# FUNDAMENTALS OF VIRTUALIZATION

**Virtualization**

☐ Introduction of what virtualization is.

☐ Understanding of where virtualization can be applied.

☐ Knowledge of virtualization could help for BG.

□ What virtualization is.

□ Where virtualization could be used.

□ Why virtualization is used.

□ Which kind of virtualization options could be used.

□ Main aspects to remember when deploying a

  virtualization system.

# Contents

- **What virtualization is.**

- Where virtualization could be used.

- Why virtualization is used.

- Which kind of virtualization options could be used.

- Main aspects to remember when deploying a virtualization system.

□ It is possible to search information from several sites:

■ Wikipedia:
http://en.wikipedia.org/wiki/Virtualization

■ Forums:
http://virt.kernelnewbies.org/

■ On-line courses:
http://www.govirtual.org/docs/DOC-1024

■ Etc.

☐ Main aspects that we are going to find:
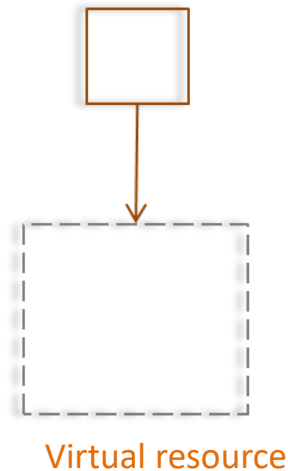
1. Virtualization term is not new:
   - It has been used since 60's

   *Vintage !* ☺

2. It has been applied to different aspects and areas of computing:
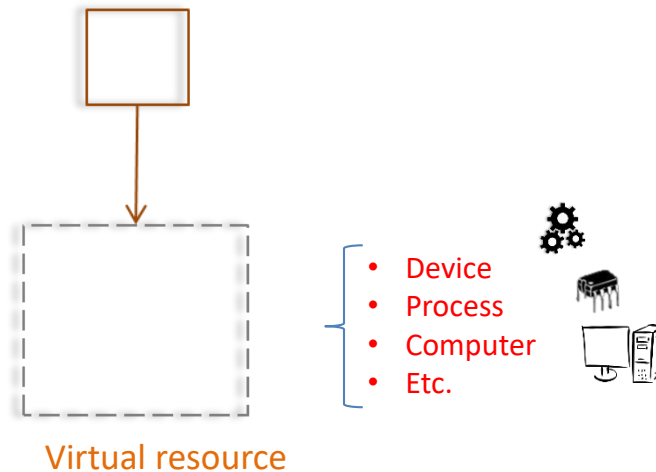   - Components, servers, personal computers, etc.

☐ **Virtualization** is a broad term that refers to the **abstraction of computer resources**.

☐ A technique for **hiding the physical characteristics** of computing **resources from the way** in which **other** systems, applications or end users **interact** with those resources.
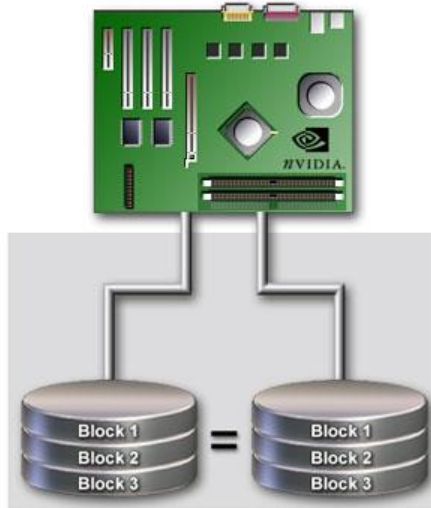
Virtual resource

□ What virtualization is.

□ **Where virtualization could be used.**

□ Why virtualization is used.

□ Which kind of virtualization options could be used.

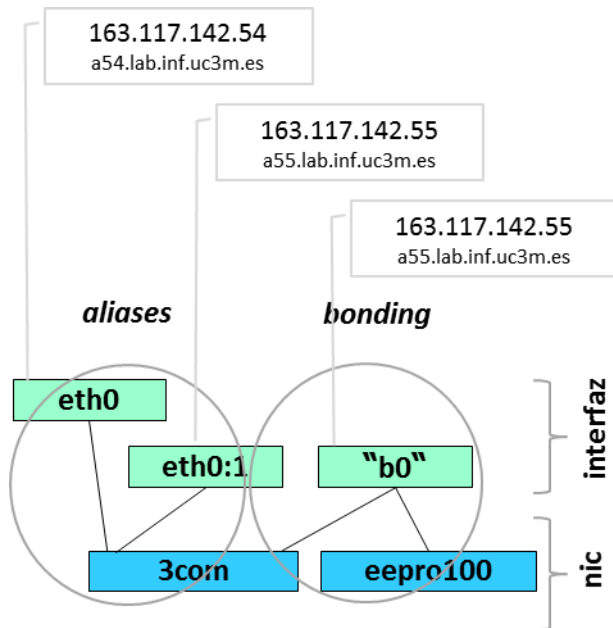□ Main aspects to remember when deploying a

virtualization system.

☐ **Virtualization** is a broad term that refers to the abstraction of computer resources.

☐ A technique for hiding the physical characteristics of computing resources from the way in which other systems, applications or end users interact with those resources.

Virtual resource

- Device
- Process
- Computer
- Etc.

- Storage device:
  - **E.g.: RAID**

- Networking:
  - E.g.: IP bonding, IP aliasing, NAT

- Processor:
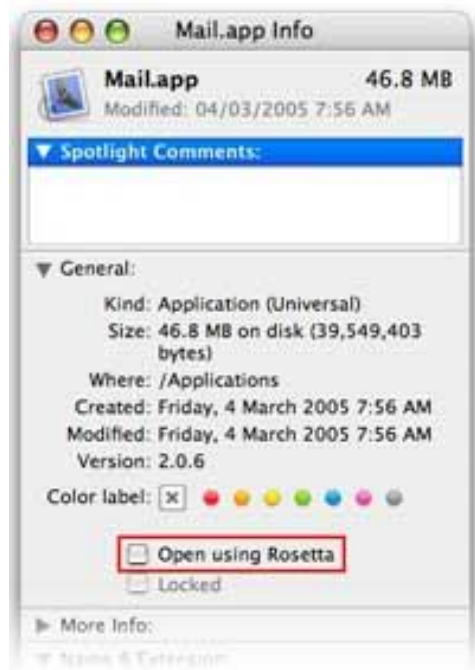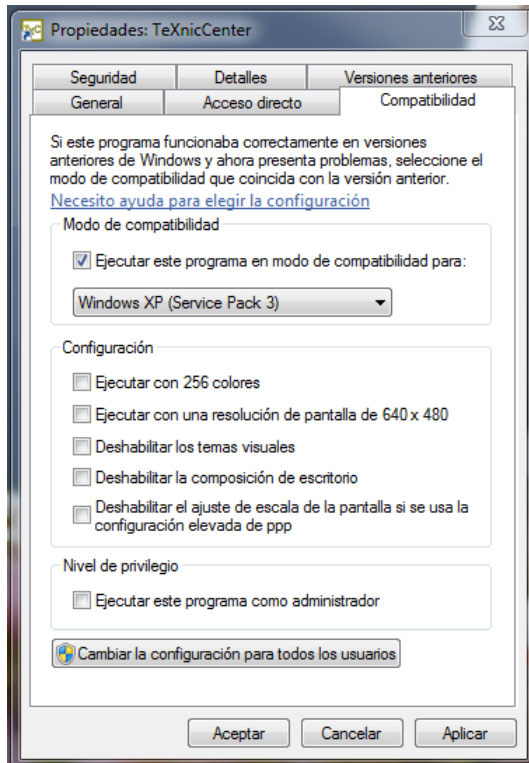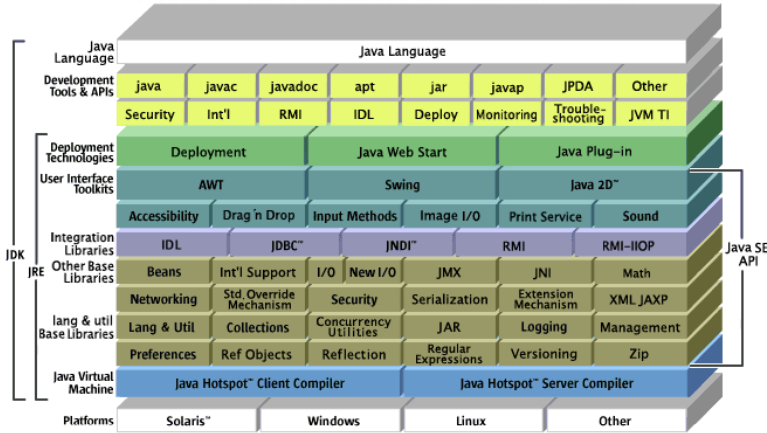  - E.g.: emulate a different instruction set

- Storage device:
  - E.g.: RAID

- Networking:
  - **E.g.: IP bonding, IP aliasing, NAT**

- Processor:
  - E.g.: emulate a different instruction set

Rosetta was used by Apple in the
PowerPC to Intel transition.

– Storage device:
  - E.g.: RAID

– Networking:
  - E.g.: IP bonding, IP aliasing, NAT

– Processor:
  - **E.g.: emulate a different instruction set**

– Emulation of the application private environment:

- **E.g.: Windows Vista/7 compatibility mode**

– Language level virtualization:
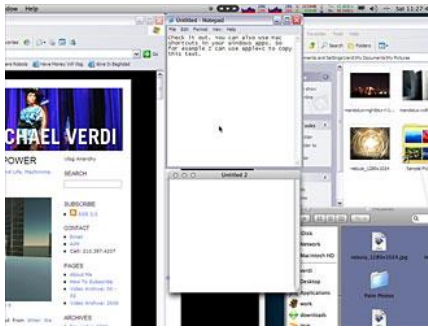
- E.g.: Java and .NET

Java™ Platform Standard Edition

- Emulation of the application private environment:
  - E.g.: Windows Vista/7 compatibility mode
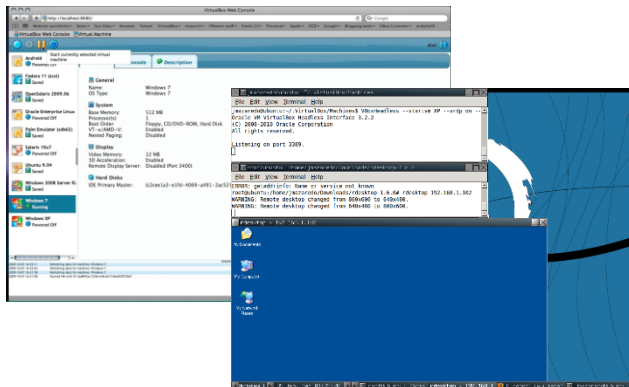
- Language level virtualization:
  - **E.g.: Java and .NET**

- Guest desktop as a window:
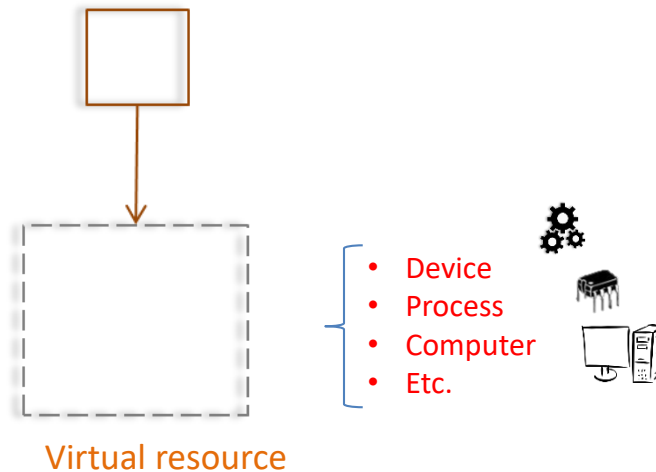  - E.g.: VMWare, VirtualBox, etc.



- A window for each guest application:
  - E.g.: Coherence mode, fluid view, etc.



- Remote desktop window:
  - E.g.: XEN, VMWare ESX, VirtualBox Headless, etc.

☐ **Virtualization** is a broad term that refers to the abstraction of computer resources.

☐ A technique for hiding the physical characteristics of computing resources from the way in which other systems, applications or end users interact with those resources.
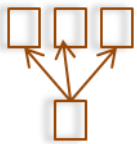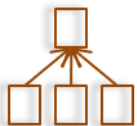
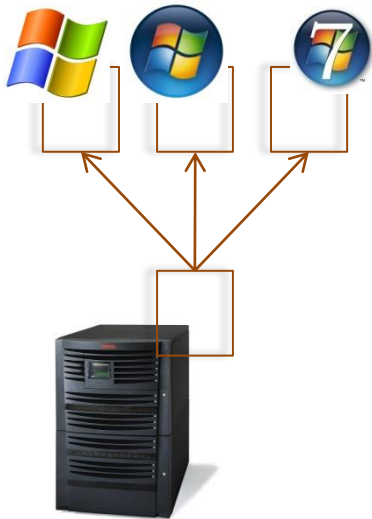Virtual resource

- Device
- Process
- Computer
- Etc.

□ It includes:

▫ To make a single physical resource (such as a server, an operating system, an application, etc.) be exposed as a different logical resource.
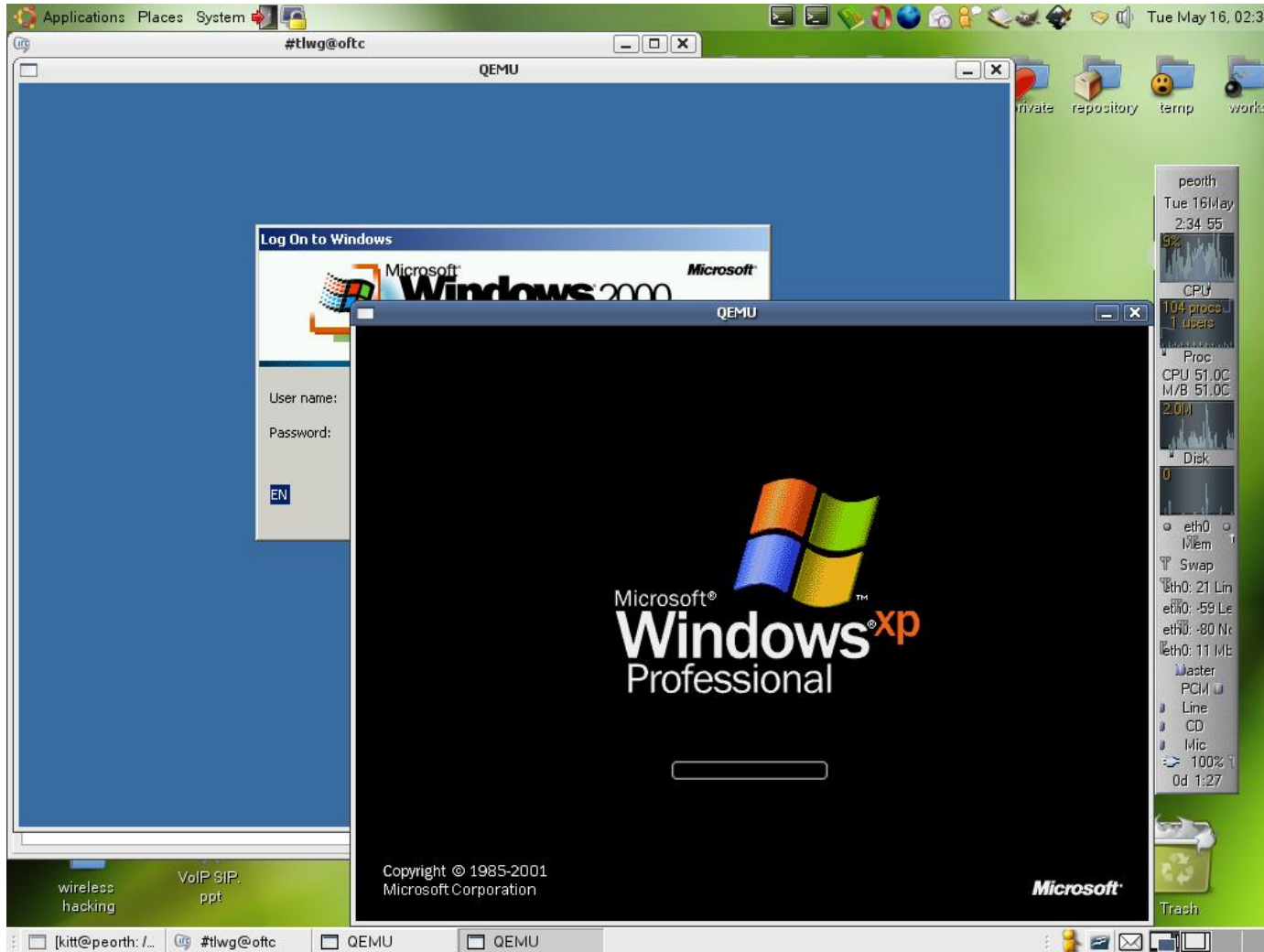
▫ To make a single physical resource (such as storage devices, servers, etc.) be exposed as multiple logical resources.

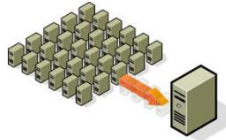▫ To make multiple physical resources (such as storage devices, servers, etc.) be exposed as a single logical resource.

- We will focus on the platform virtualization in terms of virtual machines.

- The real system will be named as host system, the virtualization system will be named guest.

http://www.kitty.in.th/files/376/qemu.png

- What virtualization is.

- Where virtualization could be used.

- **Why virtualization is used.**

- Which kind of virtualization options could be used.

- Main aspects to remember when deploying a virtualization system.

- ☐ Server consolidation
- ☐ Service isolation
- ☐ Disaster recovery
- ☐ Testing or training
- ☐ Application portability

□ Server consolidation:

◘ To reduce costs (by multiplexing resources).

◘ Simplifying the administration and management.

**1000 €**    **1000 €**    **1000 €**    **1000 €**

# Server consolidation:

- To reduce costs (by multiplexing resources).
- Simplifying the administration and management



**3000 €**

□ Improve security:

   ◘ Insolate services in different computers.

   ◘ Different security policy for each computer.

□ Improve security:

◘ Insolate services in different computers.

◘ Different security policy for each computer.

- Improve disaster recovery:
  - Hot-spare machine(s).
  - Automatic work re-routing while rebooting/fixing.

□ Improve disaster recovery:

  ◘ Hot-spare machine(s).
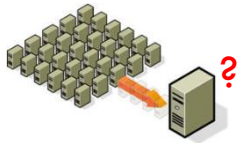
  ◘ Automatic work re-routing while rebooting/fixing.

- Improve disaster recovery:
  - Hot-spare machine(s).
  - Automatic work re-routing while rebooting/fixing.

□ Better testing environment:

■ It enables the execution in other work environment.
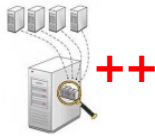
■ It improves the restoring process (easier/bit-faster).

□ Better testing environment:

  ◘ It enables the execution in other work environment.

  ◘ It improves the restoring process (easier/bit-faster).

Server Virtualization Usage

□ Complex dimensioning

□ More resources per node are needed

□ Double administration level

□ Some performance loss

- An appropriated sizing is required:
  - The (virtual) servers might change it requirements (memory, cpu, …)
  - An unappropriated sizing has impact in <u>all</u> (virtual) servers.

□ More resources per host are needed:

 ◘ 8 1GiB servers consolidated on 1 8GiB server

 ◘ Not always is easy (and cheaper) to buy one server with n network cards, n GB of RAM, n TB of disk, etc.

□ More resources per host are needed:

- ▪ 8 1GiB servers consolidated on 1 8GiB server

- ▪ Not always is easy (and cheaper) to buy one server with n network cards, n GB of RAM, n TB of disk, etc.

□ Double administration level:

  ◘ (Real) host computers

  ◘ (Virtual) guest computers

  ◘ Management of host/guest relationship

    ▪ If a host computer has problems, all guest computer has to be migrated.



ssh virtual

ssh real

- A "little" loss of performance:
  - In CPU could be low: between 3% and 12%
  - Graphic card and buses bandwidth?
  - Hard disk shared among several guest computers?

□ What virtualization is.

□ Where virtualization could be used.

□ Why virtualization is used.

□ **Which kind of virtualization options could be used.**

□ Main aspects to remember when deploying a

virtualization system.

☐ Good news: many options…

User Mode Linux

QEMU
open source processor emulator

IBM System z™

bochs 2.3

OpenVZ

vmware®

Xen™

KVM

□ **Bad** news: many options…


User Mode Linux


QEMU
open source processor emulator


IBM System z™


bochs 2.3


OpenVZ


vmware®


Xen™


KVM

□ **To know how internally works.**

  ◘ Dependencies, restrictions, etc.

□ **To know the important details about virtualization system architectures.**

  ◘ To group solutions by common characteristics.

- A new layer is added between the operating system and hardware
  - It will 'talk' with all kind of hardware
  - Arbitrate hardware resources across all operating systems

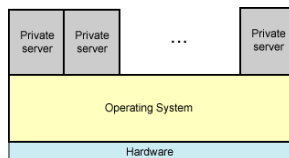- Each operating system is executed in privileged mode but the monitor/hypervisor intercepts its requests to server them.

1

□ A new layer is added between the operating system and hardware

3 □ It will 'talk' with all kind of hardware

4 □ Arbitrate hardware resources across all operating systems

□ Each operating system is executed in privileged mode but the monitor/hypervisor intercepts its requests to server them.

2

□ Hardware Emulation



□ Full Virtualization



□ Para-virtualization



□ Containers

□ A virtual machine on the host system is created to emulate the target hardware.

□ Advantage: you can execute software for CPU1 on CPU2 without modifications.

□ Disadvantage: s-l-o-w (about x100)

❑ **Bochs**





*Linux/W95*          *W95/WXP*

❑ **Qemu**

- Hardware is shared among all guest operating systems through a **hypervisor.**

- Advantage: the operating system do not need to be modified.

- Disadvantage: it is necessary to intercept the access of the operating system to the hardware:
  - Hardware support
  - On-the-fly binary patching

❏ VMware
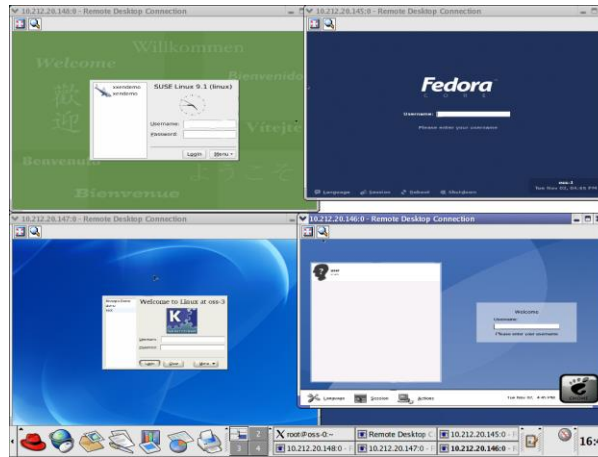




*WXP/MacOS*                    *WXP/Linux*

❏ z/VM



*z/Linux sobre z/VM*

□ Similar to the full-virtualization but the guest software collaborate with hypervisor.

□ Advantage: The operating system works with the hypervisor (less wasted time by interception mechanism).



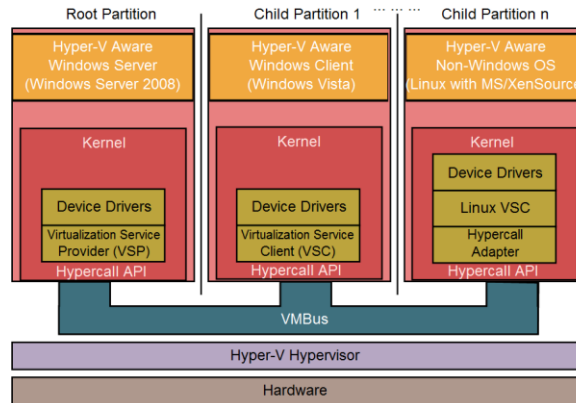□ Disadvantage: The operating system has to be modified in order to interact with the hypervisor.

❑ XEN



*Linux(S,F,D)/Linux(F)*

❑ Hiper-v



*Windows S2008+Vista+Suse / Windows*

❑ User Mode Linux (UML)



*Linux/Linux*

SIMULACIÓN DE UN CLUSTER USANDO *USER MODE LINUX*

AUTOR: VICTOR INIESTA SAMPAYO

2004

☐ Different approach: the operating system provides virtual copies of itself.

☐ Advantage: it is NOT possible to execute different operating systems.

☐ Disadvantage: best performance and more number of virtual machines executing (with less memory).



Private server | Private server | ... | Private server

Operating System

Hardware

# ❑ OpenVZ



**Density of OpenVZ in a
768 MiB (¾ Gb) RAM computer**

☐ **Hardware Emulation**



☐ **Full Virtualization**



☐ **Para-virtualization**



☐ **Containers**

http://www-128.ibm.com/developerworks/library/l-linuxvirt/index.html

□ A new layer is added between the operating system and hardware

3 □ It will 'talk' with all kind of hardware

4 □ Arbitrate hardware resources across all operating systems

□ Each operating system is executed in privileged mode but the monitor/hypervisor intercepts its requests to server them.

2

- Regular operating system has been designed to be executed in hardware in privileged mode.
  - In x86 processors, on ring 0

- But now the privileged code has to be executed without been privileged anymore (the hypervisor is now privileged)
  - Binary patching
  - New virtualization instructions

http://pdc-amd01.poly.edu/~wein/cs6243/ppts/CPUVirtualization.pptx

- ☐ <span style="color:purple">Binary translation/patching</span>:
  - ▣ Patching the instructions on the fly.
  - ▣ The guest code is analyzed and the privileged instructions are replaced with hypervisor calls.
  - ▣ Speed-up by caching the patched fragments.
  - ▣ <span style="color:green">Advantage</span>: Can be used on any kind of CPU.
  - ▣ <span style="color:red">Disadvantage</span>: S-l-o-w.

- □ Special hardware Instructions :
  - ▣ Ring '-1' where the hypervisor is executed.
    - ■ It reduce the performance penalty of dynamic on-the-fly translation.
  - ▣ Intel and AMD have developed instructions set extensions for virtualization. There are similar but no compatible.
  - ▣ Advantage: Fast request to the hypervisor.
  - ▣ Disadvantage: It needs special CPU support.

(intel)   Intel has:
  - VT-x as extensions IVT for IA-32 (Vanderpool)
  - VT-i as extensions IVT for IA-64 (Silvervale)
  - VT-d in 32/64 for Directed I/O

AMD   AMD has:
  - AMD-V (Pacifica) for 32/64
  - IOMMU as Directed I/O or PCI-Passthrough

http://en.wikipedia.org/wiki/Virtualization_Technology

Apps | Apps | Apps | …
Guest OS | Guest OS | Guest OS | …
Monitor/Hypervisor ⭐
Hardware
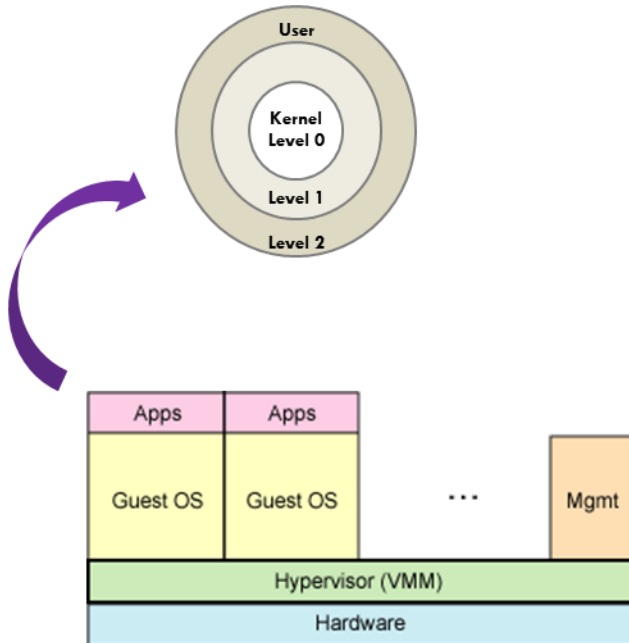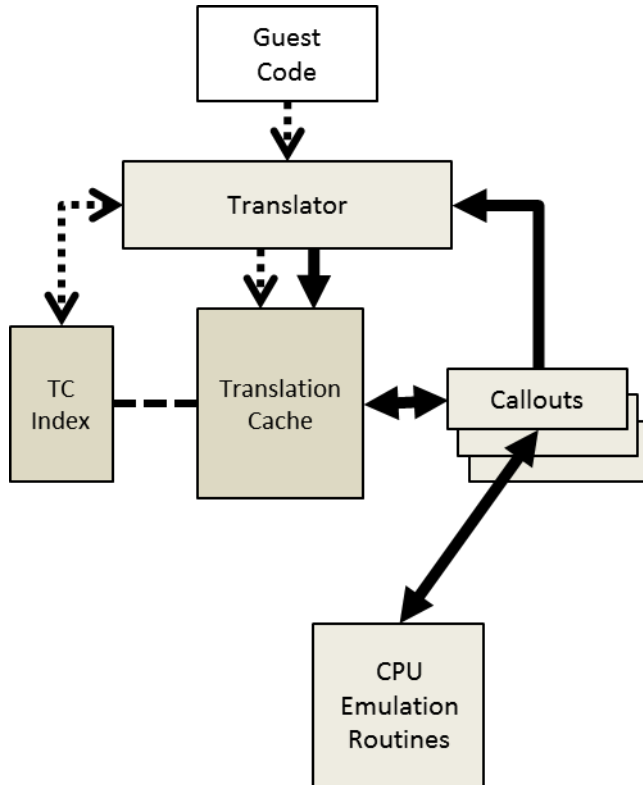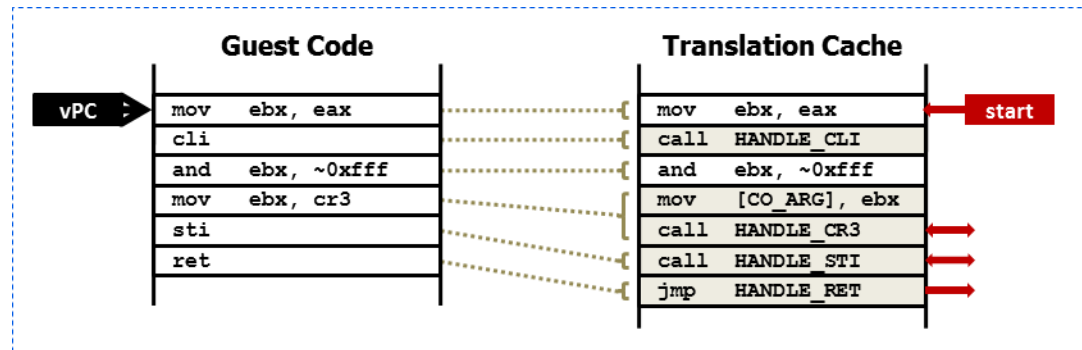
□ A new layer is added between the operating system and hardware

3 ▫ It will 'talk' with all kind of hardware

4 ▫ Arbitrate hardware resources across all operating systems

□ Each operating system is executed in privileged mode but the monitor/hypervisor intercepts its requests to server them.
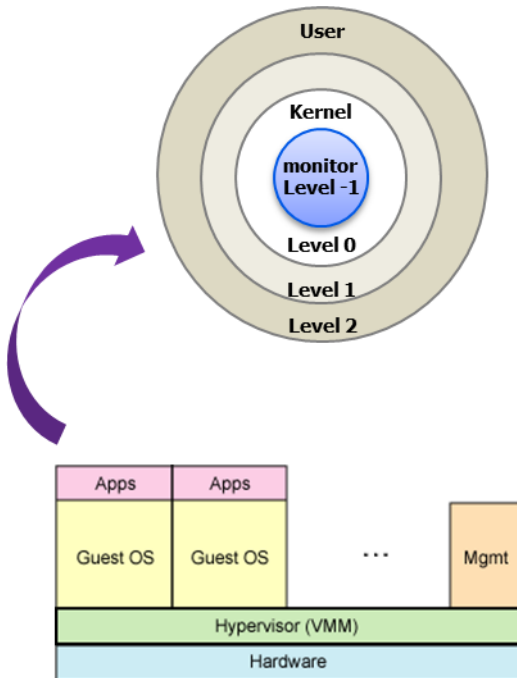
□ The monitor/hypervisor have to be able to work with all kind of hardware

- ◘ It has to have drivers for all hardware.
- ◘ It is very difficult to get drivers for all existing (and new) hardware.

□ But we can use a modified operating system as monitor/hypervisor:

- ◘ Hosted or Split
- ◘ "Pure" Hypervisor

□ Hypervisor

▫ To remove from an existent operating system everything but what is needed to transform it into a hypervisor.

▫ The hypervisor boots first, and then every virtual machine that uses it:

■ A: less interferences between guests

■ D: no so easier to install

| User-space (applications) | User-space (applications) |
|---|---|
| Guest OS (Virtual machine) | Guest OS (Virtual machine) |
| Hypervisor (Virtual machine monitor) ||
| Hardware ||

- XEN is now included in the Linux kernel (4.6 in progress)
- XEN could be described as a Linux system to which all has been remove but the base to be used as hypervisor.
  - Initially designed as para-virtualization system.

□ Hosted or Split

▫ Transform an existent operating system into a hypervisor.

▫ A V.M. is a process in the host system:

- D: Double scheduling
- D: Expensive access to the hardware
- A: Easy to install (like a familiar application)

- KVM is include in the Linux kernel since version 2.6.20
- KVM transforms the Linux kernel into an hypervisor as a module
  - Other guest operating system can be executed in user-space.
  - It use a modified QEMU process.

- A new layer is added between the operating system and hardware
  - It will 'talk' with all kind of hardware
  4 - Arbitrate hardware resources across all operating systems

- Each operating system is executed in privileged mode but the monitor/hypervisor intercepts its requests to server them.

- The monitor/hypervisor must be able to deal with all types of hardware:
  - It must have driver for all devices.
  - It provides access to the underlying hardware.

- Expose the hardware to the guest operating system:
  - Hypervisor device emulation.
  - User-space device emulation.
  - Gateway to device
  - SR-IOV and MR-IOV

http://www.ibm.com/developerworks/linux/library/l-pci-passthrough/

□ Hypervisor device emulation.

◘ E.g.: VMware workstation

◘ Advantage: easy to migrate

□ **User-space device emulation.**

▫ E.g.: KVM

▫ Advantage: easy to migrate (even to other hypervisor) and safe (no privileged)

□ **Gateway to device.**

  ▫ E.g.: VMware, XEN, etc.

  ▫ **Advantage**: efficient

□ **SR-IOV and MR-IOV**.

  ◘ Single-Root I/O Virtualization (one server)

  ◘ Multi-Root I/O Virtualization (blades)

❑ **Hardware emulation**

❑ **Full Virtualization**

❑ **Para-Virtualization**

❑ **O.S. Virtualization level**

1



❑ **Binary translation (patching)**

❑ **New instructions**

2



❑ **Hosted or split**

❑ **Hypervisor**

3



❑ **Hypervisor device emulation**

❑ **User-space device emulation**

❑ **Gateway to device**

❑ **SR-IOV and MR-IOV**

4

☐ What virtualization is.
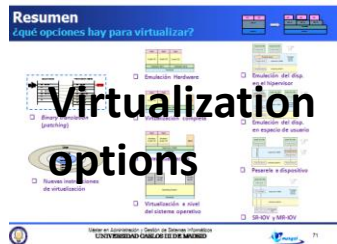
☐ Where virtualization could be used.

☐ Why virtualization is used.

☐ Which kind of virtualization options could be used.

☐ **Main aspects to remember when deploying a virtualization system.**

Virtualization options

**To understand** →

Criteria:
- Features
- Performance
- Information

**To apply** →

Requirements

OpenVZ

**Xen** **vmware** **QEMU** open source processor em

**IBM System z™**

**KVM**

Requirements

**Virtualization options**

To understand

Criteria:
* Features
* Performance
* Information

To apply

| | full virt | paravirt | containers (OS virt) | license | architectures | performance | SMP guests | CPU / memory hotplug | standalone host | notes |
|---|---|---|---|---|---|---|---|---|---|---|
| XEN | ✔ | ✔ | | GPL | i686, x86-64, IA64, PPC | paravirt very fast, full virt medium | | | | full virt needs VT / AMD-V |
| KVM | ✔ | ✔ | | GPL | i686, x86-64 | paravirt very fast, full virt medium | | | | full and para virt need VT / AMD-V |
| lguest | | ✔ | | GPL | i686 | slow/medium | | | | |
| rhype | | ✔ | | GPL | i686, x86-64, PPC | fast | (?) | | | research project |
| MoL | ✔ | | | GPL | PPC | fast | | | | 32 bit only |
| UML | | ✔ | | GPL | i686, x86-64, PPC | slow | | | | upstream |
| L4Linux | | ✔ | | GPL | i686, ARM | medium | | | | |
| qemu | ✔ | | | GPL | i686, x86-64, IA64, PPC | slow/medium, fast with kQEMU | | | | |
| OpenVZ | | | ✔ | GPL | i686, x86-64, IA64, PPC, SPARC | native | | | | live migration |
| Linux-VServer | ✔ | | ✔ | GPL | i686, x86-64, IA64, PPC | native | | | | poor performance isolation |
| VMware | ✔ | | | proprietary | i686, x86-64 | medium | | | | |
| LPAR | | ✔ | | proprietary | s390 | native | | | | |
| z/VM | ✔ | ✔ | | proprietary | s390 | very fast | | | | typically runs under LPAR |

http://virt.kernelnewbies.org/TechComparison

Hypervisor comparison

Performance loss

- Support

- Documentation

- Forums

- Recent deployments

OpenVZ

Requirements

Virtualization options

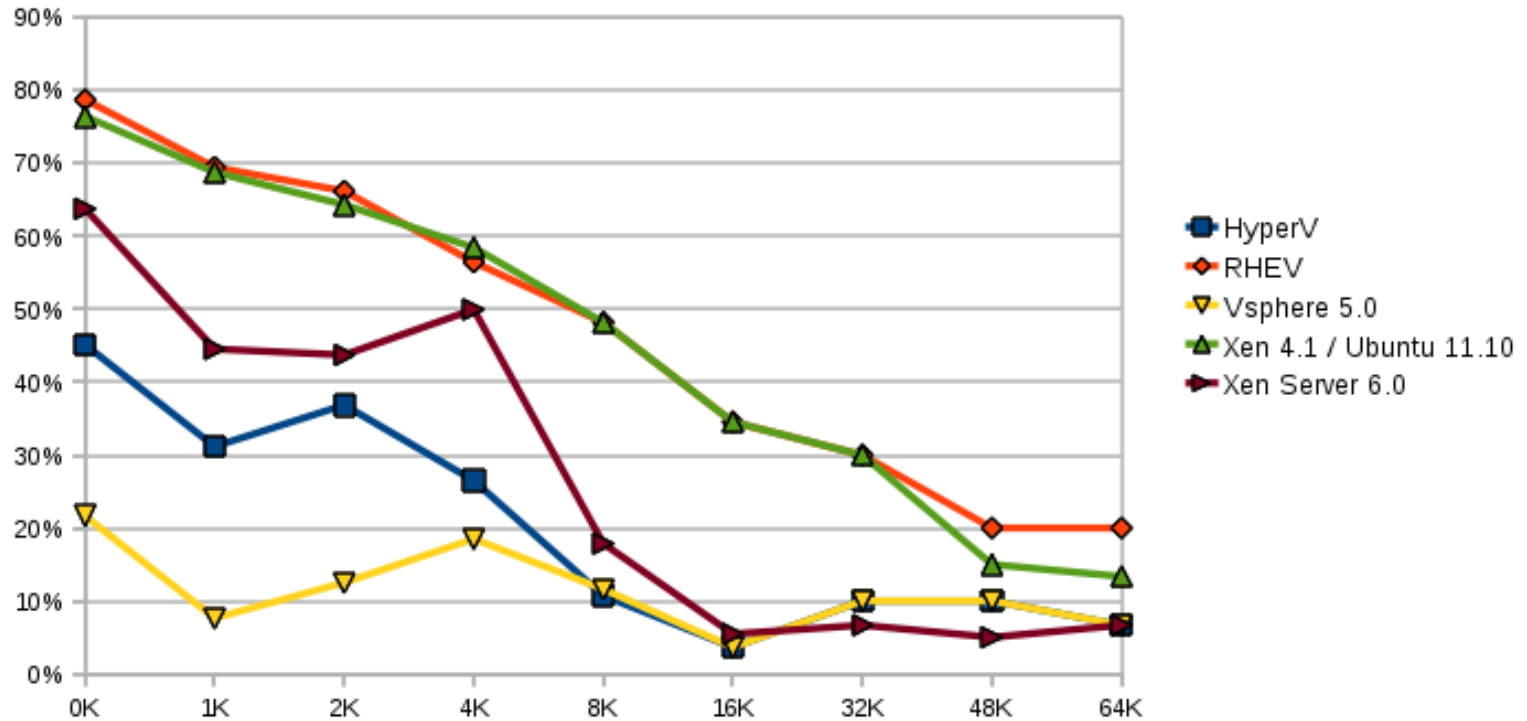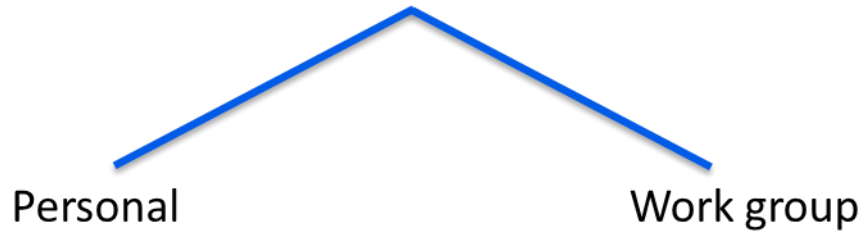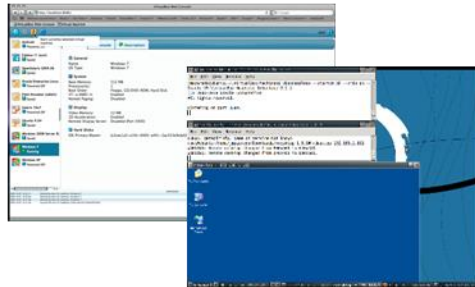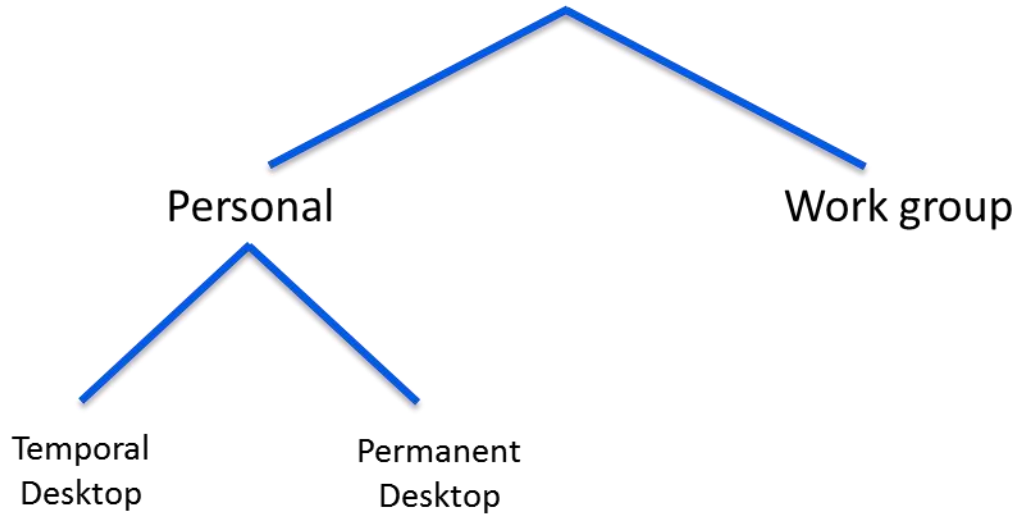**To understand** →
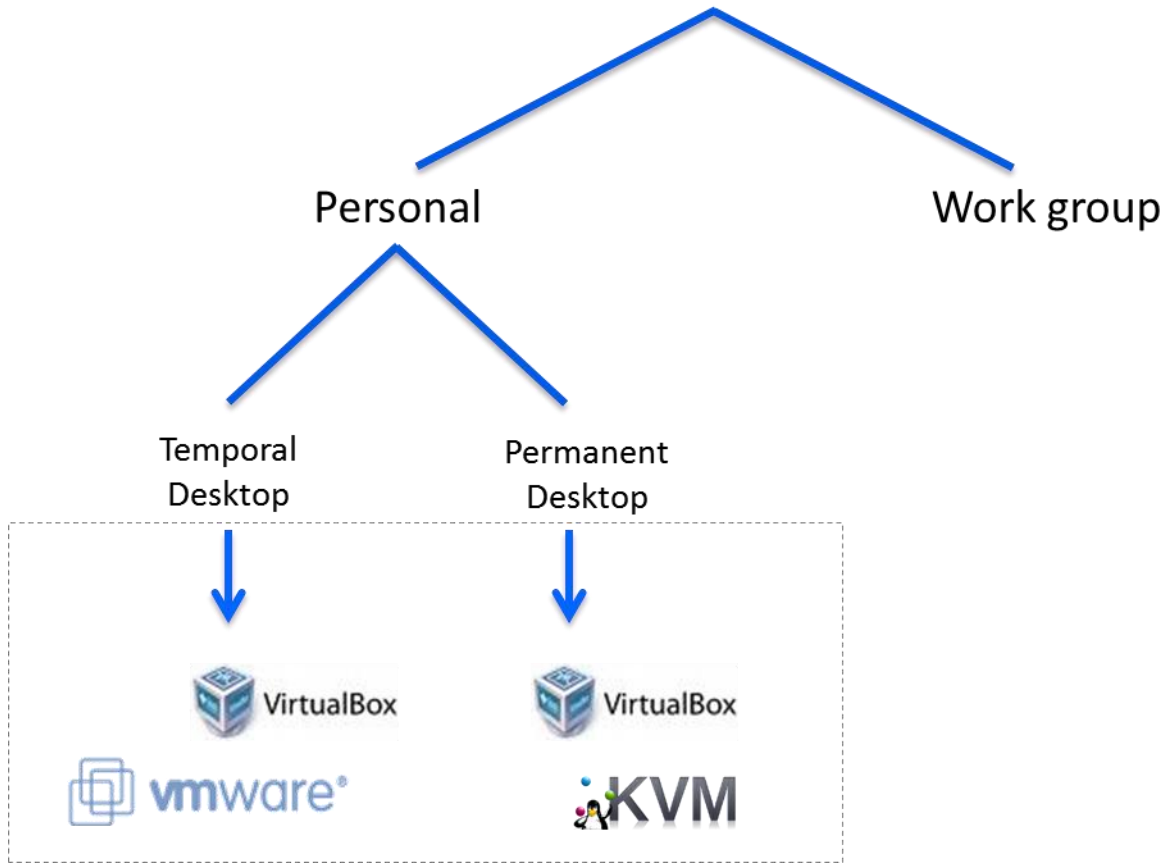
Criteria:
- Features
- Performance
- Information

**To apply** →

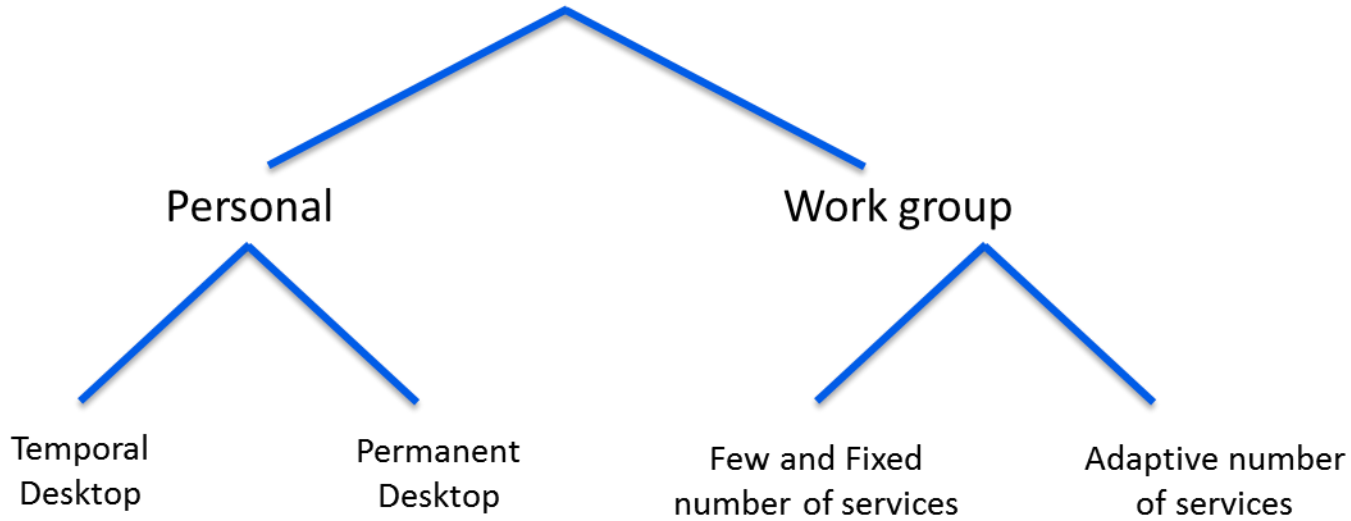Personal          Work group

```
                    Personal                           Work group
                   /        \                         /          \
            Temporal      Permanent          Few and Fixed      Adaptive number
            Desktop       Desktop            number of services  of services
```

```
                                    /\
                                   /  \
                                  /    \
                           Personal    Work group
                            /\              /\
                           /  \            /  \
                 Temporal      Permanent  Few and Fixed    Adaptive number
                 Desktop       Desktop    number of services  of services
                                                               /\
                                                              /  \
                                                    Own infrastructure   Rented infrastructure
```

# FUNDAMENTALS OF VIRTUALIZATION ON BIG DATA SYSTEMS

**Lesson 3**

**Virtualization**