

Component-based system for management of multilevel virtualization of networking resources

Robert Boczek

Dawid Ciepliński

June 12, 2011

Contents

1	Introduction	5
2	Context	7
2.1	QoS-aware networking	7
2.2	Resource virtualization approaches	7
2.3	Multilevel network virtualization	7
2.3.1	Virtual network resources	7
2.3.2	Fine-grained QoS control	7
2.3.3	Virtual appliances	7
2.3.4	„Network in a box” concept	7
2.4	Applications and benefits of virtual infrastructures	7
2.4.1	Testing and simulations	7
2.4.2	Improving server-side infrastructure scalability	7
2.4.3	Infrastructure as a service	7
2.4.4	The role of resource virtualization in the SOA stack	7
3	Requirements analysis	9
3.1	Functional requirements	9
3.1.1	Instantiation	9
3.1.2	Discovery	9
3.1.3	Accounting	9
3.2	Non-functional requirements	9
3.3	Underlying environment characteristics	9
3.4	General approach and problems it imposes	9
3.4.1	Load balancing / Deployment	9
3.4.2	Infrastructure isolation	9
3.4.3	Broadcast domain preservation	9
3.4.4	Constraints	9
4	Solaris OS as a resource virtualization environment	11
4.1	General information	11
4.2	Lightweight OS-level virtualization with Solaris Containers	11
4.2.1	General information	11
4.2.2	Isolation of processes	12

4.2.3	Advantages of Containers technology when compared to non-virtualized environments	12
4.2.4	Virtual appliances	13
4.3	Crossbow - network virtualization technology	14
4.3.1	Crossbow architecture	14
4.3.2	Virtualization lanes	15
4.3.3	Dynamic polling	16
4.3.4	Virtual switching	16
4.3.5	Crossbow elements	18
4.4	Resource access control	19
5	The system architecture	21
5.1	High-level design	21
5.2	System components and their responsibilities	21
5.2.1	Assigner	21
5.2.2	Supervisor	21
5.2.3	Worker	21
5.3	Crossbow resources instrumentation	21
5.4	Domain model and data flows	21
6	Implementation	23
6.1	Implementation environment	23
6.2	Domain model transformation details	23
6.3	Low-level functions access	23
6.4	Building and running the platform	23
7	Case Study	25
7.1	Multimedia server	25
7.1.1	Scenario description	25
7.1.2	Resource access requirements	25
7.1.3	Providing tunable and scalable virtual infrastructure	25
8	Summary	27
8.1	Conclusions	27
8.2	Achieved goals	27
8.3	Further work	27

Chapter 1

Introduction

Chapter 2

Context

Chapter overview

2.1 QoS-aware networking

2.2 Resource virtualization approaches

2.3 Multilevel network virtualization

2.3.1 Virtual network resources

2.3.2 Fine-grained QoS control

2.3.3 Virtual appliances

2.3.4 „Network in a box” concept

2.4 Applications and benefits of virtual infrastructures

2.4.1 Testing and simulations

2.4.2 Improving server-side infrastructure scalability

2.4.3 Infrastructure as a service

2.4.4 The role of resource virtualization in the SOA stack

Summary

Chapter 3

Requirements analysis

Chapter overview

3.1 Functional requirements

3.1.1 Instantiation

3.1.2 Discovery

3.1.3 Accounting

3.2 Non-functional requirements

3.3 Underlying environment characteristics

3.4 General approach and problems it imposes

3.4.1 Load balancing / Deployment

3.4.2 Infrastructure isolation

3.4.3 Broadcast domain preservation

3.4.4 Constraints

Summary

Chapter 4

Solaris OS as a resource virtualization environment

Chapter overview

The chapter provides an overview of Oracle Solaris 10 and evaluates it as a platform for resource virtualization with respect to networking.

Three main components are presented and discussed: Containers - OS-level virtualization technology, Crossbow - network virtualization environment and the system's resource access control utilities.

4.1 General information

resource management (allocation/assignment) accounting

4.2 Lightweight OS-level virtualization with Solaris Containers

The concept of lightweight (OS-level) virtualization is commonly known and supported by most (TODO really?) modern operating systems. The solutions are either integrated into the system's kernel (Solaris Containers, AIX Workload partitions, BSD jails TODO cite) or provided by third-party manufacturers as kernel patches and utility software (OpenVZ and LXC for Linux OS).

4.2.1 General information

Containers (or zones) were introduced as of Solaris Operating System 10. They provide a way of partitioning systems' resources and allow for isolated application execution environment.

There are two types of zones: global and non-global. Global zone is the default one and is used to execute applications as well as to administrate the system. Non-global zones can be created from within the global zone only. (TODO cite sag (how zones work)) A single operating system instance with one global zone can host as many as 8192 non-global zones.

Zones can be assigned system resources such as CPU power, the amount of random-access memory or disk space quota. Network isolation is also supported at two levels: basic isolation at the

IP layer and advanced fine-grained network isolation and virtualization with Crossbow technology. (TODO verify the sentence!)

Each non-global zone can run a different set of applications, with optional system call translation (*Branded Zones Technology* TODO cite sag (branded zones technology)) thus emulating different operating environment. The user is able to create a branded container with Linux system call translation and run Linux-specific applications within the container without code recompilation.

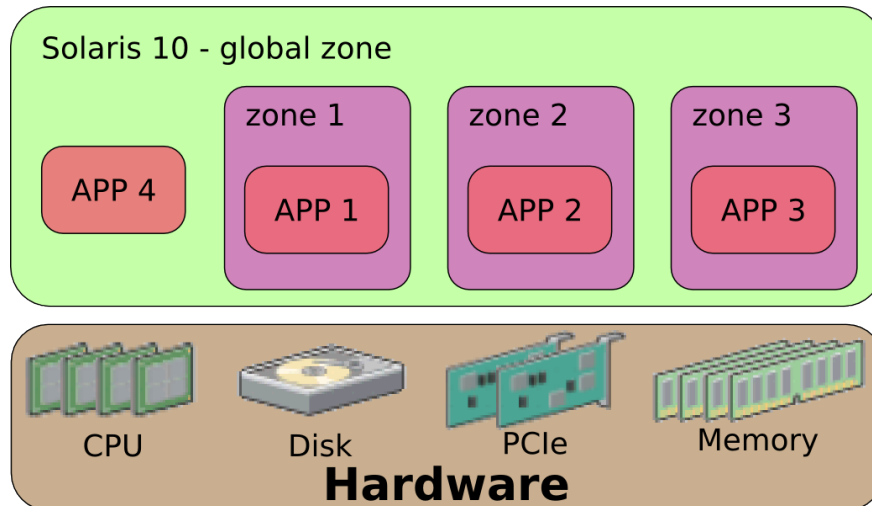


Figure 4.1: Solaris Containers hierarchy (source: wiki)

4.2.2 Isolation of processes

The Containers technology offers a high level of application security and isolation. This is accomplished by imposing software bounds on the resource usage and introduction of additional abstraction layer over hardware.

Every process and its children are bound to concrete zone and the assignment cannot be changed. Moreover, it is impossible for processes in distinct zones to monitor each other operation. They are not visible to each other and no standard interprocess communication can take place (standard==all except networking, TODO verify).

An application failure possibly affects only the processes in containing zone. (TODO easier recovery, independent container management)

4.2.3 Advantages of Containers technology when compared to non-virtualized environments

- low footprint and resource consumption

The technology allows for running a number of isolated instances of operating system sharing CPU time, physical network bandwidth, filesystem contents and binary code. The instances are able to execute applications with minimum overhead [?].

- snapshots, restoration, mobility (migration)

Container's contents are easily serializable with ZFS (TODO more?) snapshots.

Lightweight virtual machines can be transferred between physical nodes using network connection or other means of data transfer.

- flexibility and dynamic fine-grained resource allocation

The zones can be created and destroyed on demand, which allows for 1:1 zone:service mapping.

Resources (such as CPU, memory, network bandwidth) can be assigned dynamically without interrupting the container operation.

4.2.4 Virtual appliances

Virtual appliance is pre-built, pre-configured, ready-to-run enterprise application packaged along with an optimized operating system inside a virtual machine (definition from [?]). Solaris Containers fulfil all the requirements (TODO what requirements?) to be used as virtual appliance building block.

The main problem virtual appliances can solve is the complexity and duration of application deployment process. In general, a service deployment can be described as comprising the following stages: preparation (learning the dependencies), pre-installation, installation and post-installation. With traditional (non-virtualized) approach, these stages have to be repeated every time a service is deployed.

preparation → pre-installation → installation → post-installation

Figure 4.2: Traditional application deployment stages.

Virtual appliance approach makes it possible to reduce deployment time significantly [?]. This is achieved by performing most of the deployment stages once and storing the configured environment in a virtual appliance. The appliance can now be published in publicly-available repository for actual deployment on a host system.

1. preparation, pre-installation, installation, post-installation
2. (virtual) appliance preparation
3. appliance publication
4. appliance retrieval and activation
5. configuration adjustment

Figure 4.3: Deployment process with virtual appliances. Stage 1 is executed once.

It is possible to prepare sets of virtual appliances containing traditional services (such as application servers, database servers or media servers) as well as highly specialized networking-focused appliances that can act as routers, firewalls or load balancers. These Virtual (Network) Appliances,

together with other components provided by Solaris 10, can be leveraged to build fully virtual network topologies.

4.3 Crossbow - network virtualization technology

It is generally acknowledged that Crossbow was invented in China in 341 B.C but it was in middle ages when it earned its recognition. Very easy in use and simultaneously very effective. The Solaris Crossbow mechanism for QoS are just like real crossbow's very efficient in comparison to other existing QoS mechanisms and this similarity indicates the project name origin.

4.3.1 Crossbow architecture

One of the most important condition in terms of network virtualization is that network traffic should be insulated between virtual machines. This kind of isolation can be achieved by having a dedicated physical NIC, network cable and port from the switch to the virtual machine itself. Moreover, switch must also ensure sustainability on every port. In every other case virtual machines will definitely interfere between each other. In a particular case when we have to share physical NIC between virtual machines the most promising solution is to virtualize NIC hardware and the second layer of the OSI/ISO stack where sharing is fair and interferences will be avoided. These approach was adapted in the Crossbow architecture in OpenSolaris OS. Traffic separation is achieved by fundamental blocks of new architecture which are Virtual NICs (VNICs) created by dividing NIC into many VNICs. A VNIC can be created over NIC or Etherstub (more about them later) and be dynamically controlled by the bandwidth and CPU resources assigned to it. The crossbow architecture has introduced fully paralyzed network stack structure. Each stack could be seen as fully independent lane (without any shared locks, queues, and CPUs) therefore network isolation is guaranteed. Key concept is hardware classification performed by the NIC over which VNIC was created. Each lane has a dedicated buffer for Transmit (Tx) and Receive (Rx) ring. In case when load exceeds assigned limit packets must be dropped as it is wiser to drop them then to expend OS CPU resources.

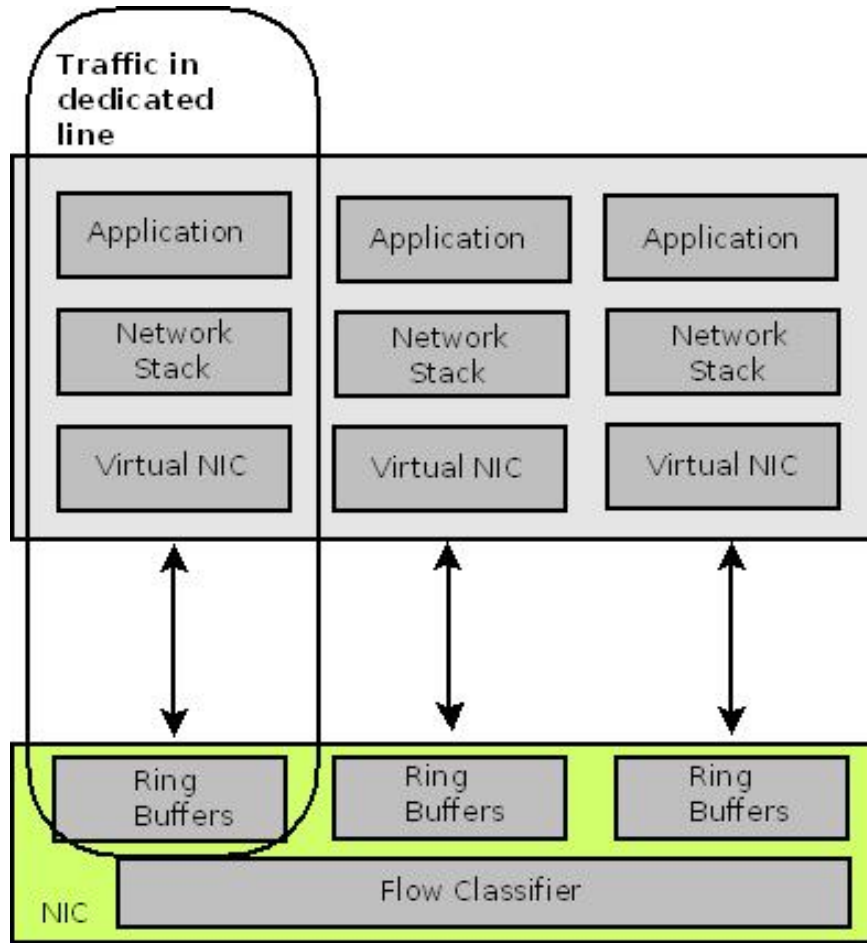


Figure 4.4: Dedicated lanes in the Crossbow architecture

4.3.2 Virtualization lanes

Virtualization lane is the most key component in the Crossbow architecture. Each lane consists some dedicated hardware and software that might be used to some concrete type of traffic. It usually would be composed of:

1. NIC resources(receive and transmit rings, interrupts, MAC address slots)
2. Driver resources(DMA bindings)
3. MAC layer resources (data structures, execution threads, locks)

A virtualization lane can be one of two types, hardware-based or software-based.

Hardware-based virtualization lanes

This type requires ability to partitioning resources from NIC. The minimum requirement is that a hardware-based lane should must have a dedicated receive ring. Other resources such as transmit

lane can be exclusive or shared between lanes. Each virtual machine could have one or more lanes assigned and the incoming packets would be distributed among them based on even scheduling unless some administrative policies were created, such as priority or bandwidth limit.

Software-based virtualization lanes

In case when NIC runs out of hardware-based virtualization lane, receive and transmit rings may be shared by multiple VNICs. The number of software-based virtualization lanes also often called softtrings is unlimited. The main disadvantage of software-based lanes is the lack of fairness and isolation which in fact is provided in hardware-based lanes. The received and sent rings may work also in mix mode, whereas some of the rings may be assigned to software and some may be assigned to hardware based lanes.

4.3.3 Dynamic polling

The Crossbow architecture proposed two types of working mode. Currently used mode is determined by traffic and load. Under low load, where the rate of arriving packets is lower than time of packet processing lane works in the interrupt mode which means that receive ring generates an interrupt when new packet arrives. However, when the backlog grows, the line switches to dynamic polling mode in which a kernel thread goes down to the receive ring in the NIC hardware to extract all outstanding packets in a single chain. Key aspect is that every virtualization lane works independently and transparently from each other. Usually only three threads are used per lane:

1. Poll thread which goes to the NIC hardware to get all packet chain
2. Worker thread which is responsible for protocol processing (IP and above) or delivers packets to virtual machine. Thread performs also any additional transmit work which is a natural requirement some concrete protocol, such as processing TCP packets that require sending ACK packets.
3. Transmit thread that is activated when if packets are being sent after transmit side flow control relief discharge, or after retrieving transmit descriptor. Application or virtual machine can transmit any packets without performing queuing because of flow control or context switching.

4.3.4 Virtual switching

Virtual switches are always created implicitly when the first VNIC is defined under existing NIC and could never be accessed directly nor be visible by any user (even administrator).

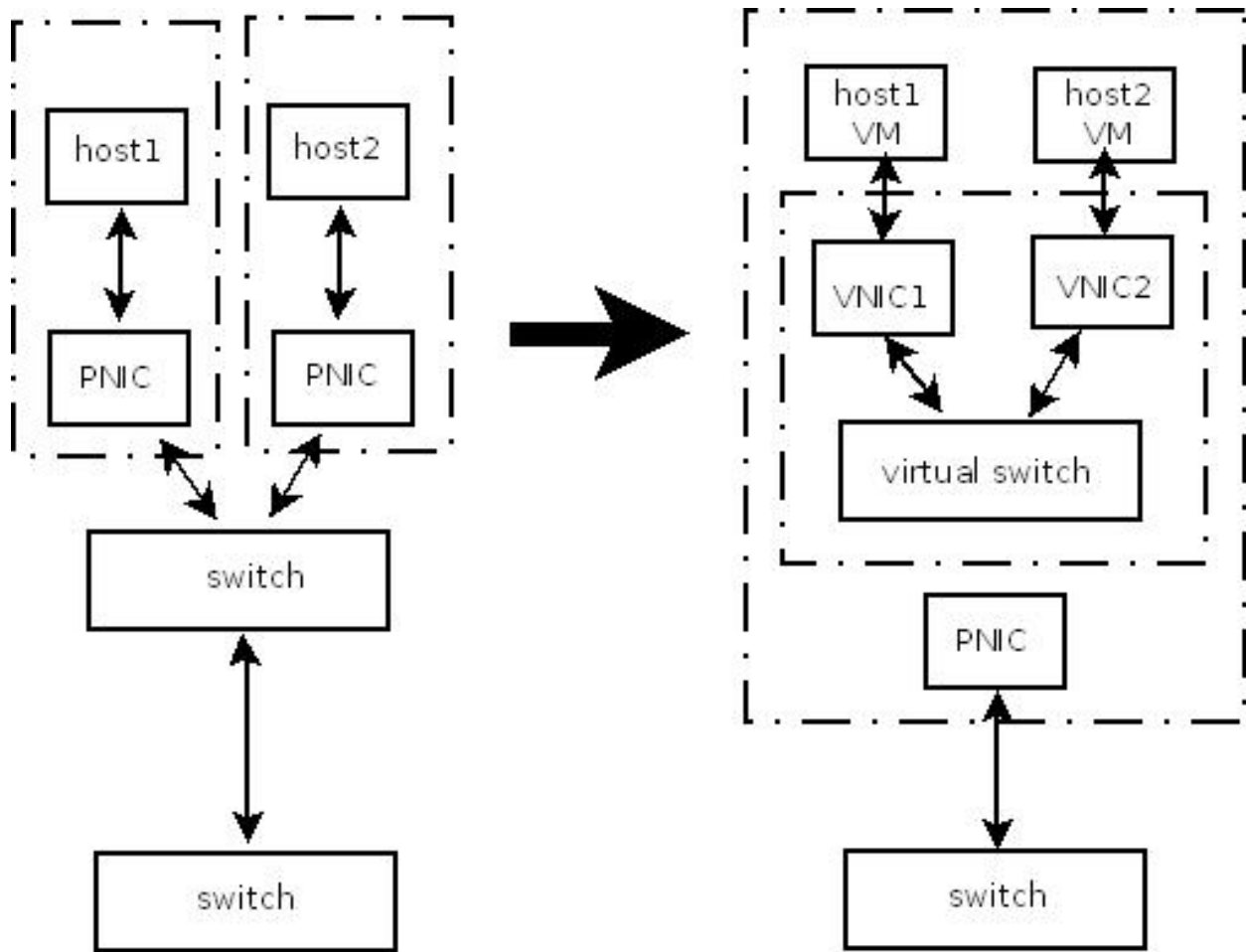


Figure 4.5: Mapping between physical and virtual network building elements

Semantics assured by virtual switches are the same as provided by physical switches:

1. VNICs created on top of the same NIC can send etherstub packets to each other
2. Broadcast packets received by the underlying NIC are distributed to every single VNIC that was defined on the top of this NIC
3. Broadcast packets sent by one of the VNICs is distributed to all VNICs defined on the top of the same NIC and to the NIC for further transmission as well
4. In terms of multicast network traffic multicast group membership is monitored and used for distributing packets to appropriate VNIC

Connectivity between VNICs is available only when they were defined on the top of the same NIC.

4.3.5 Crossbow elements

VNics

Virtual NICs (VNICs) each containing their own lane are the key element in crossbow architecture. There is no difference between NIC and VNIC in administration, as they are all treated as data links. Every VNIC has an assigned lane and flow classifier which classifies received packets by VNIC's MAC address and sometimes by the VLAN tag. If created with a VLAN tag, protocols like GVRP or MVRP may be used to register the VLAN tag with the physical switches too.

In terms of sharing bandwidth, Crossbow enables administrative control of bandwidth for every single VNIC. The bandwidth of the link is implemented by regulating the periodic intake of incoming packets per dedicated lane. The network stack allows only as many packets as it was assigned to specific VNIC. The lane picks more packets when the next period begins. In case of regulating the speed of transmitted bandwidth it is much easier as the network stack can either control the application that is generating the stream of packets or just drop the excessive amount of packets. These mechanisms are also used in flows QoS described and discussed later in this paper.

Etherstubs

As it was mentioned before, the MAC layer provides the virtual switching capabilities which allow VNICs to be created over existing physical NICs. In some cases, creating virtual networks without the use of a physical NIC is more welcomed than creating over physical NICs. In that case VNICs would be defined on the top of pseudo NICs. The Crossbow provides these kind of elements which are called Etherstubs. These components could be used instead of NICs during creation of VNICs.

Flows

Flows are additional instruments created to allow easier network traffic administration. They might be used in order to provide administer bandwidth resource control and priority for protocols, services, containers. Defined flow is a set of attributes based on Layer 3 and Layer 4 headers of the OSI/ISO model which are then used to identify protocol, service or virtual machine. Flows assigned to link must be independent therefore before adding new one its correctness is checked. Input and output packets are matched to flows in very efficient manner with minimal performance impact.

flowadm is the console command used to create, modify, remove or display network bandwidth and priority limits assigned to a particular link.

Examples

dladm is the admin command for crossbow datalinks elements management. Below a few examples of VNICs, Etherstubs management commands are presented and how bandwidth and priority values might be assigned to these elements.

1. `# dladm create-vnic vnic1 -l e1000g0` - creates new VNIC **vnic1** over existing NIC **e1000g0**
2. `# dladm create-etherstub ether00` - creates new Etherstub **ether00**
3. `# dladm show-linkprop vnic11` - lists all properties assigned to **vnic11** link
4. `# dladm set-linkprop -pmaxbw=1000 vnic11` - assigns 1Mbps bandwidth limit to **vnic11** link

5. # `dladm set-linkprop -ppriority=low vnic11` - assigns low priority to **vnic11** link

These were just basics, for more examples see **man dladm**

flowadm is the admin command for flow management. It might be used as follow:

1. # `flowadm show-flow -l e1000g0` - displays all flows assigned to link **e1000g0**
2. # `flowadm add-flow -l e1000g0 -a transport=udp udpflow` - creates new flow assigned to link **e1000g0** for all udp packets

To see more see **man flowadm**

4.4 Resource access control

Nowadays existing operating systems must provide mechanisms for response to the varying resource demands per workload which is an aggregation of processes of an application. By default resource management features are not used and system gives equal access to resources. When it is necessary it is possible to modify this default behaviour with respect to different workloads. This management allows you to:

1. Restrict access to specific resource
2. Offer resources to workloads on a preferential basis
3. Isolate workloads from each another

Resource is any part of computing system that may be modified in order to change application behaviour. Resource management enables more effective resource utilization and avoid wasting available ones due to load variability. Reserving additional capability during off-peak periods and effective sharing of resources definitely increases application performance.

Solaris Operating System introduced three types of resource management control mechanisms:

1. constraints - allows defining set of bounds on used resources for a workload
2. partitioning - enables binding subset of system's available resources to specific workload
3. scheduling - involves predictable algorithm making sequence of allocation decisions at specific intervals

Summary

Solaris 10 OS seems to be ideal cross-platform choice for customers dealing with management of high level services, complex system administration and high costs. It is the only open operating system which has proven results running from every critical enterprise databases to high performance Web farms that is why Solaris OS is becoming strategic platform for today's constantly growing demands towards operating systems.

Chapter 5

The system architecture

Chapter overview

The Domain model and data flows section describes the transformations performed by the system's components in order to instantiate/deploy an object model. These include simple one-node instantiation as well as more complex multi-node instantiations.

5.1 High-level design

5.2 System components and their responsibilities

5.2.1 Assigner

5.2.2 Supervisor

5.2.3 Worker

5.3 Crossbow resources instrumentation

5.4 Domain model and data flows

Summary

Chapter 6

Implementation

Chapter overview

6.1 Implementation environment

6.2 Domain model transformation details

6.3 Low-level functions access

6.4 Building and running the platform

Summary

Chapter 7

Case Study

Chapter overview

7.1 Multimedia server

7.1.1 Scenario description

- similar to DiffServ (traffic classes, selectors, filters, priority, queuing)
- DiffServ doesn't specify anything virtual
- DSS and adaptive codecs
- 2 classes: VOD + streaming
- _ unicast _ vs multicast streaming
- access rules for resources of different quality
- enabling QoS for defined classes
- priorities + limiting the bandwidth (per user)!
- 3 users: 2 streaming, 1 VOD

7.1.2 Resource access requirements

7.1.3 Providing tunable and scalable virtual infrastructure

Summary

Chapter 8

Summary

Chapter overview

Bibliography [?] test.

8.1 Conclusions

8.2 Achieved goals

8.3 Further work

In terms of the future work there are many many improvements that might be implemented. Probably the largest component we'd planned to implement was automatic resource assigner, which would run and perform automatic assigning resources to nodes that run under least load. This assigner with attached rule based system should gather data about the load on each node and based on that should decide what and where instantiate. What we have managed to complete is manual assigner, where you have to select on which node you would like to have your virtual resources created.