

Distributed Systems 600.437 Distributed Operating Systems

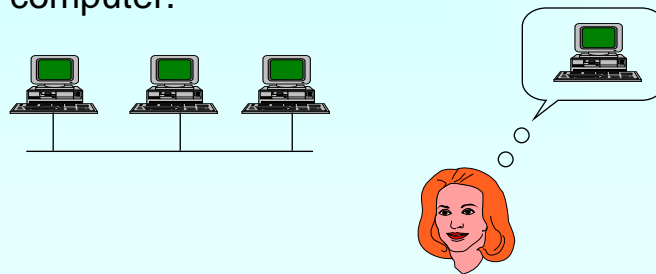
Department of Computer Science
The Johns Hopkins University

Distributed Operating Systems

Lecture 11

A Distributed Operating System

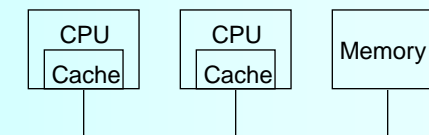
An operating system which manages a collection of independent computers and makes them appear to the users of the system as a single computer.



Hardware Architectures

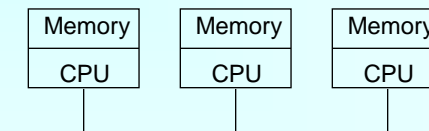
- Multiprocessors
 - Tightly coupled.
 - Shared memory.

Parallel architecture



- • Multicomputers.
 - Loosely coupled.
 - Private memory.
 - Autonomous.

Distributed architecture



Software Architectures

- Multiprocessor OS
 - Looks like a virtual uniprocessor, contains only one copy of the operating system, communication via shared memory, single run-queue
- Network OS
 - Does not look like a virtual uniprocessor, contains n copies of the operating system, communication via shared files, n run-queues
- Distributed OS
 - Looks like a virtual uniprocessor (more or less), contains n copies of the operating system, communication via messages, n run-queues

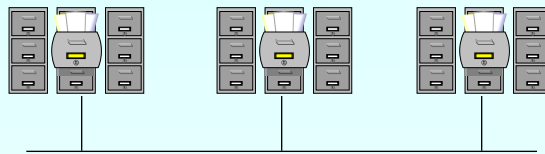
Design Issues

- Transparency
 - Location (processes, cpu(s), files)
 - Replication of files
 - Parallelism
- Performance
 - Throughput / response time
 - Load balancing (static/dynamic)
 - Communication is slow compared to processing speed : Fine grain / coarse grain
- Scalability
- Reliability
- Architecture flexibility.



Distributed File Systems

- File and directory naming
- Semantics of file sharing
- Implementation considerations
 - Caching
 - Update protocols
 - Replication



File and Directory Naming

- Machine + path naming `/machine/path`.
 - one name space, but not transparent.
- Mounting remote file systems onto the local file hierarchy.
 - The view of the file system **may be** different at each computer.
- A single name space that looks the same on all machines.
 - Full naming transparency.

File Sharing Semantics

- One-copy semantics (unix semantics).
 - Updates are written to the single copy and are available immediately.
- Serializability.
 - Transaction semantics (locking files - share for read and exclusive for write).
- Session semantics.
 - Copy the file on open, work on local copy, and copy back on close.

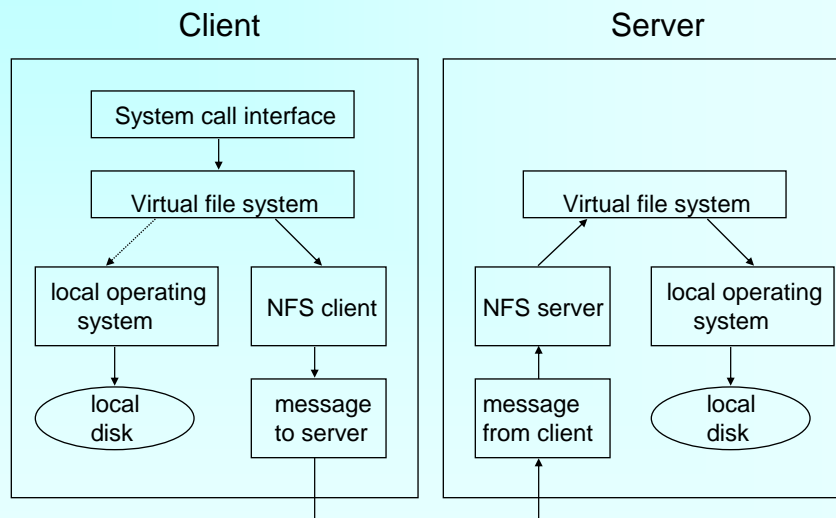
Sun-NFS

- Supports heterogeneous systems.
- Architecture
 - Server exports one or more directory trees for access by remote clients.
 - Clients access exported directory trees by mounting them to the client local tree.
 - Diskless clients can mount exported directory to their root directory.
 - Auto-mount (on the first access).

Sun-NFS (cont.)

- Protocols
 - Mounting protocol
 - Directory and file accessing protocol
 - Stateless
 - No open / close messages.
 - Each read / write message contain the full path and file description position.
- Semantics
 - Not entirely Unix since there is no way to lock files
 - Timing issues

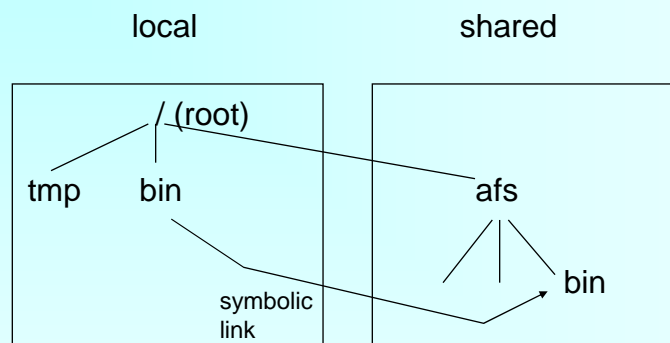
Sun-NFS Structure



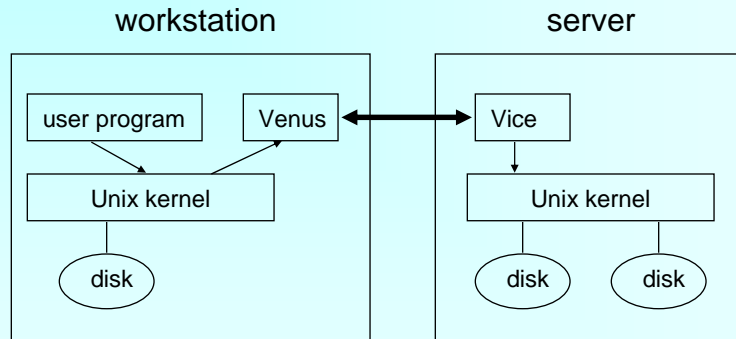
Andrew File System

- Supports information sharing on a large scale (thousands of workstations).
- Uses session semantics.
- The entire file is copied to the local machine (Venus) from the server (Vice) when open. If the file is changed - it is copied to the server when closed.
- The method works because in practice most files are changed by one person only (non-database).
- Measurements show that only 0.4% of changed files were updated by more than one user during one week.

Andrew File System (cont)



AFS Structure



AFS File Validation

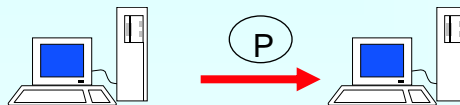
- **Older AFS versions:**
 - On open: the Venus access the vice to see if its copy of the file is still valid. This causes a substantial delay even if the copy is valid.
 - The Vice is stateless
- **Newer AFS versions:**
 - The Vice maintains lists of valid copies. If a file is modified, the Vice invalidates other copies.
 - On open: if the Venus has a valid copy it can open it immediately.
 - If Venus crashes it has to invalidate its version or check their validity.

The Coda File System

- Descendant of AFS that is substantially more resilient to server and network failures.
- Support for “mobile” users.
- Directories are replicated in several servers (Vice)
- When the Venus is disconnected, it uses local versions of files. When Venus reconnects, it reintegrates using optimistic update scheme.

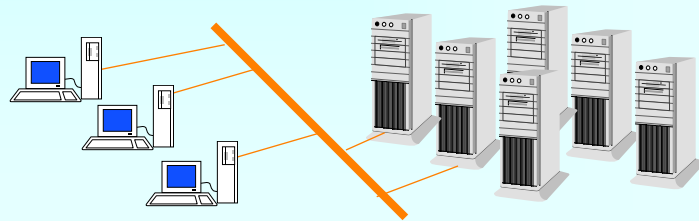
Process Migration

- Load Balancing
 - Static load balancing - CPU is determined at process creation.
 - Dynamic load balancing - processes dynamically migrate to other computers to balance the CPU (or memory) load.
- Migration architecture
 - One image system
 - Point of entrance dependent system (the deputy concept)

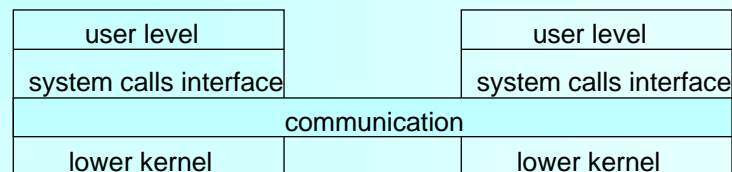


A Mosix Cluster

- Mosix: Kernel level enhancement to Linux that provides dynamic load balancing in a network of workstations.
- Dozens of PC computers connected by local area network (Fast-Ethernet and up).
- Any process can migrate *anywhere anytime*.
- www.mosix.org



An Architecture for Migration

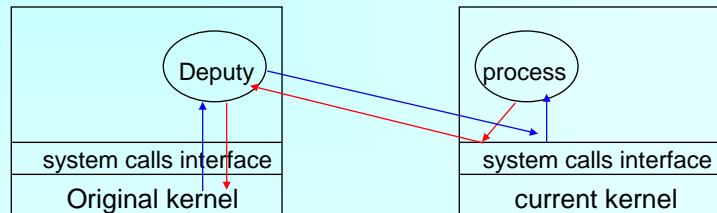


Architecture that fits one system image.
Needs location transparent file system.

(Mosix previous versions)

Architecture for Migration (cont.)

The process migrated from the original computer to the current computer and now is initiating a system call.



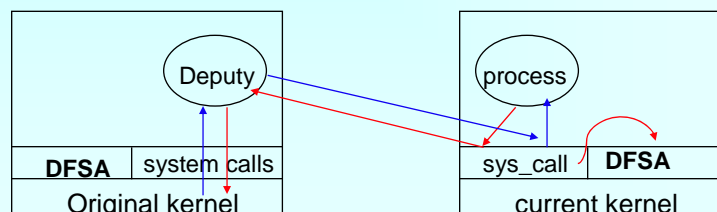
- Architecture that fits entrance dependent system.
 - Easier to implement based on current Unix.
- (Mosix current versions)

Mosix: File Access

Regularly, each file access must go back to deputy...

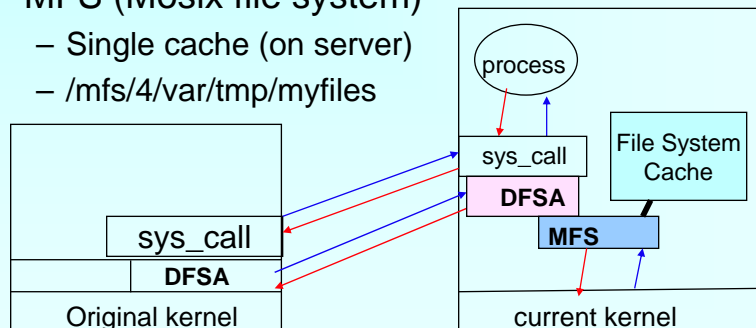
➔ Very Slow for I/O apps.

Solution-- Allow processes to access a distributed file system through the current kernel.

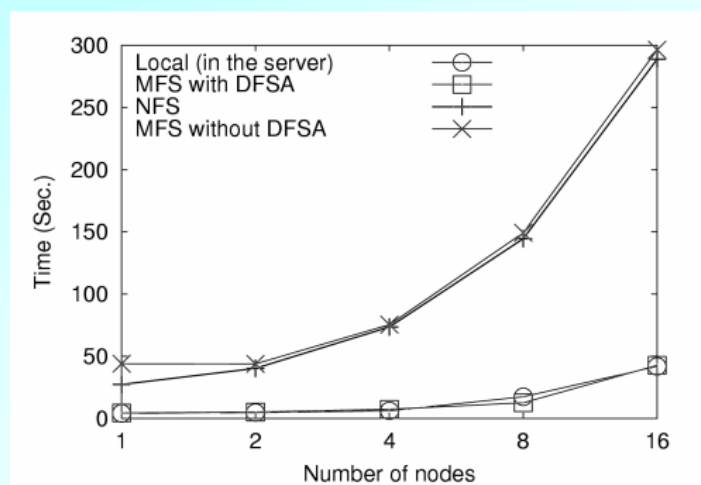


Mosix: File Access

- DFSA (direct file system access)
 - Requirements (cache coherent, monotonic timestamps, files not deleted until all nodes finished)
 - Bring the process to the files.
- MFS (Mosix file system)
 - Single cache (on server)
 - /mfs/4/var/tmp/myfiles



Mosix: DFSA/MFS Performance



Other Considerations for Migration

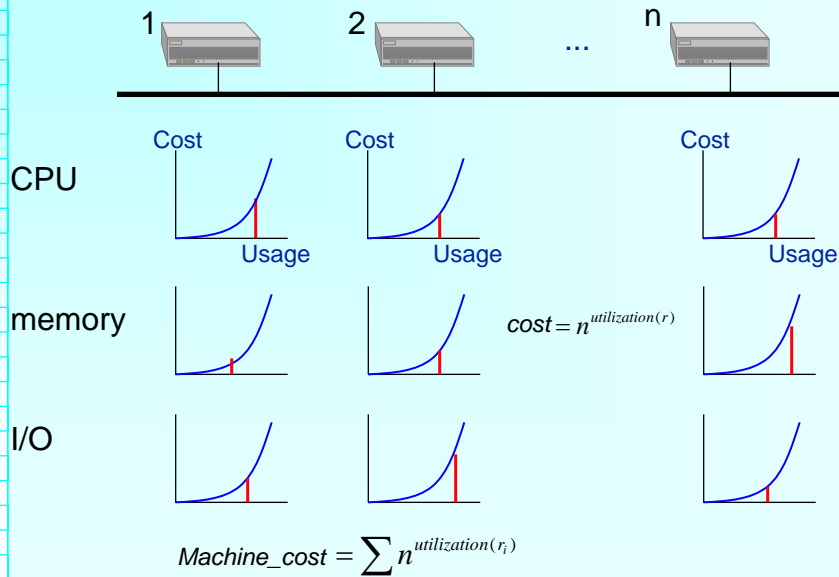
- Not only CPU load!!!
- Memory.
- I/O - where is the physical device?
- Communication - which processes communicate with which other processes ?

Resource Management of DOS

- A new *online* job assignment policy based on *economic principles*, *competitive analysis*.
- Guarantees near-optimal *global lower-bound* performance.
- Converts usage of *heterogeneous* resources (CPU, memory, IO) into a single, homogeneous cost using a *specific* cost function.
- Assigns/migrates a job to the machine on which it incurs the *lowest* cost.



DOS Resource Management



Yair Amir & Raluca Musaloiu-E

Spring 2008/ Lecture 11

27

Case Study – PVM and MOSIX

- PVM - Parallel Virtual Machine
 - Static Assignment of jobs to machines.
 - Default - round robin assignment policy.
 - Widely used.
- Mosix
 - Dynamic job migration.
 - Main objective is load balancing, with some ad-hoc heuristics for memory depletion.

Yair Amir & Raluca Musaloiu-E

Spring 2008/ Lecture 11

28

Enhanced PVM and Enhanced Mosix

- Enhanced PVM is similar to PVM. The decision where to statically place a new job is made according to the cost benefit framework.
- Enhanced Mosix is similar to Mosix. Job migration decisions are made according to the cost-benefit framework.

This framework takes into account:

- CPU load
- Memory utilization

Easy to add I/O.

Harder to add IPC.

Evaluation Methodology

- **Cluster:**

Machine Type	# of these Machines	Processing Speed	Installed Memory
Pentium Pro	3	200 MHz.	64 MB of RAM
Pentium	2	133 MHz.	32 MB of RAM
Laptop w/ Ethernet	1	90 MHz.	24 MB of RAM.

- **Simulation:**

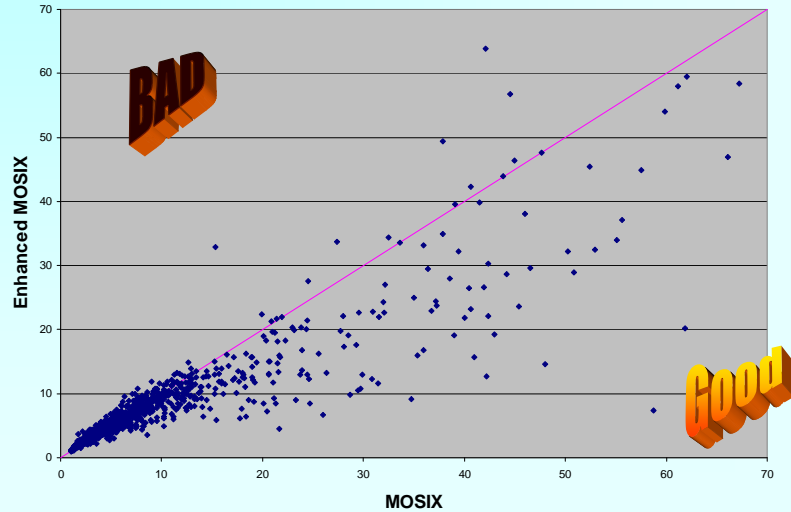
- 3,000 identical executions (scenarios) per policy, each represents 10,000 - 20,000 sec.

- **Validation (real life executions):**

- 50 identical executions per policy (similar to the simulations).

Simulation: Mosix / EMosix

Average Slowdown (lower is better)



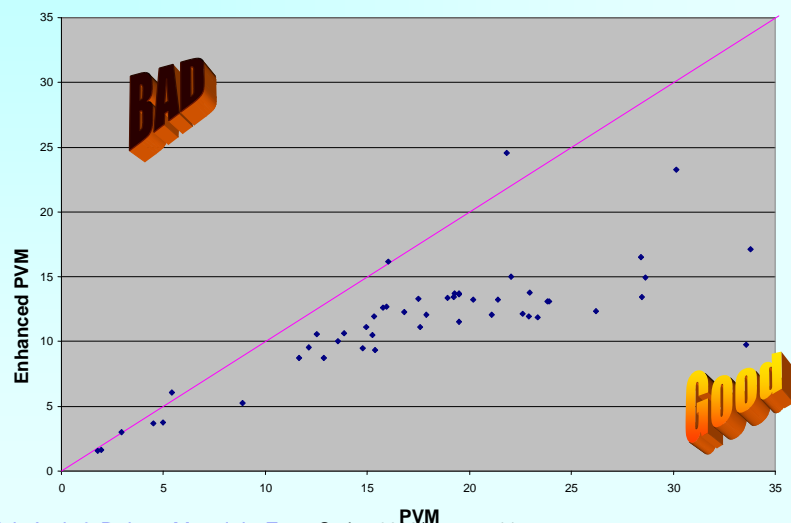
Yair Amir & Raluca Musaloiu-E

Spring 2008/ Lecture 11

31

Real Life: PVM / EPVM

Average Slowdown (lower is better)



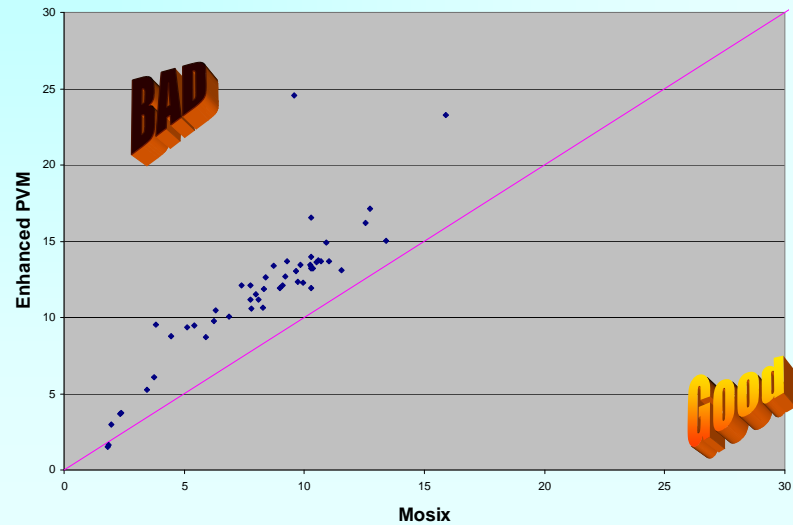
Yair Amir & Raluca Musaloiu-E

Spring 2008/ Lecture 11

32

Interesting: Mosix / EPVM

Average Slowdown (lower is better)



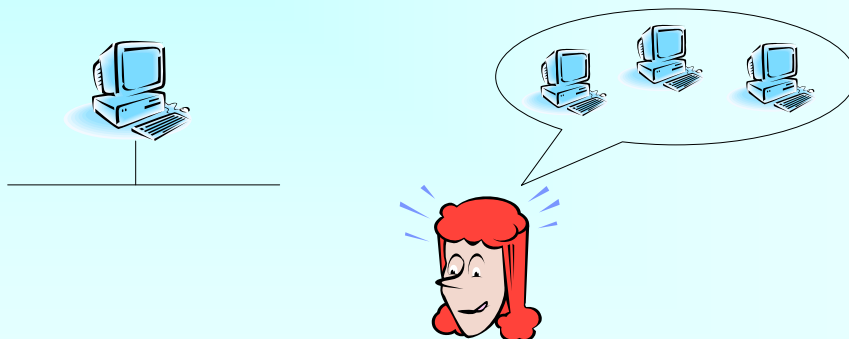
Yair Amir & Raluca Musaloiu-E

Mosix Spring 2008/ Lecture 11

33

Virtualization

- Creates the illusion of **multiple** virtual machines on a **single** physical machine.



Yair Amir & Raluca Musaloiu-E

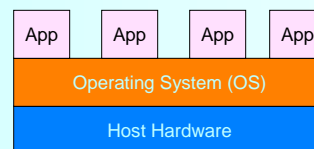
Spring 2008/ Lecture 11

34

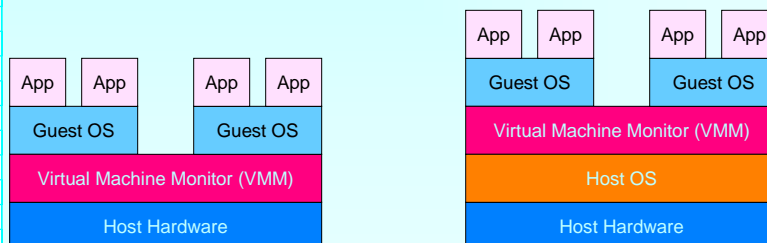
Why Virtualize?

- Security
- Reliability and availability
- Cost
- Load balancing
- Legacy applications

Regular machine



Virtual Machine



Virtualization through Emulation

- VM emulates complete hardware.
- Allows unmodified guest OS for a different PC to be run.

Bochs



Microsoft Virtual PC for Mac



QEMU
open source processor emulator

Full / Native Virtualization

- VM simulates enough hardware to allow an unmodified guest OS to be run in isolation.
- Host and guest architectures are the same.

Parallels

VMware
Workstation

VMware
Server

Mac-on-Linux



QEMU
open source processor emulator

Paravirtualization

- Virtual hardware architecture differs from the underlying physical architecture.
- VMM presents a software interface to the guest OS. Guest OS must be modified to use VMM interface.
- Does not require changes to the application.



OS Level Virtualization

- No guest OS. Host and VMs share the same kernel.
- No need to manage the OS.



(used by PlanetLab)

Application Level Virtualization

- Create isolated virtual environments to deploy and execute applications.
- Applications need to be packed for virtualization environment.



Virtualization in the Past

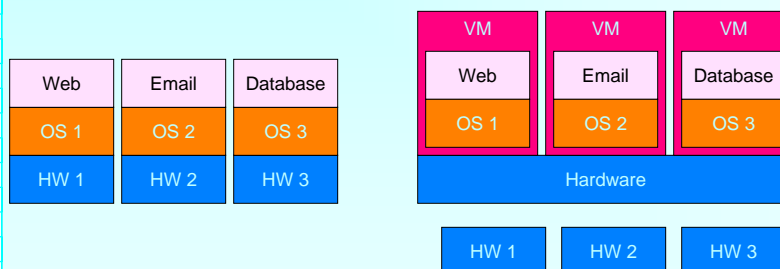
- Long time ago (**more than 30 years**): IBM VM/370.
- 80's: the cost of the hardware decreased, making the virtualization less interesting (moved away from mainframes).

Virtualization today: Server Consolidation

- Need for application isolation. Run one application per server.
- Lots of underutilized hardware (estimated resource utilization is 10-15%).

Server Consolidation

- Use virtualization to utilize physical resources as efficiently as possible.



Less hardware, space, power consumption.

Beyond Server Consolidation: Application Deployment

- Traditional approach: predict resource usage, acquire and release resources (servers acquisition can take months!).
- Use virtualization to provide a uniform [application deployment](#) environment for developers.
- Applications and OS are no longer tied to the physical hardware. Basic deployment unit is now the VM.
- VMs can be dynamically moved based on their resource requirements.

Even more: Utility Computing

- No need to own the hardware. Instead, pay only what you use!
- Computing and Storage is used as a [service](#).
- Example: Amazon EC2, Amazon S3
 - 10¢ for for 1 EC2 Compute Unit (1 hour equivalent CPU capacity of a 1.0-1.2 GHz 2007 Opteron or 2007 Xeon processor).

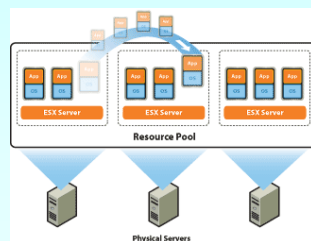
Challenges

- How to accurately characterize applications' resource requirements.
 - What and how much: CPU, memory, disk, network utilization...
 - Resource usage may vary over time.
- How to distribute the VMs over the physical resources.
 - Various strategies: keep VMs together, separate VMs.

Challenges (2)

- How to balance the server workloads at runtime (migrate VM to another physical server).

VMWare DRS



- Resource pools.
- Automate hardware maintenance.
- Manual or automatic mode.
- Power management.

Image from <http://www.vmware.com/products/vi/vc/drs.html>