# Database Performance in a Virtualized World

NoCOUG Winter Conference 2012

Eric Jenkinson

# Agenda

- Organizational Challenges
- Types of virtualization
- CPU Scheduling and Resource Allocation
- What you as a DBA need to do to thrive in a virtualized environment

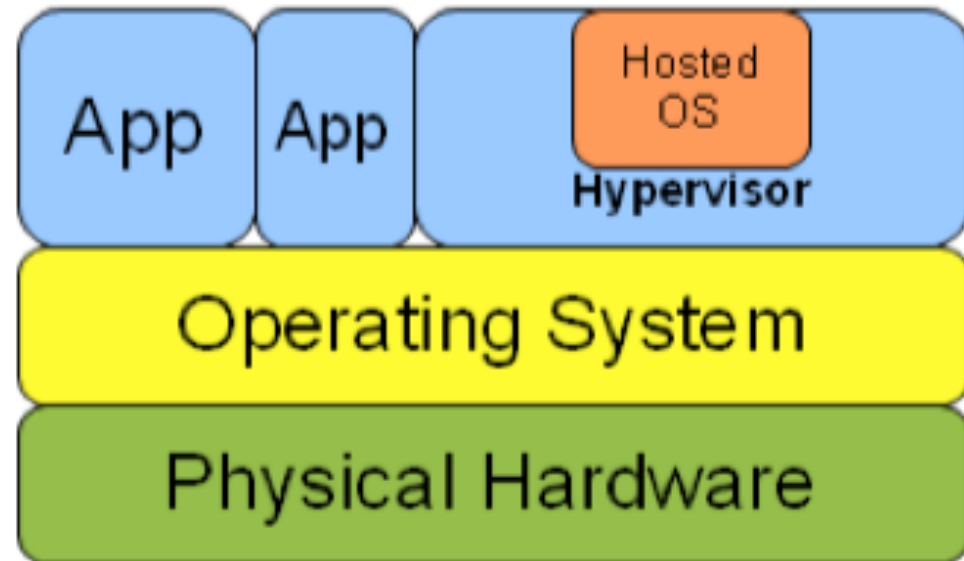# Organizational Issues in Virtualized Environments

- Knowledge of virtualization technology is not well understood outside of the server group
- Limited visibility into the virtualization technology stack
- "Throwing Hardware at the Problem" is easier.
- VM Sprawl
- DBAs, especially Oracle DBAs, are less likely to adopt virtualization

# Hypervisor

- Provides an abstraction of the physical hardware to the guests
- Manages the execution of the guest OS
- Manages the physical hardware resources
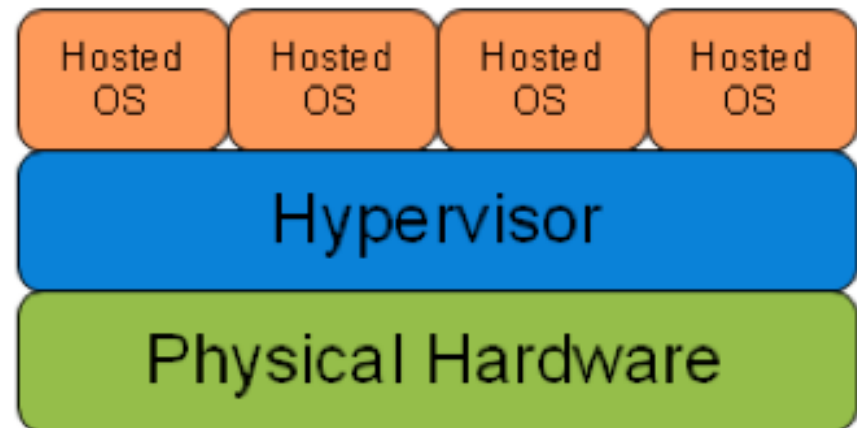- Two hypervisor types
  - Type 2
  - Type 1

# Type 2 Hypervisor

- Hosted
- Hypervisor runs as an application
- Hypervisor does not have direct control of hardware
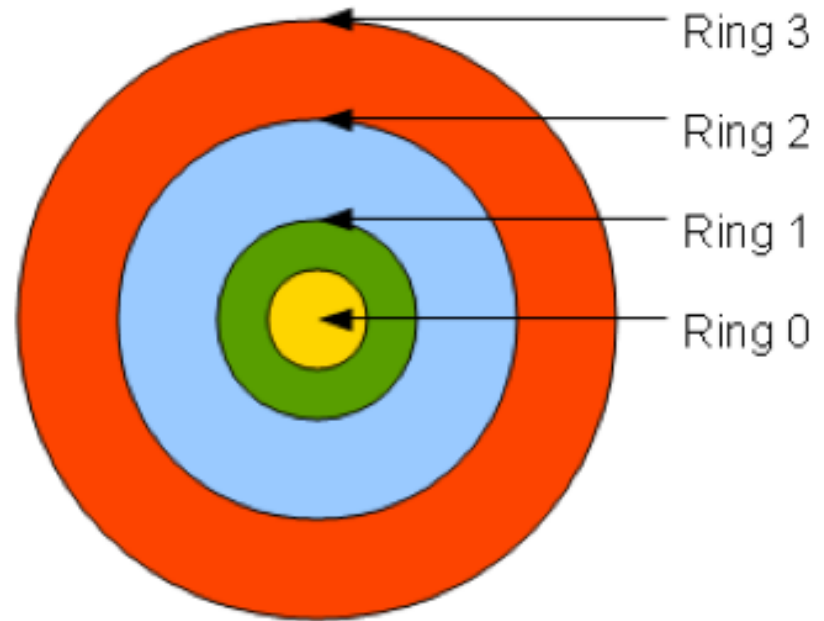- Examples
  - VMWare Server
  - Oracle VirtualBox

# Type1 Hypervisor

- Bare metal
- Has full control of hardware
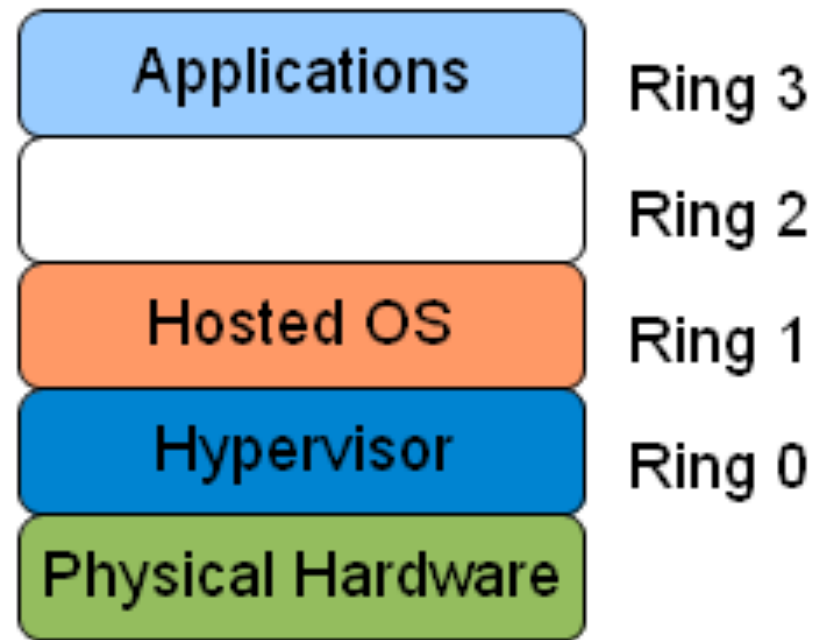- Examples
  - VMWare ESX, ESXI
  - Oracle VM

# CPU Rings

- Ordered from most privileged (ring 0) to least privileged (ring 3)

- OS and device drivers operate in ring 0
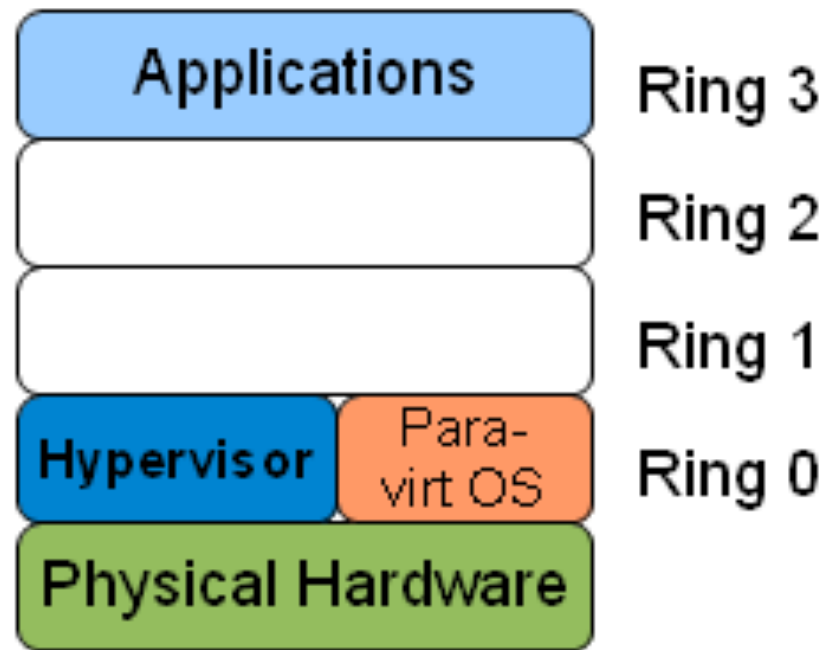
- Applications run in ring 3

# Full Virtualization

- Guest OS is unaware of virtualization

- Hypervisor traps privileged OS calls and reprocesses them

- Guest OS kernel is not modified

- Support
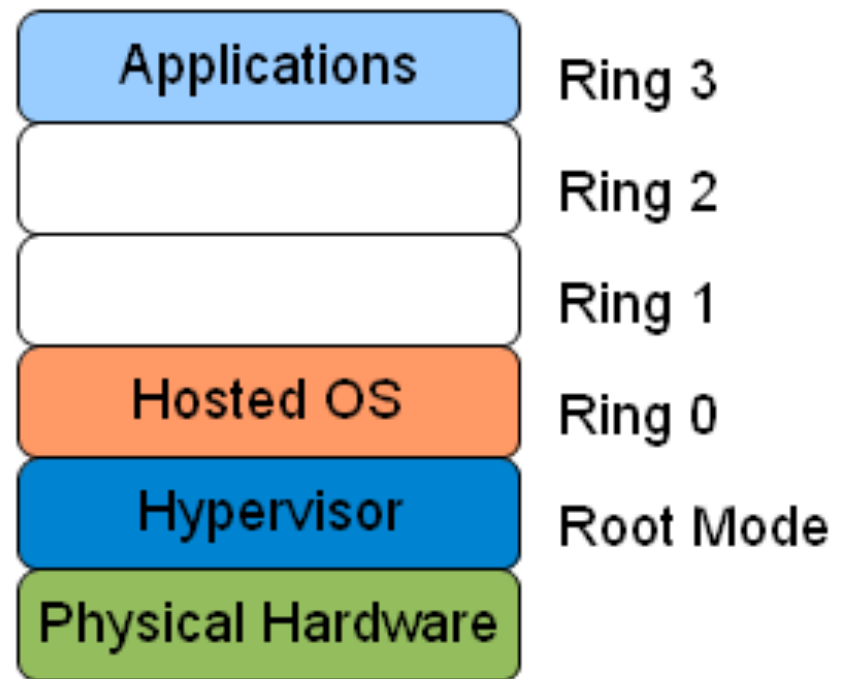  - VMWare on old hardware

| | |
|---|---|
| Applications | Ring 3 |
| | Ring 2 |
| Hosted OS | Ring 1 |
| Hypervisor | Ring 0 |
| Physical Hardware | |

# Paravirtualization

- Guest OS is aware of virtualization

- Guest OS Kernel modified to make hyper-calls instead of privileged calls

- Paravirtualization support in Linux Kernel 2.6.23 and higher

| | |
|---|---|
| Applications | Ring 3 |
| | Ring 2 |
| | Ring 1 |
| Hypervisor / Para-virt OS | Ring 0 |
| Physical Hardware | |

# Hardware Assisted Virtualization

- CPU support for Virtualization
- Root Mode ring below ring 0
- Privileged calls trapped and sent hypervisor
- Guest OS does not need to be modified
- Support
  - VMWare and Oracle VM

| | |
|---|---|
| Applications | Ring 3 |
| | Ring 2 |
| | Ring 1 |
| Hosted OS | Ring 0 |
| Hypervisor | Root Mode |
| Physical Hardware | |

# Hardware Assisted Virtualization with Paravirtualized Drivers

- Hybrid approach
- Guest OS does not have to be modified
- CPU virtualization support is required
- Paravirtualized drivers are required
- Support
  - VMWare and Oracle VM

# VMWare ESX Architecture

- VMKernel
  - CPU scheduler
  - Memory
  - Device Drivers
- Virtualization types
  - Binary Translation
  - Paravirtualized Drivers
  - Hardware Assisted Virtualization

# Xen / Oracle VM Architecture

- Three components
  - Xen Hypervisor
  - Domain 0 (dom0)
  - Domain U (domU)
- Two domain types
  - Privileged
    - dom0
  - Unprivileged
    - domU (Guest VMs)

# Resource Over Commitment

- Number of vCPUs can exceed the number of physical processors
- Sum of memory allocated to VM can exceed the amount of memory of the host
  - Oracle VM only through Max Memory

# ESX CPU Relaxed Co-Scheduling

- A vCPU can be in one of three states
  - Waiting for a CPU to become available
  - Has a CPU and executing
  - Has a CPU and idle
- Relaxed CPU co-scheduling
  - Per vCPU
  - vCPUs that advance too much are individually stopped

# ESX CPU Relaxed Co-Scheduling

- A vCPU is making progress if running or idle at the guest level
- Progress of each vCPU is tracked individually
  - Skew is measures as the difference between the slowest vCPU and other vCPUs
  - Skew does not grow if the vCPUs make equal progress during the co-scheduling period
- Skew enforcement
  - vCPUs that advance too much are stopped once the skew is reduced the stopped vCPUs may start individually

# Oracle VM – Xen Credit Scheduler

- Proportional fair CPU scheduler
- Each domain is assigned a weight and cap
  - Weight: a domain with 256 received twice as much CPU as a domain with 128
  - Cap maximum amount of CPU a domain can consume even if there are idle CPUs
- Automatically load balances vCPU across all available CPUs on SMP host

# Oracle VM – Xen Credit Scheduler

- Each physical CPU manages a run queue of vCPUs sorted by priority
  - Priority UNDER and OVER
- When inserting a vCPU to a queue it is put after all vCPU of equal priority
- When vCPU runs it consume credits. Until all credits are consumed its priority is UNDER
- Fair CPU scheduling, I/O can be skewed

# Memory Management

- Transparent Page Sharing
- Memory Ballooning
  - Requires Guest Additions to be installed
- Memory Compression
  - Compress memory pages that need to be swapped to disk
- ESX Swapping – Demand Paging
- Oracle VM only has memory ballooning at this time
  - Page sharing and demand paging is in Xen unstable

# Distributed Resource Scheduling (DRS)

- Both VMWare and Oracle VM have the ability to move VMs across physical servers

- The goal is to provide consistent resources to running VMs

- Moves VMs from heavily loaded servers to servers with a lighter load

- With out rules or affinity groups in place, DRS can be a source of "random" performance issues

*So what am I as as a DBA supposed to do with this?*

# My Experience

- Most performance problems fall under these areas
  - Poor knowledge of the virtualization stack
  - Poor or no VM placement policies
  - Poor or no resource prioritization
  - Little to no visibility into the virtualization stack
- The rest are the same problems that can exist in a purely physical environment

# Visibility is a Must

- VMWare: Request a read only account in vSphere
- Oracle VM: Oracle Enterprise Manager
  - Can be a problem with Oracle VM 2
- VMWare: esxtop OVM: top, vmstat, sar with paravirtualized kernel
- Third Party tools
  - Quest Spotlight for Oracle
  - Confio Ignite for VM

# Recognize Default Settings

# Recognize Default Settings

# Recognize Default Settings

- VMWare (CPU and Memory)
  - Reservation: minimum amount allocated/available at VM power on
  - Limit: maximum amount of the resource
  - Shares: Priority in acquiring the resource
- Oracle VM (CPU only)
  - Priority: The higher the priority, the more physical CPU cycles given to the VM
  - Processor Cap: The maximum amount of CPU a VM can consume

# The Problem with Defaults/Unlimited

# The Problem with Defaults/Unlimited

# The Problem with Defaults/Unlimited

# Virtual CPU Recommendations

- Set values for Limit/Shares or Priority/Cap to match business value

- Use only the vCPUs required and no more

- Monitor stolen time (OVM) and ESX Ready time to ascertain competition between VMs

- Watch out for CPU over commitment with VM that have many vCPUs

# Virtual Memory Recommendations

- VMWare
  - Use Reservation to avoid ballooning and swapping
    - SGA + PGA + processes overhead
  - Ensure VMWare Tools are installed (and up to date) to provide ballooning
- Oracle VM
  - Set Memory = SGA + PGA + process overhead

# Storage I/O

- Avoid sparse or dynamic growth virtual disks
- Follow Oracle and Storage vendor's best practices for Oracle Databases
- Use Storage IO Control (VMWare) to prioritize VM access to datastore
- Use dedicated datastores (VMWare) to avoid sharing disk workloads

# Network I/O

- Avoid having multiple high storage I/O VMs on the same physical host

- Insure paravirtualized drivers are installed

- Host server should have 1Gb min 10Gb recommended network adapter

# When Requesting a VM

- Request paravirtulaized drivers / VMWare tools to be installed
- Plan CPU and memory requirements and priority
  - Avoid "Cookie Cutter VMs"
- Know the business importance of this database
- Find where the VM is going to be placed and who its neighbors are

# Questions/Answers

# Thank you!