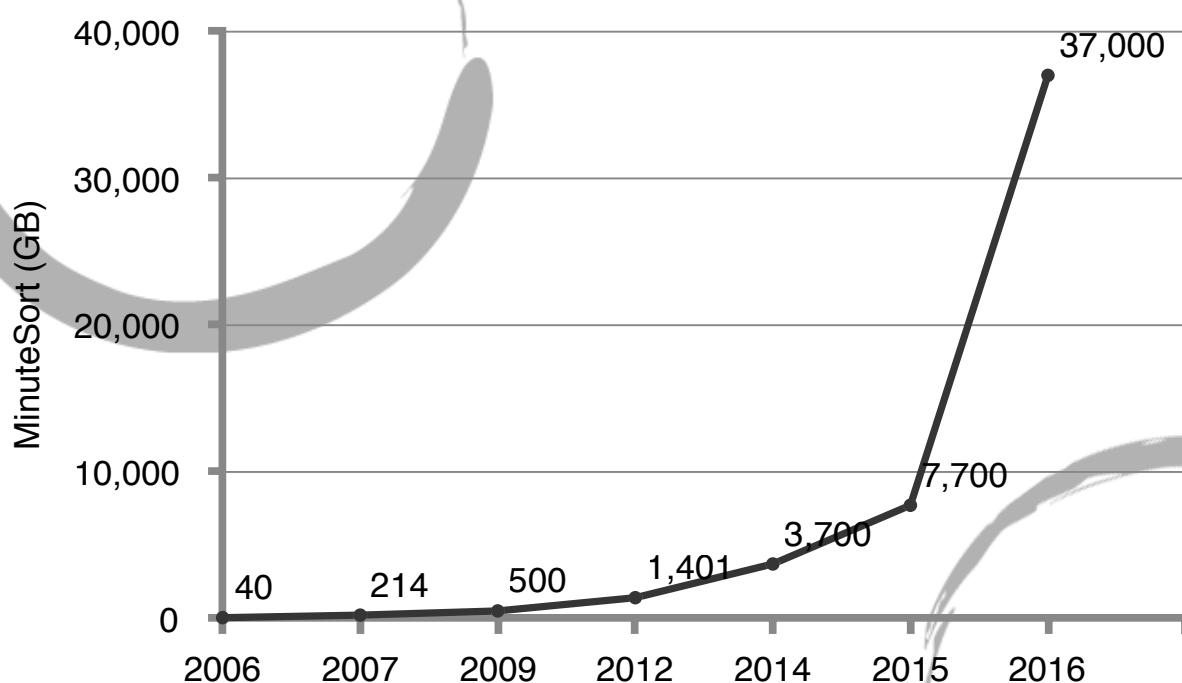
Minute Sort: How much data one can sort in one minute?

Minute	2016, 37 TB Tencent Sort 512 nodes x (2 OpenPOWER 10-core POWER8 2.926 GHz, 512 GB memory, 4x Huawei ES3600P V3 1.2TB NVMe SSD, 100Gb Mellanox ConnectX4-EN) Jie Jiang, Lixiong Zheng, Junfeng Pu, Xiong Cheng, Chongqing Zhao Tencent Corporation Mark R. Nutter, Jeremy D. Schaub	2016, 55 TB Tencent Sort 512 nodes x (2 OpenPOWER 10-core POWER8 2.926 GHz, 512 GB memory, 4x Huawei ES3600P V3 1.2TB NVMe SSD, 100Gb Mellanox ConnectX4-EN) Jie Jiang, Lixiong Zheng, Junfeng Pu, Xiong Cheng, Chongqing Zhao Tencent Corporation Mark R. Nutter, Jeremy D. Schaub
--------	--	--

Gray: How much time it takes to sort 100TB?

Gray	Daytona 2016, 44.8 TB/min Tencent Sort 100 TB in 134 Seconds 512 nodes x (2 OpenPOWER 10-core POWER8 2.926 GHz, 512 GB memory, 4x Huawei ES3600P V3 1.2TB NVMe SSD, 100Gb Mellanox ConnectX4-EN) Jie Jiang, Lixiong Zheng, Junfeng Pu, Xiong Cheng, Chongqing Zhao Tencent Corporation Mark R. Nutter, Jeremy D. Schaub	Indy 2016, 60.7 TB/min Tencent Sort 100 TB in 98.8 Seconds 512 nodes x (2 OpenPOWER 10-core POWER8 2.926 GHz, 512 GB memory, 4x Huawei ES3600P V3 1.2TB NVMe SSD, 100Gb Mellanox ConnectX4-EN) Jie Jiang, Lixiong Zheng, Junfeng Pu, Xiong Cheng, Chongqing Zhao Tencent Corporation Mark R. Nutter, Jeremy D. Schaub
------	--	--

Status of Minute Sort



Status of Minute Sort

- 2016:
 - Winner: Tencent Corporation, China
 - 37 TB TB/min using 512 nodes x (2 OpenPOWER 10-core POWER8 2.926 GHz, 512 GB memory, 4x Huawei ES3600P V3 1.2TB NVMe SSD,
- 2015:
 - Winner: FuxiSort by AliBaba
 - 7.7 TB using 3,134 nodes x (2 Xeon E5-2630 2.30Ghz, 96 GB memory, 12x2 TB SATA HD, 10 Gb/s Ethernet) + 243 nodes x (2 Xeon E5-2650v2 2.60Ghz, 128 GB memory, 12x2 TB SATA HD, 10 Gb/s Ethernet)
- 2014:
 - Winner: DeepSort by Zheng Li, Juhan Lee, Samsung.
 - 3.7 TB using 384 nodes of 2x2.1GHz Intel Xeon, 64GB memory, and 8 HDs
- 2012:
 - Winner: Flat Datacenter Storage from Microsoft Research
 - 1,401 GB using 256 nodes
- 2009:
 - Winner: Hadoop from Yahoo
 - 500GB using 1406 nodes x (2 quad core Xeons, 8GB memory, 4 SATA HDs)
- 2007: 214BGB
- 2006: 40GB
- 2004: 34GB

Scalability Problems

- Decentralized algorithms should be used.
 - No machine has complete information about the system state.
 - Machines make decisions based only on local information.
 - Failure of one machine does not ruin the algorithm.
 - There is no implicit assumption that a global clock exists.
- Geographical scalability:
 - LAN vs. WAN.

Performance

- Distributed systems vs. centralized systems
 - Google search
- How can the performance of a distributed system be as good as a centralized system?
 - Batch if possible
 - Cache whenever possible
 - Minimize copying of data
 - Minimize network traffic
 - Take advantage of fine-grain parallelism for multiprocessing

L Barroso, J Dean, U Hölzle, Web Search for a Planet: The Architecture of the Google Cluster, - IEEE Micro, 2003.

Comparison between Systems

Item 因為有 centralize	Distributed OS		Network OS	Middleware-based OS
	Multiproc	Multicomp		
Degree of transparency	Very High	High	Low	High
Same OS on all nodes	Yes	Yes	No	No
Number of copies of OS	1	N	N	N
Basis for communication	Shared memory	Messages	Files	Model specific
Resource management	Global, central	Global, distributed	Per node	Per node
Scalability	No	Moderately	Yes	Varies
Openness	Closed	Closed	Open	Open

Additional Readings

- David L. Cohn, William P. Delaney, Karen M. Tracey, ARCADE: A Platform for Heterogeneous Distributed Operating Systems, *Proceedings of the Symposium on Experiences with Distributed and Multiprocessor Systems*, 1989, 373-390.
- Dougulis, F., Ousterhout, J.K., Kaashoek, M.F., and Tanenbaum, A.S., *Comparison of Two Distributed Systems: Amoeba and Sprite*, *Computing Systems Journal* 4(Fall), 1991, 353-384.
- L Barroso, J Dean, U Hoezle, *Web Search for a Planet: The Architecture of the Google Cluster*, IEEE Micro, Volume: 23, Issue: 2, March-April, 2003, 22- 28.
- Thain, D., Tannenbaum, T. and Livny, M. (2005), Distributed computing in practice: the Condor experience. *Concurrency and Computation: Practice and Experience*, 17: 323–356. doi: 10.1002/cpe.938
- B Hayes, *Cloud Computing*, Communication of ACM, Vol. 51, Issues 7, July 2008.
- Anil Madhavapeddy, Richard Mortier, Charalampos Rotsos, David Scott, Balraj Singh, Thomas Gazagnaire, Steven Smith, Steven Hand, and Jon Crowcroft. 2013. Unikernels: library operating systems for the cloud. *SIGARCH Comput. Archit. News* 41, 1 (March 2013), 461-472.
- Jiamang Wang, Yongjun Wu, Hua Cai, Zhipeng Tang, Zhiqiang Lv, Bin Lu, Yangyu Tao, Chao Li, Jingren Zhou, and Hong Tang, *FuxiSort*, file: <http://sortbenchmark.org/FuxiSort2015.pdf>

Before next class

- Reading assignment:
 - Anil Madhavapeddy, Richard Mortier, Charalampos Rotsos, David Scott, Balraj Singh, Thomas Gazagnaire, Steven Smith, Steven Hand, and Jon Crowcroft. 2013. Unikernels: library operating systems for the cloud. *SIGARCH Comput. Archit. News* 41, 1 (March 2013), 461-472.
 - Brian N. Bershad, Thomas E. Anderson, Edward D. Lazowska, and Henry M. Levy. 1990. Lightweight remote procedure call. *ACM Trans. Comput. Syst.* 8, 1 (February 1990), 37-55.