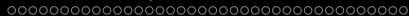


Virtualization – Fundamentals

- By using **virtualization**, the resources of a computer system can be split and used by multiple independent operating system instances
- Several fundamentally different approaches and technologies exist to implement virtualization
- Each **virtual machine** (VM). . .
 - behaves like any other computer, with own components
 - runs inside an isolated environment on a physical machine
- Inside a VM, an operating system with applications can be installed, exactly like on a physical computer
 - The applications do not notice that they are located inside a VM
- Requests of the operating system instances are captured by the virtualization software and converted for the existing physical or emulated hardware
 - The VM itself does not become aware of the virtualization layer between itself and the physical hardware



History of Virtualization

- Virtualization is not a new concept
 - Introduced in the 1960s by IBM for mainframes
- 1970/71: IBM introduced the Virtual Machine Facility/370 (VM/370)
 - On this platform, multi-user operation is implemented by using multiple single-user mode instances, which are executed in virtual machines
 - Each VM is a complete duplicate of the underlying physical hardware

Sources

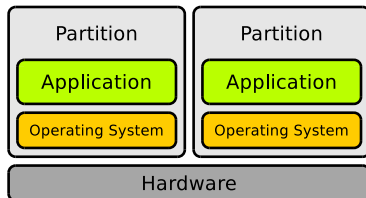
- Creasy RJ. **The origin of the VM/370 time-sharing system.** IBM Journal of Research and Development 25 (1981), No. 5, 483–490
- Amit Singh. **An Introduction to Virtualization.** 2004
<http://www.kernelthread.com/publications/virtualization/>

Virtualization Concepts

- Different virtualization concepts exist:
 - Partitioning
 - Hardware emulation
 - Application virtualization
 - Full virtualization (Virtual Machine Monitor)
 - Paravirtualization (Hypervisor)
 - Hardware virtualization
 - Operating system-level virtualization / Container / Jails
 - Storage virtualization (SAN)
 - Network virtualization (VLAN)
 - ...

Partitioning

- If partitioning is used, the total amount of resources can be split to create subsystems of a computer system
 - Each subsystem may contain an executable operating system instance
 - Each subsystem can be used like an independent computer system
- The resources (CPU, main memory, storage...) are managed by the **firmware** of the computer and assigned to the VMs
- Partitioning is used, e.g. in IBM mainframes (zSeries) and midrange systems (pSeries) with Power5/6/7 CPUs
 - Resource allocation is possible during operation without having to restart
 - On a modern mainframe computer several hundred to thousands of Linux instances to operate simultaneously
- Modern CPUs only support the partitioning of the CPU itself and not of the entire system (Intel Vanderpool, AMD Pacifica)
 - Partitioning is not used for desktop environments



Selection of Emulators

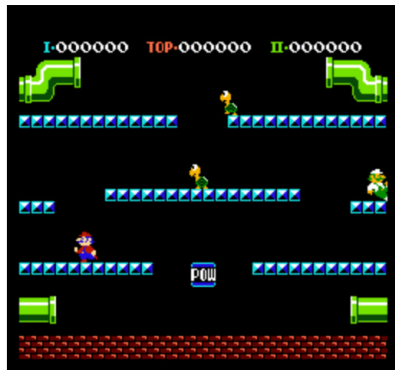
| Name | License | Host | Emulated architecture | Guest system |
|---------------------|----------------|--|---------------------------------------|--------------------------------|
| Bochs v2.3.6 | LGPL | Linux, Solaris, MacOS, Windows, IRIX, BeOS | x86, AMD64 | Linux, DOS, BSD, Windows, BeOS |
| QEMU v0.9.0 | GPL | Linux, BSD, Solaris, BeOS, MacOS-X | x86, AMD64, PowerPC, ARM, MIPS, Sparc | Linux, MacOS-X, Windows, BSD |
| DOSBox v0.72 | GPL | Linux, Windows, OS/2, BSD, BeOS, MacOS-X | x86 | DOS |
| DOSEMU v1.4.0 | GPL | Linux | x86 | DOS, Windows bis 3.11 |
| PearPC v0.4.0 | GPL | Linux, MacOS-X, Windows | PowerPC | Linux, MacOS-X, BSD |
| Baseilisk II v0.9-1 | GPL | Linux, various UNIX, Windows NT4, BeOS, Mac OS, Amiga OS | 680x0 | MacOS ≤ 8.1 |
| Wabi v2.2 | proprietary | Linux, Solaris | x86 | Windows 3.x |
| MS Virtual PC v7 | proprietary | MacOS-X | x86 | Windows, (Linux) |
| M.A.M.E. v0.137 | MAME-Lizenz | Linux, Windows, DOS, BeOS, BSD, OS/2 | various Arcade | various Arcade |
| SheepShaver | GPL | Linux, MacOS-X, BSD, Windows, BeOS | PowerPC, 680x0 | MacOS 7.5.2 bis MacOS 9.0.4 |
| Hercules 3.07 | QPL | Linux, MacOS-X, BSD, Solaris, Windows | IBM mainframes | IBM System/360, 370, 390 |

- The table is not complete!
- Many more emulators exist

Example of a current Emulator - JSNES

Ben Firshman
JSNES

- JSNES emulates the Nintendo Entertainment System (NES)
- The emulator is implemented in JavaScript and executes in the browser
- <http://fir.sh/projects/jsnes/>
- Free Software (GPLv3)



Mario Bros.

Running: 44.40 FPS

Latest Development: Browser emulates PC – jslinux

<http://www.webmonkey.com/2011/05/yes-virginia-that-is-linux-running-on-javascript/>

Date: May 18th 2011

Author: Scott Gilbertson

JavaScript never seems to get any respect. It's not a real programming language, detractors complain, it's just some script language that runs in the web browser. We're not sure what makes JavaScript less „real“ to some, but thanks to today's web browsers, JavaScript has become a very powerful language. Powerful enough to run Linux in your web browser. French developer Fabrice Bellard has built a **JavaScript-based x86 PC emulator capable of running Linux inside a web browser.**

If you'd like to try it out, point Firefox 4 or Chrome 11 to the demo page. Keep in mind that this is just Linux, no X Window or other graphical interface, just the command line, a small C compiler and QEmacs, Bellard's emacs clone. Still, it's really Linux, really running in your web browser, really using JavaScript to emulate hardware.

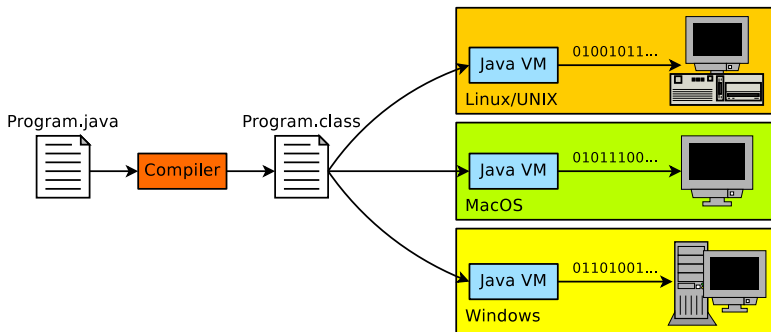
```
TCP reno registered
checking if image is initramfs...it isn't (bad gzip magic numbers); looks like a
p initrd
freeing initrd memory: 2048k freed
total HugeTLB memory allocated, 0
io scheduler noop registered
io scheduler anticipatory registered
io scheduler deadline registered
io scheduler cfq registered (default)
Real Time Clock Driver v1.12ac
JS clipboard: I/O at 0x03c0
Serial: 8250/16550 driver $Revision: 1.90 $ 4 ports, IRQ sharing disabled
serial0250: ttyS0 at I/O 0x3f8 (irq = 4) is a 16450
RAMDISK driver initialized: 16 RAM disks of 4096K size 1024 blocksize
loop: loaded (max 8 devices)
TCP cubic registered
NET: Registered protocol family 1
NET: Registered protocol family 17
Using IPI Shortcut mode
Time: pit clocksource has been installed.
RAMDISK: ext2 filesystem found at block 0
RAMDISK: Loading 2048KiB [1 disk] into ram disk... done.
EXT2-fs warning: maximal mount count reached, running e2fsck is recommended
VFS: Mounted root (ext2 filesystem).
freeing unused kernel memory: 124k freed
Booted in 8.961 s
Welcome to JS/Linux
# uname -a
Linux (none) 2.6.20 #2 Mon Aug 8 23:51:02 CEST 2011 i586 GNU/Linux
#
© 2011 Fabrice Bellard - http://bellard.org/jslinux/
```

Image Source: <http://bellard.org/jslinux/>

Application Virtualization

- Applications are executed inside a virtual environment, which uses local resources and provides all the components, which are required by the application
 - The VM is located between the executed application and the operating system
- Popular example: Java Virtual Machine (JVM)
 - The JVM is the part of the Java Runtime Environment (JRE), which executes the Java bytecode
 - The JVM is for Java programs the interface to the computer system and its operating system
- Advantage: Platform independence
- Drawback: Reduced performance, compared with native execution

Principle of the Java Virtual Machine (JVM)



- The compiler `javac` compiles source code into architecture-independent `.class` files, which contain bytecode, that can be executed in the Java VM
- The program `java` launches a Java application inside a Java VM

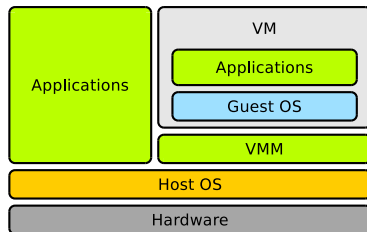
VMware ThinApp

<http://www.vmware.com/products/thinapp/>

- Further example of application virtualization: VMware ThinApp
 - Until 2008, the software was named Thinstall
- Packs Windows applications into single .exe files
- The application becomes portable and can be used without local installation
 - Applications can, e.g. be executed from an USB flash memory drive
- No entries are inserted into the Windows registry and no environment variables or DLL files are created on the system
- User preferences and created documents are stored inside a separate sandbox
- Drawback: The software only supports Microsoft Windows

Full Virtualization (1/3)

- Full virtualization software solutions provide each VM a complete virtual PC environment, including an own BIOS
 - Each guest operating system gets its own VM with virtual resources (e.g. CPU, main memory, storage drives, network adapters)
- A **Virtual Machine Monitor** (VMM) is used
 - The VMM is also called **Type-2 hypervisor**
 - The VMM runs *hosted* as an application in the host operating system
 - The VMM distributes hardware resources to VMs
- Some hardware components are emulated, because they are not designed for the concurrent access from multiple operating systems
 - Example: Network adapters
 - The emulation of popular hardware avoids driver issues



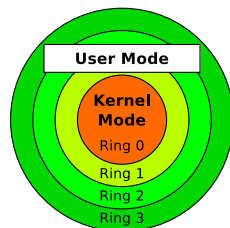
Virtualization Basics of the x86 Architecture (1/2)

- x86-compatible CPUs contain 4 privilege levels
 - Objective: Improve stability and security
 - Each process is assigned to a ring permanently and can not free itself from this ring

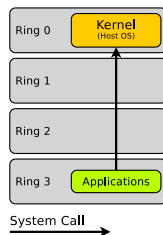
Implementation of the privilege levels

- The register CPL (Current Privilege Level) stores the current privilege level
- Source: Intel 80386 Programmer's Reference Manual 1986
<http://css.csail.mit.edu/6.858/2012/readings/i386.pdf>

- In ring 0 (= **kernel mode**) runs the kernel
 - Processes in this ring have full access to the hardware
 - The kernel can address physical memory (\implies Real Mode)
- In ring 3 (= **user mode**) run the applications
 - Processes in this ring can only access virtual memory (\implies Protected Mode)



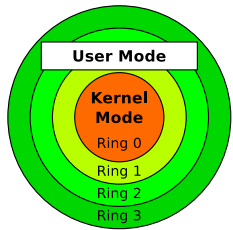
Without Virtualization



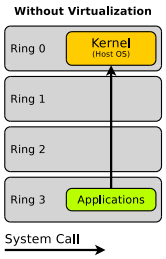
Virtualization Basics of the x86 Architecture (2/2)

Modern operating systems only use 2 privilege levels (rings)

- Reason: Some hardware architectures (e.g. Alpha, PowerPC, MIPS) support only 2 privilege levels
- Exception: OS/2 uses ring 2 for applications, which are allowed to access hardware and input/output interfaces (e.g. graphics drivers)



- If a user-mode process must carry out a higher privileged task (e.g. access hardware), it can tell this the kernel via a **system call**
 - The user-mode process generates an exception, which is caught in ring 1 and handled there



Full Virtualization Examples

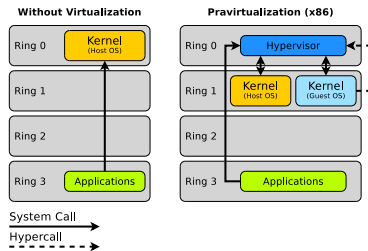
- Some virtualization solutions, which implement the VMM concept:
 - VMware Server, VMware Workstation and VMware Fusion
 - Microsoft Virtual PC (in the x86 version)
 - Parallels Desktop and Parallels Workstation
 - VirtualBox
 - Kernel-based Virtual Machine (KVM)
 - Mac-on-Linux (MoL)

Paravirtualization (1/4)

- No hardware is virtualized or emulated
 - Does not provide an emulated hardware layer to the guest operating systems, but only an application interface
- Guest operating systems use an abstract management layer (⇒ **hypervisor**) to access the physical resources
 - Hypervisor is a **meta operation system**, which is reduced to a minimum
 - The hypervisor distributes hardware resources among the guest systems, the same way, an operating system would distribute hardware resources among running processes
 - The hypervisor is a **Type-1 hypervisor** and runs *bare metal*
 - A meta operation system allows the independent operation of different applications and operating systems on a single CPU
- The hypervisor runs in located in the privileged ring 0
 - The host operating system is relocated to the less privileged ring 1
 - A host operating system is required because of the device drivers

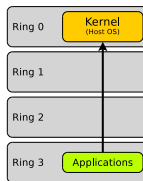
Paravirtualization (2/4)

- The host operating system is relocated from ring 0 to ring 1
 - Therefore, the kernel can not execute privileged instructions
 - Solution: The hypervisor provides **hypercalls**
- Hypercalls are similar to system calls
 - The interrupt numbers are different
 - If an application requests the execution of a system call, a replacement function in the hypervisor is called
 - The hypervisor orders the execution of the system call via the kernel API of the operating system

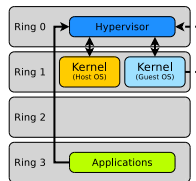


Paravirtualization (3/4)

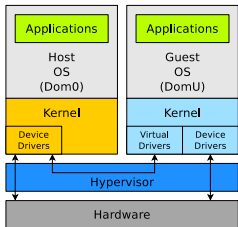
- Kernels of guest operating systems need to be modified in a way that any system call for direct access to hardware is replaced by the corresponding hypercall
- Catching and verifying system calls by the hypervisor causes just little performance loss
- Examples: Xen, Citrix Xenserver, Virtual Iron, VMware ESX Server

Without Virtualization

System Call →
Hypercall - - - - ->

Paravirtualization (x86)

Paravirtualization (4/4)



- VMs are called **unprivileged domain** (DomU)
- The hypervisor replaces the host operating system
 - But the developers can not develop all drivers from scratch and maintain them
 - Therefore, the hypervisor launches an (Linux) instance with its drivers and borrows them
 - This instance is called Domain0 (Dom0)

- Drawbacks:

- Kernels of guest operating systems must be modified (adapted) for operation in the paravirtualized context
- Rights holders of proprietary operating systems often reject an adjustment because of strategic reasons
⇒ Often works only with open source operating systems

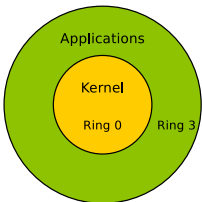
- Advantage:

- Better performance compared with VMM implementations

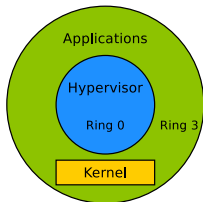
Problem: x86-64 Architecture

- The x86-64 architecture (e.g. IA64) does not implement ring 1 and 2

Without Virtualization

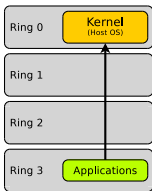


Pravirtualization (IA64)

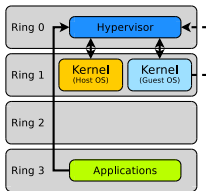


- In the x86-32 architecture, the hypervisor is located in ring 0
- In the x86-64 architecture, the operating system kernel is relocated to ring 3, where the applications are located

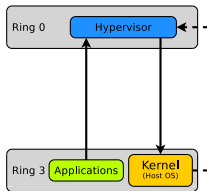
Without Virtualization



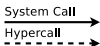
Pravirtualization (x86)



Pravirtualization (IA64)

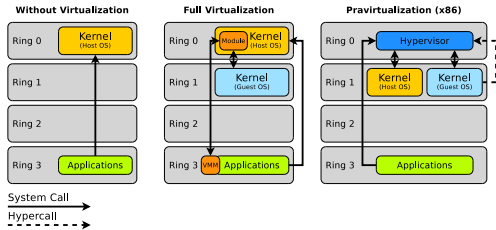


- Locating hardware drivers and applications in the same ring tends to be insecure



Summary: Virtualization vs. Paravirtualization

- **Paravirtualization** requires modified guest systems
 - Type-1 hypervisor runs *bare metal* (= replaces the host operating system)
 - Hypervisor runs in ring 0 and has full access to the hardware
 - Examples: VMware ESX(i), Xen, Microsoft Hyper-V
- **Full virtualization** supports unmodified guest systems
 - VMM (Type-2 hypervisor) runs *hosted* as an application in the host operating system
 - VMM runs in ring 3 at the level of the applications
 - Examples: VMware Workstation, KVM, Oracle VirtualBox, Parallels

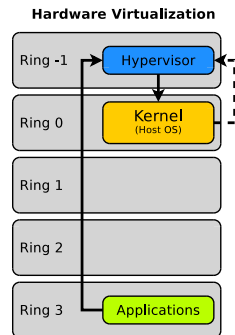


Hardware Virtualization (1/2)

- Current CPUs from Intel and AMD contain virtualization extensions for hardware virtualization
 - Advantage: Unmodified operating systems can be used as guest systems
 - The solutions from Intel and AMD are similar but incompatible
 - Since 2006, AMD64 CPUs contain the Secure Virtual Machine (**SVM**) instruction set
 - The solution is called **AMD-V** and was previously called **Pacifica**
 - The solution from Intel is called **VT-x** for IA32 CPUs and **VT-i** for Itanium CPUs
 - The solution of Intel was previously called **Vanderpool**
- Since Xen version 3, the software supports hardware virtualization
 - Windows Server since Version 2008 (Hyper-V) uses hardware virtualization
 - VirtualBox supports hardware virtualization
 - KVM can only operate with CPUs, which implement hardware virtualization

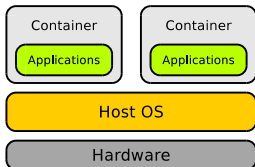
Hardware Virtualization (2/2)

- The hardware virtualization implementation contains a modification of the privilege levels
- A new ring (\Rightarrow ring -1) for the hypervisor is added
 - The hypervisor or VMM runs in ring -1 and at any time has the full control over the CPU and the resources, because with ring -1 an increased privilege level is implemented compared with ring 0
- VMs, executed inside ring 0 are called HVM
 - HVM = Hardware Virtual Machine
- Advantages:
 - Guest operating systems do not need to be modified (adapted)
 - Even proprietary operating systems (e.g. Windows) can be used as guest systems
 - In contrast to paravirtualization (IA64), the kernel is not executed in the privilege level of the applications



Operating System-level Virtualization / Containers (1/2)

- Under a single kernel, multiple identical, isolated system environments are executed
 - No additional operating system is started
 - An isolated runtime environment is created
 - All running applications use the same kernel
 - This kind of virtualization is called **Containers** in SUN/Oracle Solaris
 - This kind of virtualization is called **Jails** in BSD



- Applications only see applications from the same virtual environment
- One advantage is the low overhead, because the kernel manages the hardware as usual
- Drawback: All virtual environments use the same kernel
 - Only independent instances of the same operating system are started
 - It is impossible to start different operating systems at the same time

Operating System-level Virtualization / Containers (2/2)

- This type of virtualization is used to execute applications in isolated environments with high security
- Especially Internet service providers, which offer (virtual) root servers, or web services on multi-core processor architectures, use this type of virtualization
 - Little performance loss, high security level
- Examples:
 - SUN/Oracle Solaris (2005)
 - OpenVZ for Linux (2005)
 - Linux-VServer (2001)
 - FreeBSD Jails(1998)
 - Parallels Virtuozzo (2001, commercial version of OpenVZ)
 - FreeVPS
 - Docker (2013)
 - chroot (1982)

Docker



- Container are standard sized boxes to transport goods
- Docker offers a way to
 - package,
 - distribute and
 - run software.
- Containers can run on:
 - Linux (origin)
 - Mac (Beta, Yosemite 10.10)
 - Windows (Beta, Window 10 64 bit)
 - AWS
 - Azure

Docker Overview

Source: <http://pointful.github.io/docker-intro/>

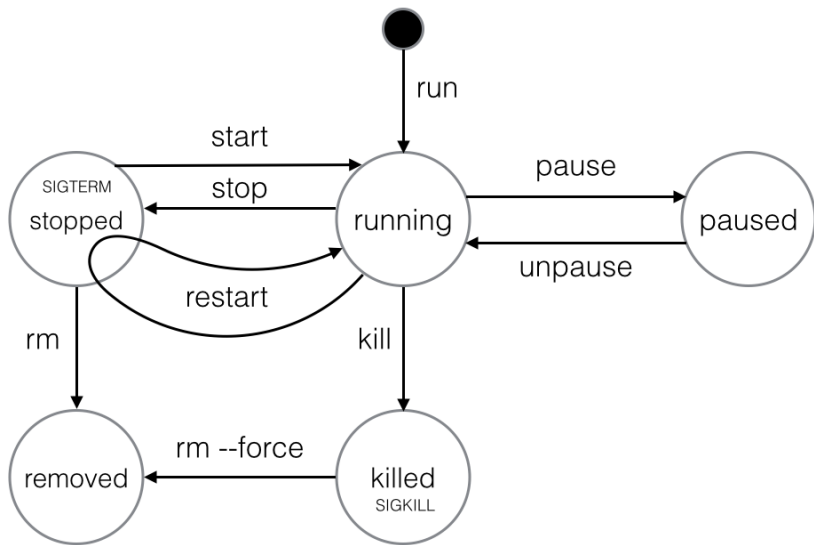
WHY

- Run everywhere
 - Regardless of kernel version (2.6.32+)
 - Regardless of host distro
 - Physical or virtual, cloud or not
 - Container and host architecture must match
- Run anything
 - If it can run on the host, it can run in the container
 - i.e. if it can run on a Linux kernel, it can run

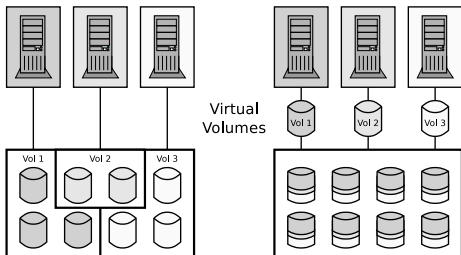
WHAT

- High Level – It's a lightweight VM
 - Own process space
 - Own network interface
 - Can run stuff as root
 - Can have its own `/sbin/init` (different from host)
 - „machine container“
- Low Level – It's chroot on steroids
 - Can also not have its own `/sbin/init`
 - Container=isolated processes
 - Share kernel with host
 - No device emulation (neither HVM nor PV) from host)
 - „application container“

Docker Container Lifecycle



Storage Virtualization



- Storage is provided to users in form of virtual drives (*volumes*)
- Logical storage is separated from physical storage

- Advantages:

- Users are independent from the physical limits of drives
- Reorganizing/expanding the physical storage does not disturb the users
- Redundancy is provided transparently in the background
- Better degree of utilization, because the physical storage can be split among the users in a more efficient way

- Drawback: Professional solutions are expensive

- Some Providers: EMC, HP, IBM, LSI and SUN/Oracle

Network Virtualization via Virtual Local Area Networks

- Distributed devices can be combined via VLAN in a single virtual (logical) network
 - VLANs separate physical networks into logical subnets (overlay networks)
 - VLAN-capable Switches do not forward packets of one VLAN into other VLANs
 - A VLAN is a network, over existing networks, which is isolated to the outside
 - Devices and services, which belong together, can be consolidated in separate VLANs
 - Advantage: Other networks are not influenced
⇒ Better security level

Helpful sources

- Benjamin Benz, Lars Reimann. *Netze schützen mit VLANs*. 11.9.2006
<http://www.heise.de/netze/artikel/VLAN-Virtuelles-LAN-221621.html>
- Stephan Mayer, Ernst Ahlers. *Netzsegmentierung per VLAN*. c't 24/2010. S.176-179

