

# Basistechnologie: Virtualisierung

Ein Vortrag zum Cloudseminar im WS09/10

Maximilian Hoecker

Fakultät für Informatik,  
Hochschule Mannheim,  
Paul-Wittsack-Straße 10,  
68163 Mannheim

`maximilian.hoecker@stud.hs-mannheim.de`

20.10.2009





# Eine wage Begriffsdefinition

## Definition von Virtualisierung

Virtualization is a framework or methodology of dividing the resources of a computer into multiple execution environments, by applying one or more concepts or technologies such as hardware and software partitioning, time-sharing, partial or complete machine simulation, emulation, quality of service, and many others.



# Aber

## Die Definition sollte man nicht zu genau nehmen

- Virtualisierung bedeutet nicht immer **Aufteilen** von Ressourcen
- Man kann z.B. X Hardwarefestplatten zu einer logischen Partition **zusammenfassen** durch einen Virtualizationlayer
- Grid Computing virtualisiert durch das Nutzen verschiedener verteilter Ressourcen

# Agenda

- 1 Begriffsabgrenzung Virtualisierung
  - Definition
- 2 Virtualisierungsarten
  - **Hardwarevirtualisierung**
  - Softwarevirtualisierung
    - Anwendungsvirtualisierung
    - Virtualisierung von Peripherie-Geräten
    - Virtualisierung auf Betriebssystemebene
- 3 Virtualisierung von Betriebssystemen/Rechnern
  - Allgemeines
  - User Mode Linux
  - Emulation auf Applikationsebene
  - Para-Virtualized-Machines @ Xen
  - Hypervisor-Virtualized-Machines @ Xen
  - Vergleich der 3 Virtualisierungsarten (Xen, VMware, UML)





# Die erste Abstraktion

- Erster Großrechner *Atlas* (1962) mit Stracheys Time Sharing
- Neuer Architekturansatz im Betriebssystem
- Dort wurde separiert zwischen:
  - Prozessen des Betriebssystems (**Supervisor**)
  - Darunter **eine Komponente** die zuständig war für ausgeführte Programme
- Ebenfalls neu waren *Virтуelle Speicher/One Level Store* und *Paging*



# Server Partitioning

- Bis heute wurde die Architektur der S/360 beibehalten und erweitert
- Aktuelle Systeme: IBM z10 Business/Enterprise 1-64 CPUs  
64GB-1,5TB RAM
- Diese Systeme haben 2 verschiedene Modi mit denen gearbeitet werden kann
  - LPAR Mode (**L**ogical **P**artition Mode)
  - VM Mode (**V**irtual **M**achine Mode)
- Diese Modi können auch gemischt werden, also auf einer Partition mehrere VMs



# Zusammenfassung Hardwarevirtualisierung

## Vorteile von Hardware-Virtualisierung gegenüber Emulation

- Der Zugriff auf die Hardware kann direkt an die phy. Adresse durchgereicht werden
- durch das fehlende Übersetzen hat man weniger Overhead
- und dadurch ist ein solches System schneller als ein Emuliertes
- man hat dadurch auch weniger Kompatibilitätsprobleme

## Vorteile von Emulation gegenüber Virtualisierung

- Man kann auf einer Hardware viele verschiedene andere Hardware emulieren (zur Konsolidierung)
- Falls ein Emulator existiert, kann die darunter liegende Hardware getauscht werden





# WINE

## WINE

- WINE Projekt startet 1993 (**W**ine **I**s **N**ot an **E**mulator)
- Keine Emulatorschicht
- x86 Architektur ist Voraussetzung
- Aufrufe werden nur zu den entsprechenden Linux-Kernel Funktionen umgeleitet





# Virtualisierung von Peripherie-Geräten

- Bei der Peripherie-Geräten wird im Betriebssystem ein Treiber/Modul geladen
- Dieses Modul simuliert ein Gerät, das am Rechner angeschlossen ist
- Beispiele:
  - Virtuelle CD-Laufwerke (Daemon Tools, Virtual CD)
  - Drucker Spooler
  - RAID/LVM
  - iSCSI
  - rCAPI
  - usw.







# Grundlagen der x86 Architektur

- Der x86 Prozessor hat ein Ringschutzkonzept
- Es gibt 4 Ringe, der innerste ist Ring 0, der äußerste Ring 3
- Je innerer der Ring desto mehr Privilegien hat der Prozess auf Hardware zuzugreifen
- In Ring 0 (Kernel-Space) läuft meistens der Kernel/Module/Treiber eines Betriebssystems
- In Ring 3 (User-Space) die Anwendungen
- Die meisten Betriebssysteme benutzen nur 2 Ringe (0 und 3)

- Ein Prozess kann zu einem Zeitpunkt nur auf einem Ring laufen
- Ein Prozess kann nicht selbst "den Ring wechseln"
- Möchte ein Prozess eine Operation aufrufen, die in einem inneren Ring liegt, erzeugt der Prozessor eine Exception, die behandelt werden muss
- Wird diese nicht behandelt: Schutzverletzung, Prozessabsturz

# Agenda

1

## Begriffsabgrenzung Virtualisierung

- Definition

2

## Virtualisierungsarten

- Hardwarevirtualisierung
- Softwarevirtualisierung
  - Anwendungsvirtualisierung
  - Virtualisierung von Peripherie-Geräten
  - Virtualisierung auf Betriebssystemebene

3

## Virtualisierung von Betriebssystemen/Rechnern

- Allgemeines
- **User Mode Linux**
- Emulation auf Applikationsebene
- Para-Virtualized-Machines @ Xen
- Hypervisor-Virtualized-Machines @ Xen
- Vergleich der 3 Virtualisierungsarten (Xen, VMware, UML)

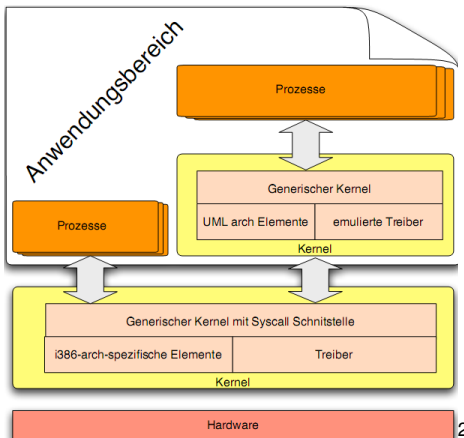


# Der Anfang von Betriebssystemvirtualisierung

## User Mode Linux

- 1999 veröffentlichte Jeff Dike einen Patch mit dem sich der Linux-Kernel sich selbst als unprivilegierten Prozess aufrufen kann
- Folge: man konnte mehrere Instanzen des Betriebssystems auf dem selben Rechner fahren
- Implementierung (User Mode) wurde mit der arch-Schnittstelle im architekturunabhängigen Teil des Linux Kernels umgesetzt
- Startet man den Kernel also im User Mode, werden alle hardwareunabhängigen Aufrufe aus der VM vom Host Kernel angenommen und bedarfsweise an den Gast-Kernel durchgereicht (gleich mehr dazu)

# Gesamtarchitektur eines UML Systems



<sup>2</sup>Quelle:

<http://www.lrr.in.tum.de/stodden/teaching/sem/virt/ss06/doc/virt06-07-20060531-kern-doc-20-20Paravirtualisierung.pdf>

# Funktionsweise der Architektur

## UML im TT-Mode

- Für jeden Gast betreibt der UML-Host einen Tracing Thread
- Für jeden Prozess eines Gasts wird ein Host Prozess angelegt
- d.h. jeder Prozess läuft im Prinzip auf dem Host Kernel (nur unterschieden durch den Tracing Thread)

## Aufgabe des Tracing Threads

- Setzt ein Gast ein Signal ab, wird dieses für den Hostkernel, falls es ein HW-unabhängiges Signal ist, immer in einen *getpid()* Call umgewandelt/"nullifiziert"
- Der gleiche Call wird aber unmodifiziert an den Gastkernel weitergereicht.

## Hää?? Wieso mit `getpid()` überschrieben?

- Die Theorie eines solchen Systems ist oft ideell
- In der Realität existiert ein UML-Debug-Handler im Host Kernel
- Die Debugschnittstelle wurde extra für UML in den Kernel eingepflegt
- Dieser Handler springt vor jedem Systemcall an und greift diesen Call auf
- Interessanterweise gibts aber keine Möglichkeit den Handler zu deaktivieren
- Es gibt auch keine Möglichkeit den Systemaufruf im Hostkernel abubrechen oder zu löschen
- Deswegen wird dieser Call mit einem `getpid()` Call überschrieben

# Vorteile und Probleme von UML

## Vorteile im Allgemeinen

- VMs laufen fast ohne Overhead (Im Vergleich zu emulierten Systemen)
- Keine eigenen Treiber auf der VM benötigt
- Keine Modifikation des Gast- oder Hostsystems nötig (da UM-Kernel wie ein normaler Prozess behandelt wird)

## Probleme im TT-Mode

- Prozessverwaltung ist mit Host gekoppelt:
  - Für jeden Prozess im Gast wird ein entsprechender Prozess im Host angelegt. (Prozesskopplung)
  - Folge bei vielen Gästen und vielen Prozessen: Performanceprobleme
- Speichervirtualisierung kaum umgesetzt: Einem Gast ist es möglich den Kernelspeicher zu überschreiben und aus dem Gast-Speicherbereich auszubrechen
- Kontextwechsel im Gastsystem können zu Leistungseinbußen führen, da jeder Wechsel explizit genehmigt werden muss vom Hostkernel (4 Stück: Syscall <-> TracingThread <-> Hostkernel) <sup>a</sup>

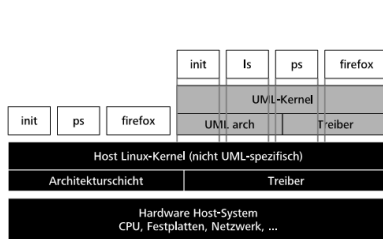
---

<sup>a</sup>Quelle: <http://www.uni-koblenz.de/vnuml/docs/vnuml/uml.pdf>

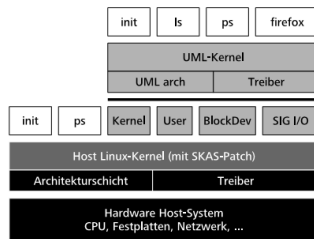
# Problemlösung mit dem SKAS Patch

## Seperate Kernel Address Space Eigenschaften:

- Wird auf dem Kernel des Hosts implementiert
- Löst die Probleme mit Prozessverwaltung, fehlende Speichervirtualisierung und auch das Kontextwechselproblem



Prozessarchitektur eines UML-Systems  
Prozesse im UML-Bereich sind für den Host-Kernel sichtbar



Prozessarchitektur eines UML-Systems mit SKAS-Patch  
UML-Bereich wird hinter vier Prozessen verborgen:  
Kernel- und User-Bereich, Blockgerätreiber und  
Signalvermittlung für Ein-/Ausgabe

## Folgen durch den SKAS Patch

- Durch die vorhandene Speichervirtualisierung wird pro Gast ein separater Kernel-Datenspeicher verwendet
- Der Host-Kernel Speicher wird durch das nun vorhandene Memory Management geschützt
- Durch eine modifizierte *ptrace()* Funktion (Prozesstracing) im Kernel werden nun nur noch 2 Kontextswiches benötigt
- Insgesamt ist das System ca. doppelt so schnell (gemessen bei Kernel-Compiling) <sup>4</sup>
- Es gibt nun nur noch 4 Prozesse auf dem Host System für alle UML Gäste
  - UML Kernel Thread - führt Code im UML Kernelspace aus
  - UML Userspace Thread - führt Code im UML Userspace aus
  - Asynchroner I/O Treiber Thread - virtuelle Block Devices
  - Emulator-Thread für WRITE-SIGIO - Hilfsprozess, der Signale für *write()* calls

---

<sup>4</sup><http://user-mode-linux.sourceforge.net/old/skas.html>



# Agenda

- 1 **Begriffsabgrenzung Virtualisierung**
  - Definition
- 2 **Virtualisierungsarten**
  - Hardwarevirtualisierung
  - Softwarevirtualisierung
    - Anwendungsvirtualisierung
    - Virtualisierung von Peripherie-Geräten
    - Virtualisierung auf Betriebssystemebene
- 3 **Virtualisierung von Betriebssystemen/Rechnern**
  - Allgemeines
  - User Mode Linux
  - **Emulation auf Applikationsebene**
  - Para-Virtualized-Machines @ Xen
  - Hypervisor-Virtualized-Machines @ Xen
  - Vergleich der 3 Virtualisierungsarten (Xen, VMware, UML)

# Virtual Machine Monitor (VMM)

- VMware veröffentlichte 1999 ihre Software VMware Workstation
- Virtualisierung eines kompletten x86 Rechners auf einem anderen x86er
- Neu: Eigenes BIOS, Eigene Hardware (begrenzt)
- Auch Neu: **Virtual Machine Monitor**, eine Anwendung überhalb des Betriebssystems, dass die vorhandene Hardwareresourcen an die VMs durch Emulation verteilt

## VMM - Scan-before-Execution / PreScan

- Der VMM überprüft parallel zur Laufzeit den Programmcode der VM **bevor** der ausgeführt wird
- wird dabei ein abzufangender Befehl erkannt, wird dort ein Breakpoint gesetzt und die Suche endet erstmal
- Gibt es konditionale Sprünge, werden beide verfolgt (bis maximale Suchtiefe erreicht)
- Falls etwas nicht gescannt werden konnte, wird der Rest der Sprünge emuliert
- Dieses Verfahren wird **Scan-before-Execution** genannt

# Weitere wichtige Teile in der VMware Architektur

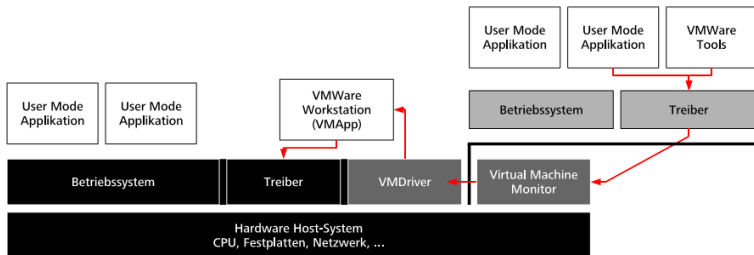
## VMDriver / VMXDriver

- Modul im Hostsystem
- dient als Kommunikationsschnittstelle für I/O Verkehr und Gast <-> VMApp
- Speicherverwaltung
- Gerätemanagement

## VMApp

- Oberfläche für den Endanwender
- Kommunikationsschnittstelle für Treiber
- Kommunikationsweg Hosttreiber <-> VMApp <-> VMXDriver

# Architekturübersicht VMware Workstation



# Agenda

- 1 **Begriffsabgrenzung Virtualisierung**
  - Definition
- 2 **Virtualisierungsarten**
  - Hardwarevirtualisierung
  - Softwarevirtualisierung
    - Anwendungsvirtualisierung
    - Virtualisierung von Peripherie-Geräten
    - Virtualisierung auf Betriebssystemebene
- 3 **Virtualisierung von Betriebssystemen/Rechnern**
  - Allgemeines
  - User Mode Linux
  - Emulation auf Applikationsebene
  - **Para-Virtualized-Machines @ Xen**
  - Hypervisor-Virtualized-Machines @ Xen
  - Vergleich der 3 Virtualisierungsarten (Xen, VMware, UML)

# Para, was?

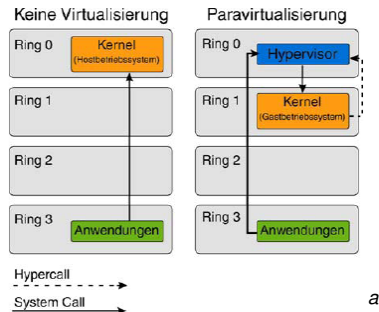
- Paravirtualisierung ist ein Verfahren zur effizienten Virtualisierung von x86 Architekturen mit einem Hypervisor und einem modifizierten Gast-Kernel
- Seit neustem auch modifizierte Kernel für Windows XP, 2003 und 2008 verfügbar

## Hypervisor

- Ein Hypervisor ist im Grunde ein Virtual Machine Monitor
- Er besteht aus einem minimalen Betriebssystem, das im Ring 0 des CPUs arbeitet
- Die Aufgaben sind ebenfalls Speicherverwaltung, Scheduling, Geräteverwaltung
- Ein Hypervisor emuliert im Paravirtualisierungsmodus keine Ressourcen, aber Geräte (Netzwerkkarte, Festplatten. . .)

# Ringbelegung und Hypercalls

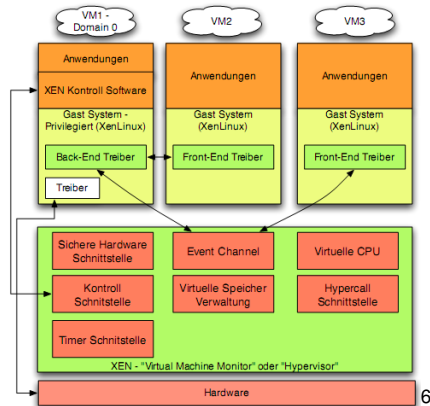
- Gastkernel weicht aus Ring 0 in Ring 1
- Problem: Gastkernel in Ring 1, will aber auf Hardware zugreifen!
- Lösung: Hypercalls



<sup>a</sup>Quelle: Informatik Spektrum 32  
3/2009 S.205



# XEN Architektur



<sup>6</sup>Quelle:

<http://www.lrr.in.tum.de/stodden/teaching/sem/virt/ss06/doc/virt06-07-20060531-kern-doc%20-%20Paravirtualisierung.pdf>

# Agenda

- 1 **Begriffsabgrenzung Virtualisierung**
  - Definition
- 2 **Virtualisierungsarten**
  - Hardwarevirtualisierung
  - Softwarevirtualisierung
    - Anwendungsvirtualisierung
    - Virtualisierung von Peripherie-Geräten
    - Virtualisierung auf Betriebssystemebene
- 3 **Virtualisierung von Betriebssystemen/Rechnern**
  - Allgemeines
  - User Mode Linux
  - Emulation auf Applikationsebene
  - Para-Virtualized-Machines @ Xen
  - **Hypervisor-Virtualized-Machines @ Xen**
  - Vergleich der 3 Virtualisierungsarten (Xen, VMware, UML)

# Unmodifizierte Gäste auf einem Hypervisor??

- Unterstützt ein Prozessor entweder die Intel **V**irtual **M**achine **E**xtenstion (VMX) oder den AMD **S**ecure **V**irtual **M**achine (SVM) Befehlssatz
  - 2 Modi des Prozessors möglich
  - VMX-Root-Modus (Hypervisor, Ring -1 am rechnen)
  - VMX-Non-Root-Modus (Gast-OS, Ring 0-3 am rechnen)

7

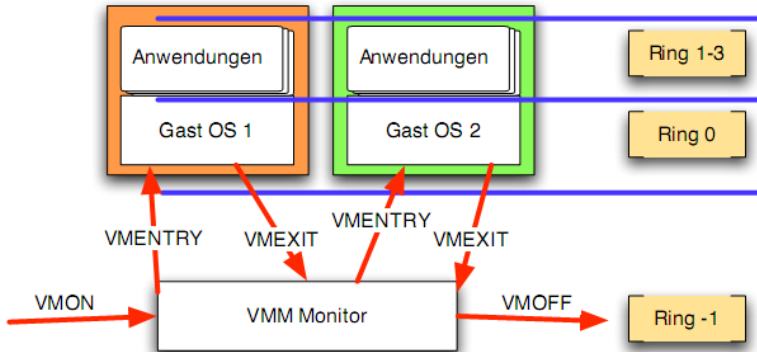
## CPU Befehle bei VT bzw. SVM

- *VMON*: CPU in Virtualisierungsmodus schicken
- *VMENTRY* (von VMM aus): Übergabe der Ringe(0-3) an Gast
- *VMEXIT* (von VMM aus): Abgabe der Ringe an VMM
- *VMOFF*: CPU aus Virtualisierungsmodus holen

<sup>7</sup>Quelle

<http://www.lrr.in.tum.de/stodden/teaching/sem/virt/ss06/doc/virt06-07-20060531-kern-doc%20-%20Paravirtualisierung.pdf>

# Ringaufbau und Zeitliches Beispiel



<sup>8</sup>Quelle:

<http://www.lrr.in.tum.de/stodden/teaching/sem/virt/ss06/doc/virt06-07-20060531-kern-doc%20-%20Paravirtualisierung.pdf>

## Folgen von VMX/SVM

- Eine VM kann nicht mehr unterscheiden, ob sie nativ oder virtualisiert betrieben wird (Ring 0 steht immer zur Verfügung)
- Kein Performanceverlust, da keine Emulation von Hardware nötig ist

## Unterschiede Vanderpool / Pacifica

- Technologien nicht kompatibel
- AMD's Pacifica virtualisiert ebenfalls den bei Intel hardwaretechnisch gelösten Speichercontroller
- Pacifica hat ebenfalls einen **Device Exclusion Vector** (DEV) integriert, der es VMs ermöglicht auch Geräte zu benutzen, die ohne CPU auf Speicher zugreifen können.

# Agenda

- 1 **Begriffsabgrenzung Virtualisierung**
  - Definition
- 2 **Virtualisierungsarten**
  - Hardwarevirtualisierung
  - Softwarevirtualisierung
    - Anwendungsvirtualisierung
    - Virtualisierung von Peripherie-Geräten
    - Virtualisierung auf Betriebssystemebene
- 3 **Virtualisierung von Betriebssystemen/Rechnern**
  - Allgemeines
  - User Mode Linux
  - Emulation auf Applikationsebene
  - Para-Virtualized-Machines @ Xen
  - Hypervisor-Virtualized-Machines @ Xen
  - **Vergleich der 3 Virtualisierungsarten (Xen, VMware, UML)**

# Relativer Performancevergleich aller Systeme zu Linux

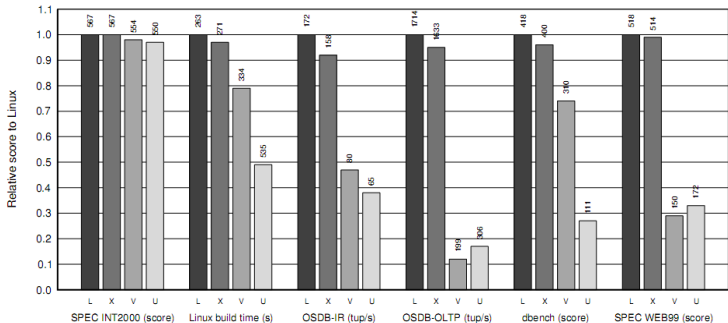


Figure 3: Relative performance of native Linux (L), XenLinux (X), VMware workstation 3.2 (V) and User-Mode Linux (U). <sup>9</sup>

<sup>9</sup>Quelle:

<http://www.cl.cam.ac.uk/research/srg/netos/papers/2003-xensosp.pdf>

# Fazit

- Virtualisierung eignet sich für viele Bereiche, bei der keine spezielle Hardware benötigt wird
- Virtualisierung birgt kaum Risiken oder Nachteile, wenn man sie richtig Einsetzt
- Administratoren brauchen fachübergreifendes Wissen  
z.B.: Netzwerk-Utills in Linux
- Große Probleme bekommt man bei allen Geräten, die nicht blockorientiert Arbeiten: ISDN Karten, USB Dongles für Lizenzen, alle A/V-USB-Geräte
- XEN ist inzwischen altbewährt, aber auf dem Rückmarsch.  
Neues Produkt ähnliche Technologie: Kernel Based Virtual Machine (KVM)



# Gibt es Fragen

