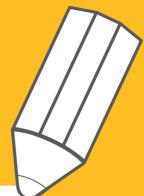


I/O Acceleration for VM

포티넷 코리아 이상훈

**DATACENTER
ARCHITECT
GROUP**





(우진) 이상훈
Fortinet SE

woojinlsh@gmail.com



KRDAG
(Korea Datacenter Architect)
Amway Korea
(Technical Architect, 인프라총괄)
Infnis
(Solution Launch Support Engineer)

Agenda

- 
1. Emulation & Hypervising
 2. I/O Acceleration
 1. Onloading
 2. Offloading
 3. I/O Virtualization Technologies
 4. 마무으리



오늘은
달라는데
걸로~!

본지의 편집방향과 일치하지 않을 수도 있습니다.

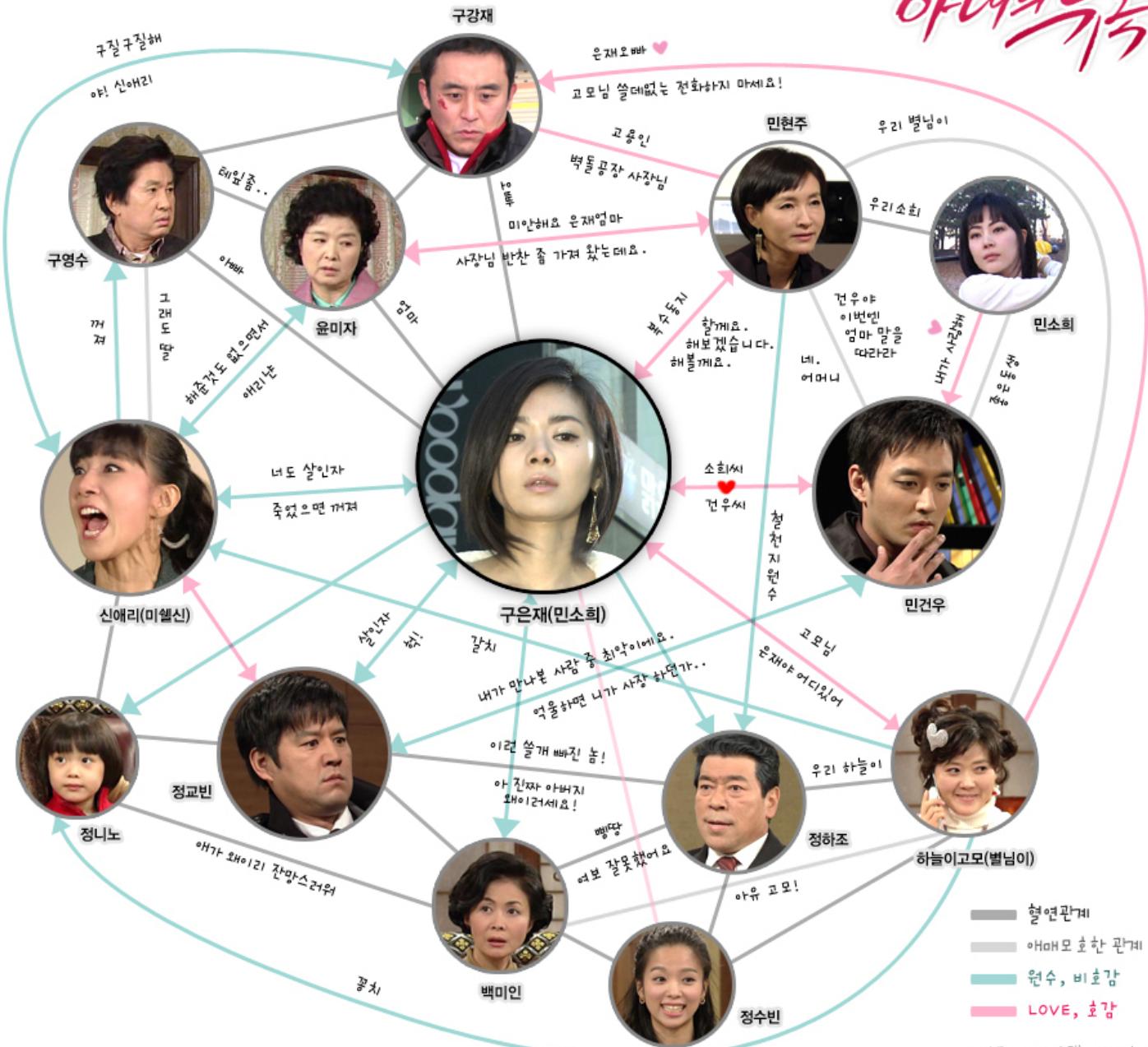
I/O Virtualization

그까이꺼

별거 있나?

인물관계도 ver.01
아내의 유품

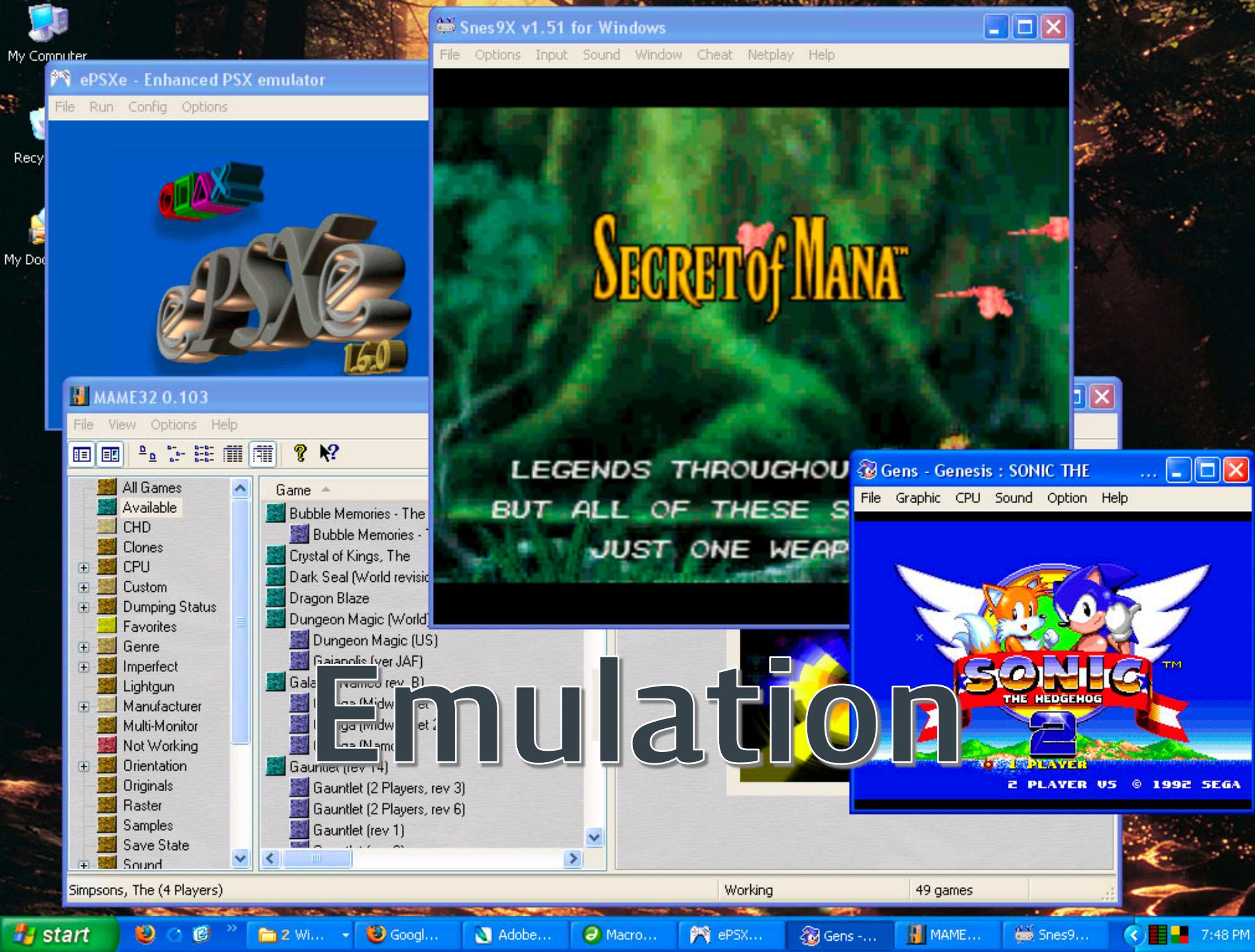
인물관계도 ver.0



그 게 좀 복 잡

1 / 3

[HTTP://GALL.DCINSIDE.COM/TEMPTATION](http://GALL.DCINSIDE.COM/TEMPTATION)

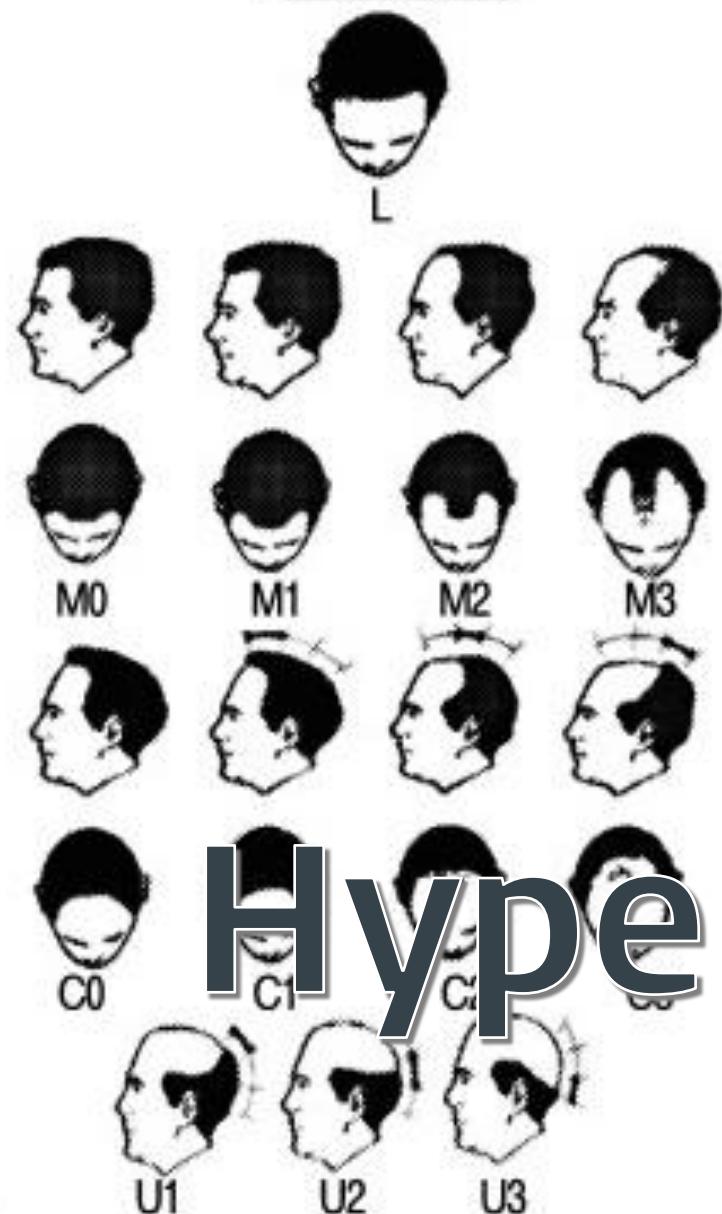


무에서
유를 창조

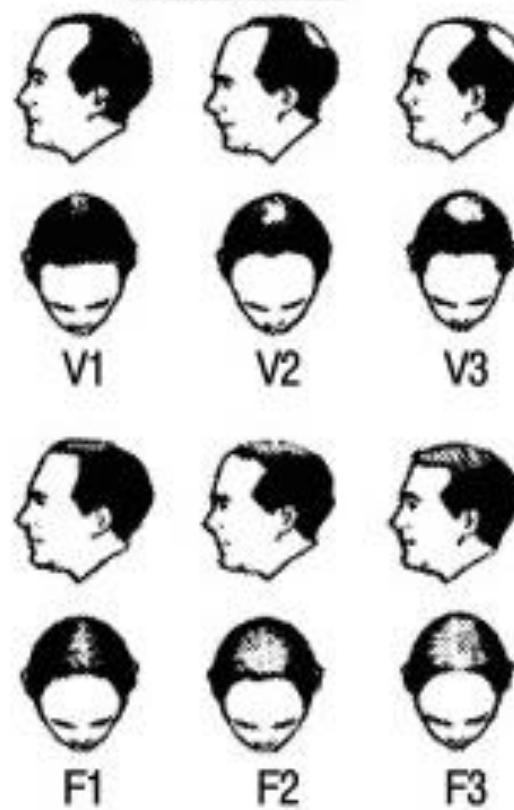


■ BASP에 따른 탈모 유형 분류

BA형



SP형



최종형

BA + SP형

Hypervisor



GIRLS



퍼프스힐

GUYS



하이힐



스틸레토힐



하이힐



웨지힐



하이힐

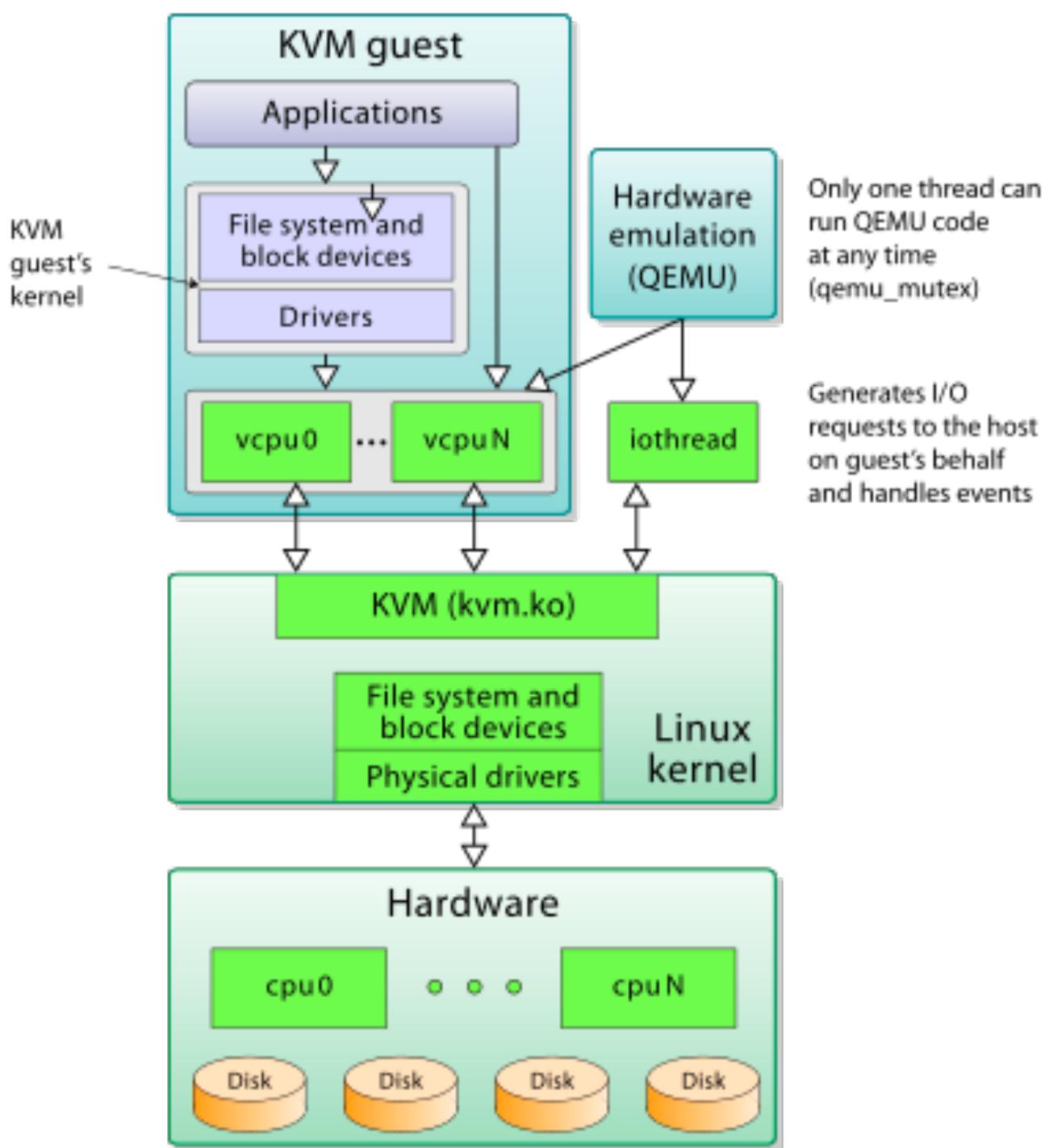


클래포힐



하이힐

Hypervisor



Parent Partition

Child Partition

VMI Provider

Virtual Machine
Management Service

VM
Worker
Processes

Applications

Windows
Kernel

Virtualization
Service Provider
(VSP)

Device
Drivers

Virtualization
Service
Consumer(VSC)

Windows
Kernel

VMBus

VMBus

User Mode

"Ring 3"

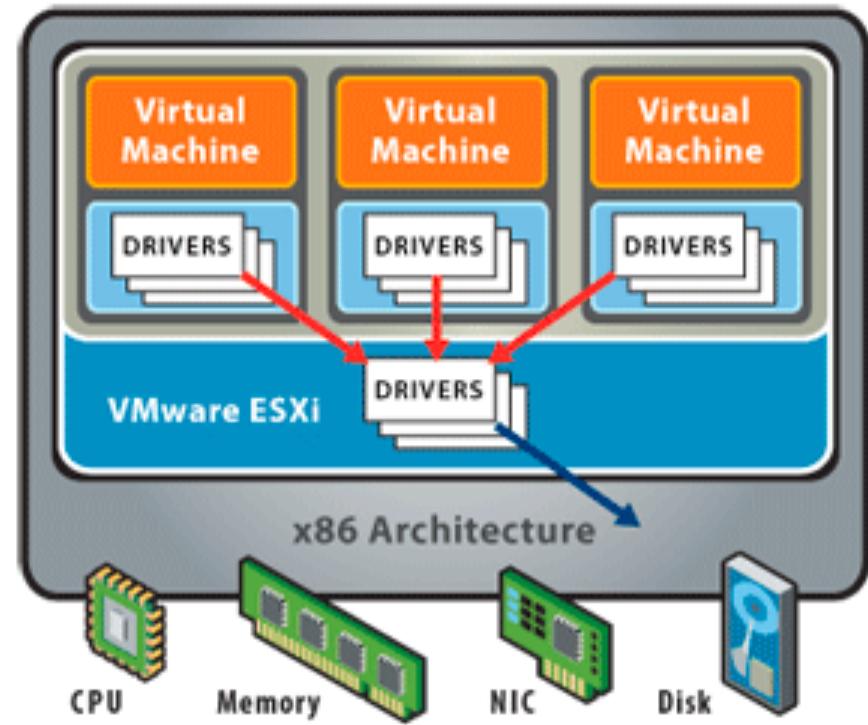
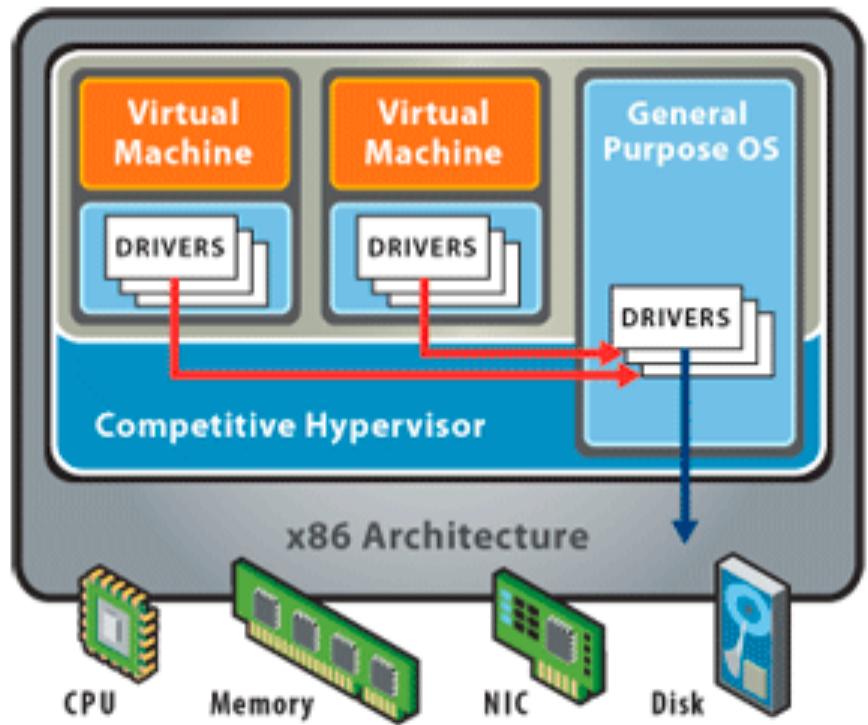
Kernel Mode

"Ring 0"

Hypervisor

"Ring -1"

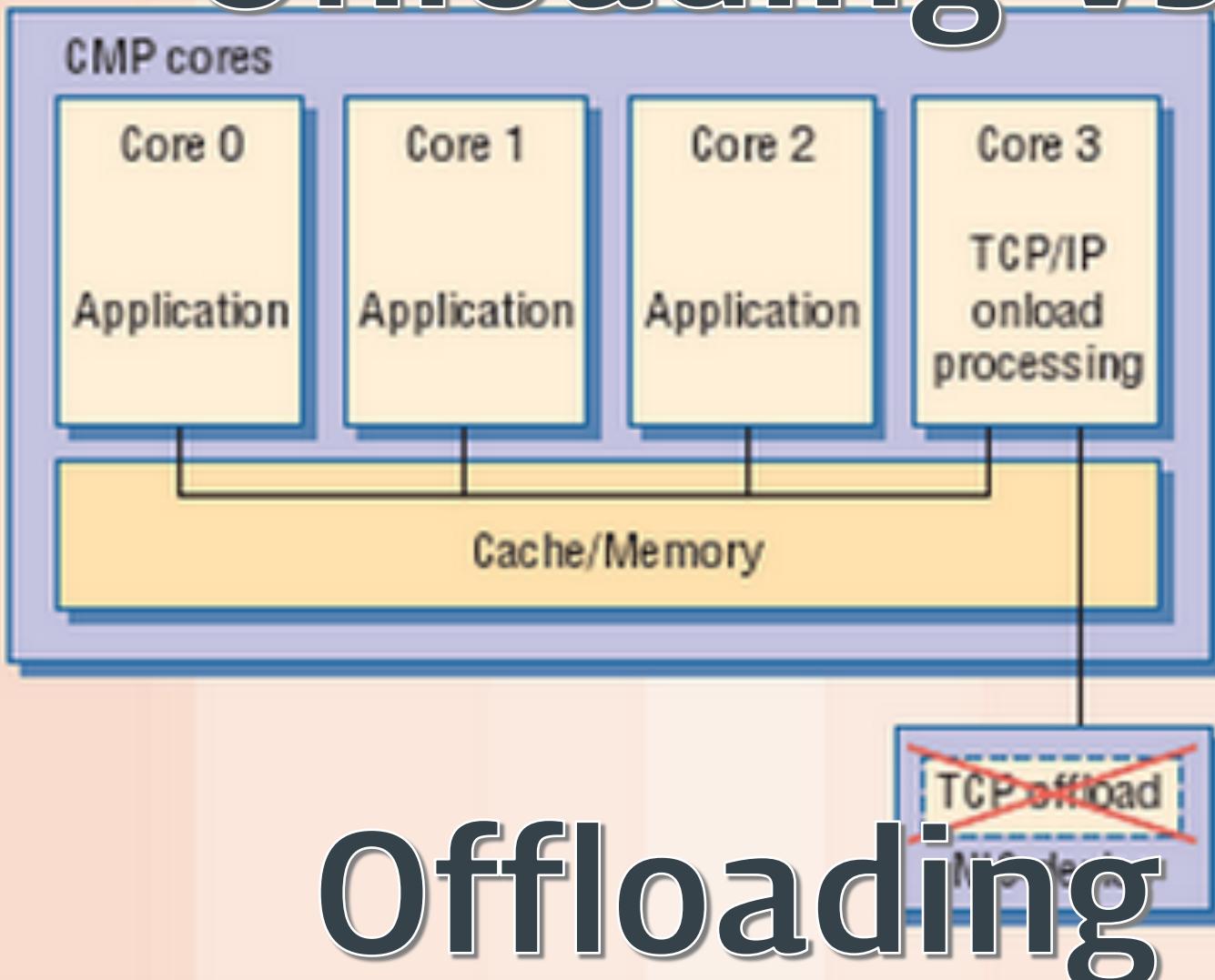
Hardware



I/O Acceleration

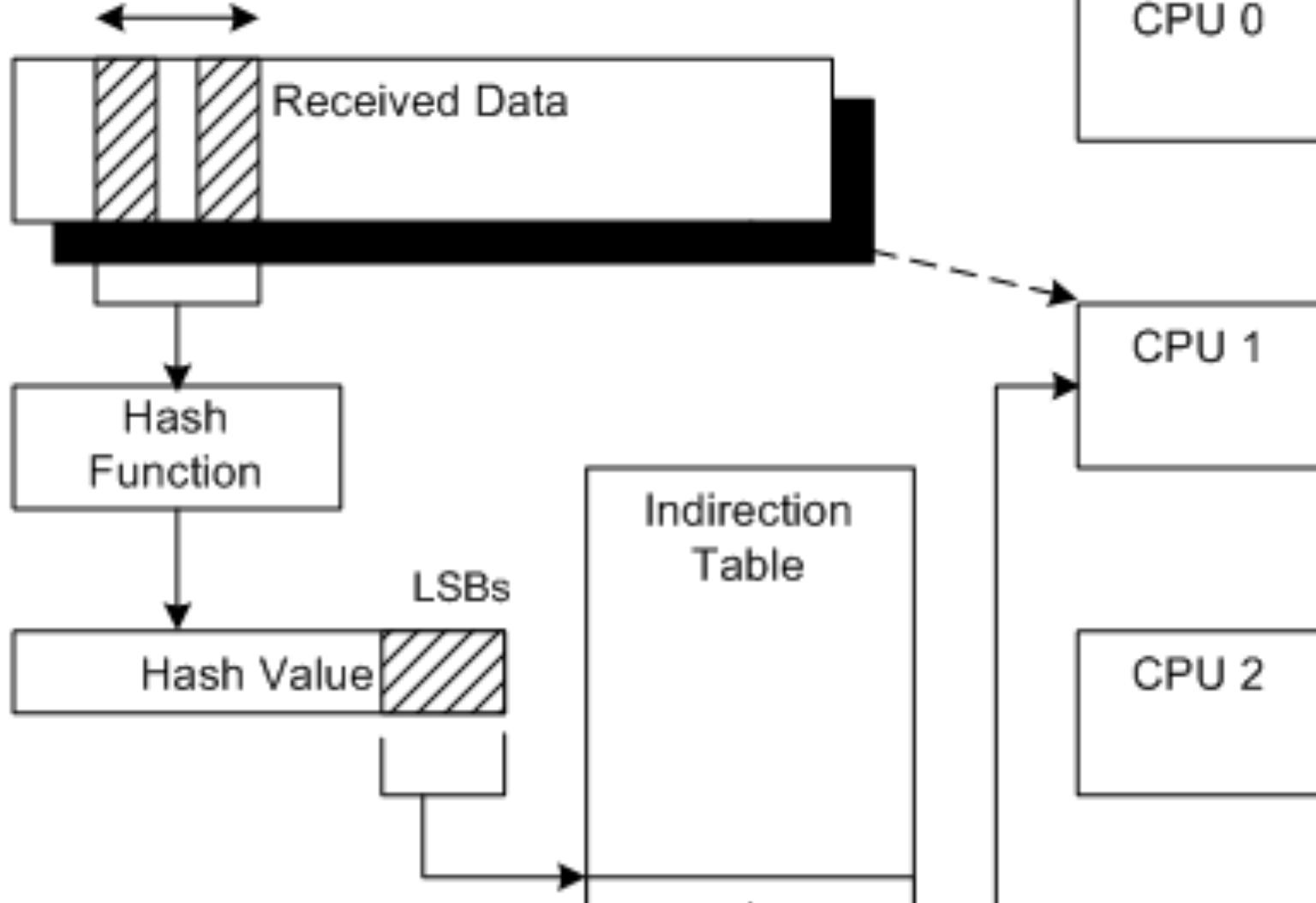


Onloading vs.



Onloading

Hash type specified



Receive Slide Scaling

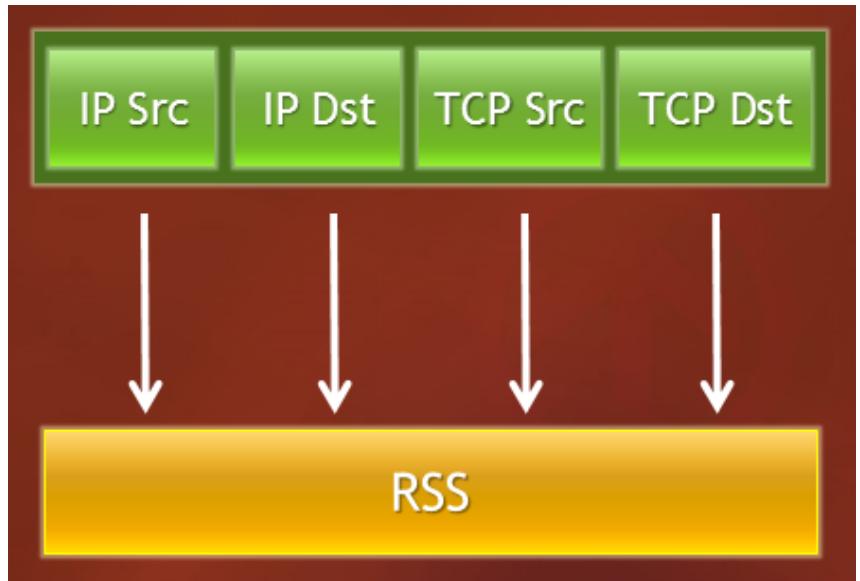


한놈만 패

A photograph showing a large group of mongooses gathered around a cobra. The mongooses are in various positions, some facing the camera and others looking towards the snake. The cobra is coiled on the ground, its hood expanded, and its head pointing towards the center of the group. The scene captures a moment of predator-prey interaction or a social gathering.

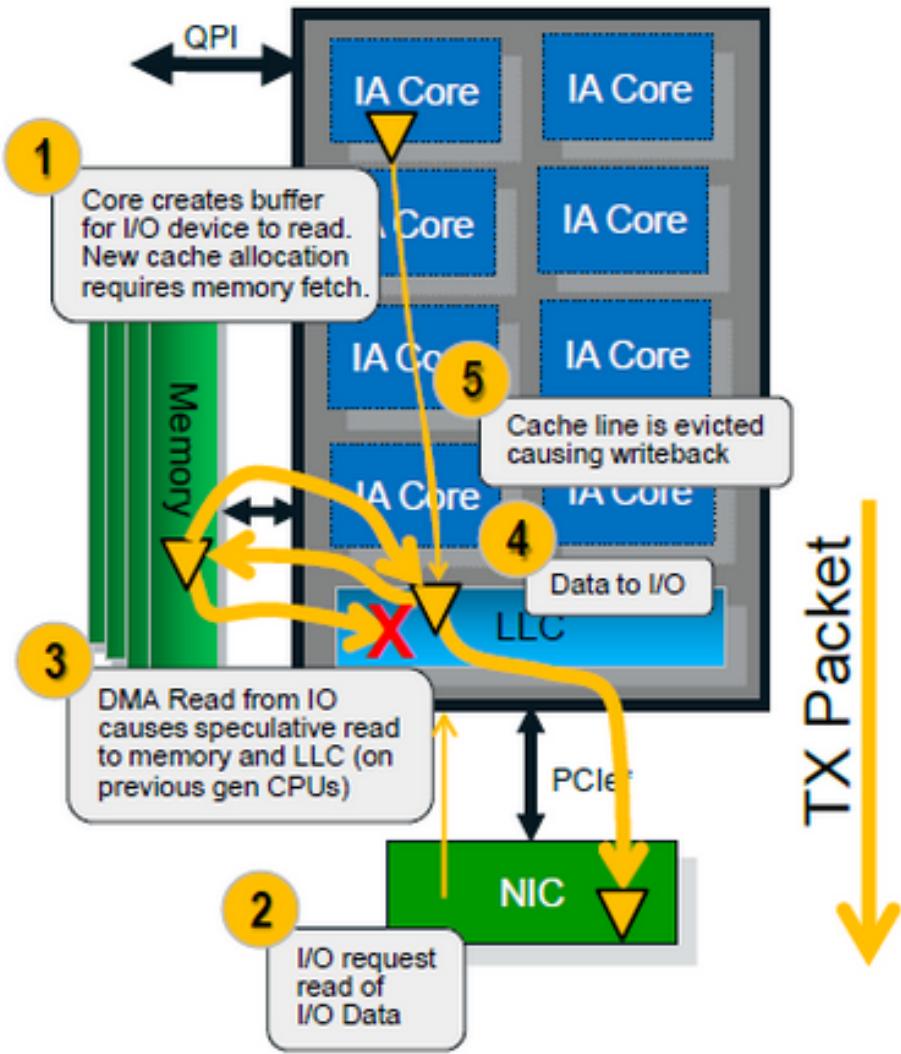
17대 1

Receive Slide Scaling (RSS)

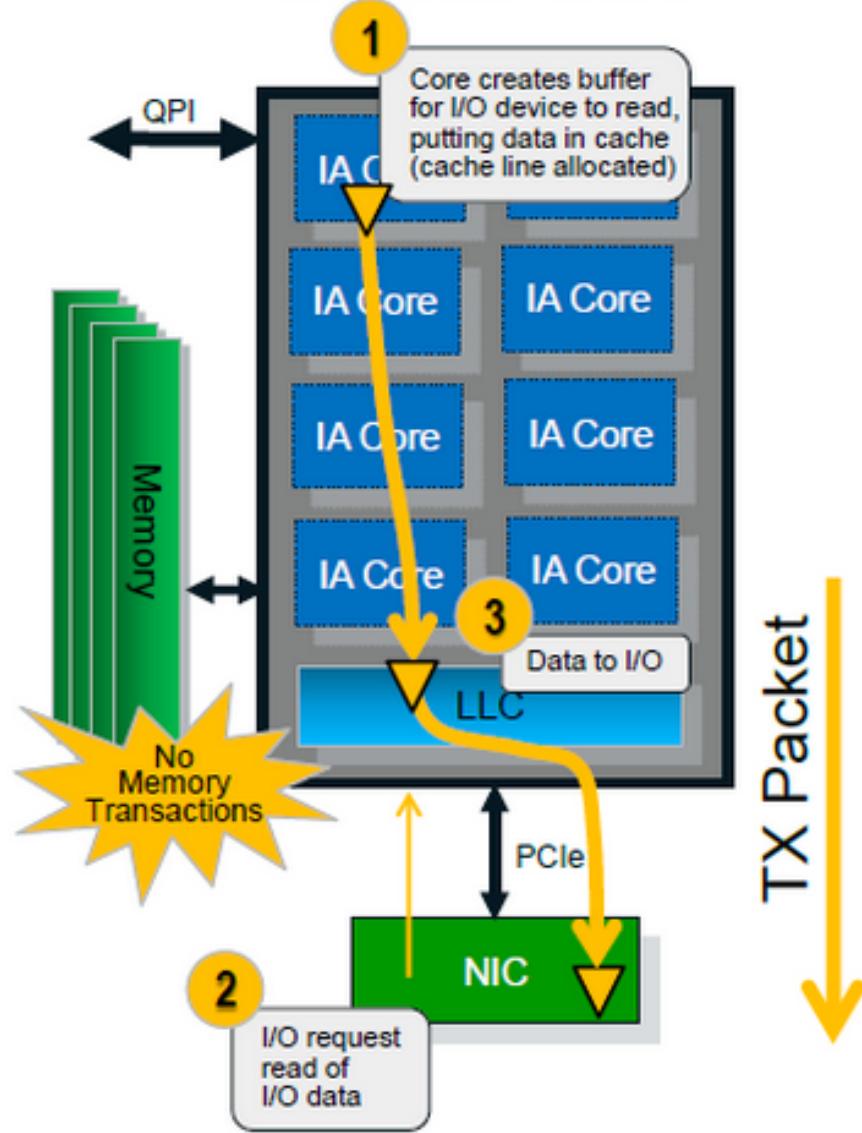


- Multicore 환경에서도 CPU 간 처리시간 차이에 따른 Packet 순서 변경으로 인한 성능 저하 문제 발생
- Flow 단위 (src, dstIP, src, dstPort)로 Hashing 하여 그 값에 따라 CPU 처리
- (HASH TYPE = IPv4, IPv4+TCP, IPv6, IPv6_TCP, IPv6_EX, IPv6_EX_TCP)
- Windows 2012의 경우 64개 이상의 Core, NUMA Architecture 지원

Without Intel® DDIO



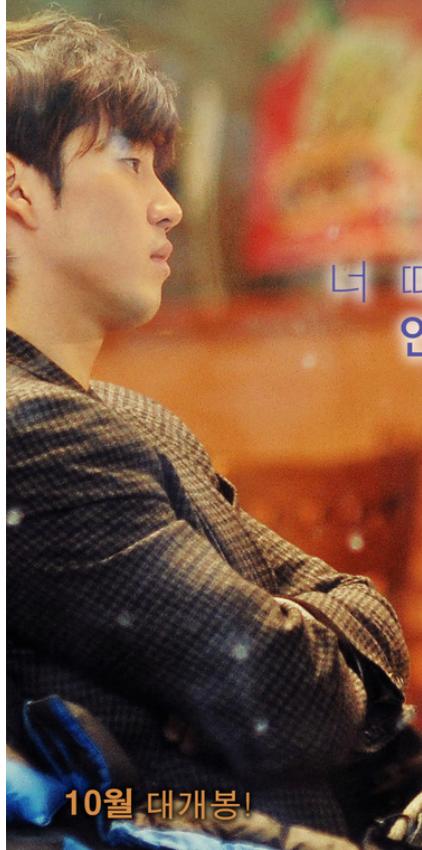
With Intel® DDIO



Come, Closer

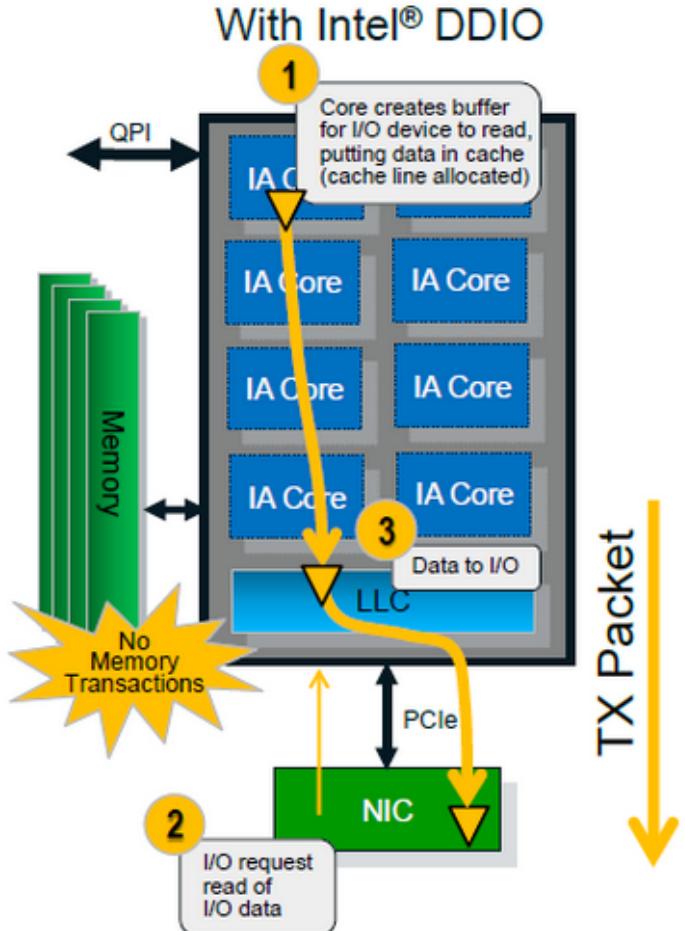
고장난 사랑에 관한 다섯 가지 증상

너 때문에 나...
연애불구야.



10월 대개봉!

Intel Data Direct I/O (DDIO)



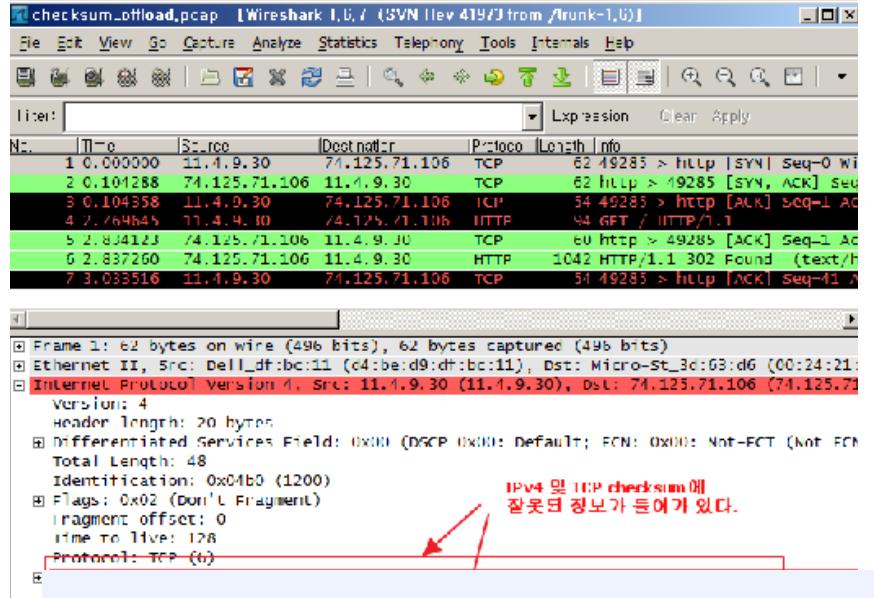
- I/O 시 Main Memory 를 사용하지 않고 CPU 의 Cache Memory를 사용함 (SRAM)
- Intel Sand Bridge Chipset 부터 지원
- LAN Card 상관 없음
- Ethernet 뿐 아니라 Infiniband, Fibre Channel, RAID 기술 까지 모두 사용 가능
- 현존하는 I/O Acceleration 중 가장 빠른 기술 (다른 가속기술과 중복 사용가능 하기 때문)

Offloading

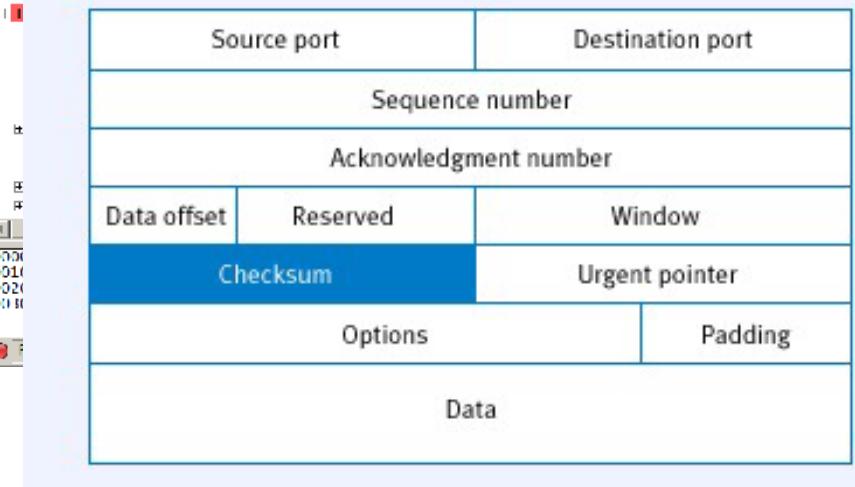


**Checksum Offload
IPSEC Offload
Segmentation Offload
Chimney**

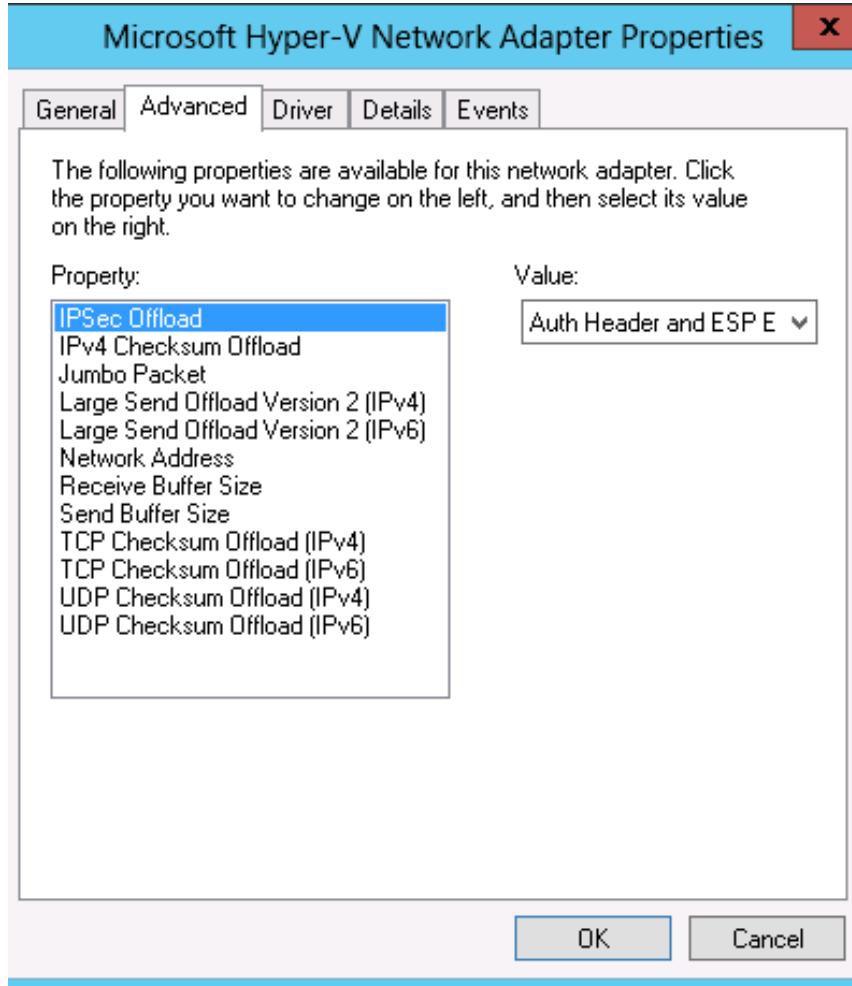
Checksum Offload



- 가장 기본적인 Offload 기능
- TCP Header 의 Checksum 을 NIC 에서 생성
- Checksum 로드를 NIC 에 할당하기 때문에 CPU 부하가 줄어듬
- Wireshark 에서 TCP Checksum 오류로 나타날 경우 Checksum 기능 제거 필요

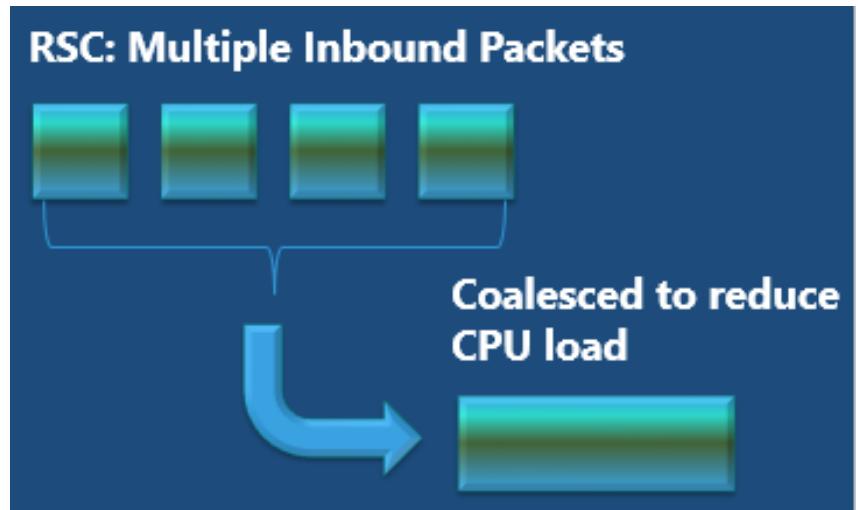


IPSEC Offload



- Server NIC 에 Encryption ASIC 이 탑재된 경우 (Wireless NIC 등)
- CAVIUM Encryption Card 처럼 별도의 가속기를 PCI 에 장착한 경우
- CPU 내부에 가속 ASIC 이 장착된 경우 (IBM PowerPC, Oracle CPU, AMD Load 에 존재)
- 대부분 IPSEC Offload 의 경우 ESP, AH 패킷에 대한 가속만 지원하고 IKE에 대한 가속은 지원하지 못함

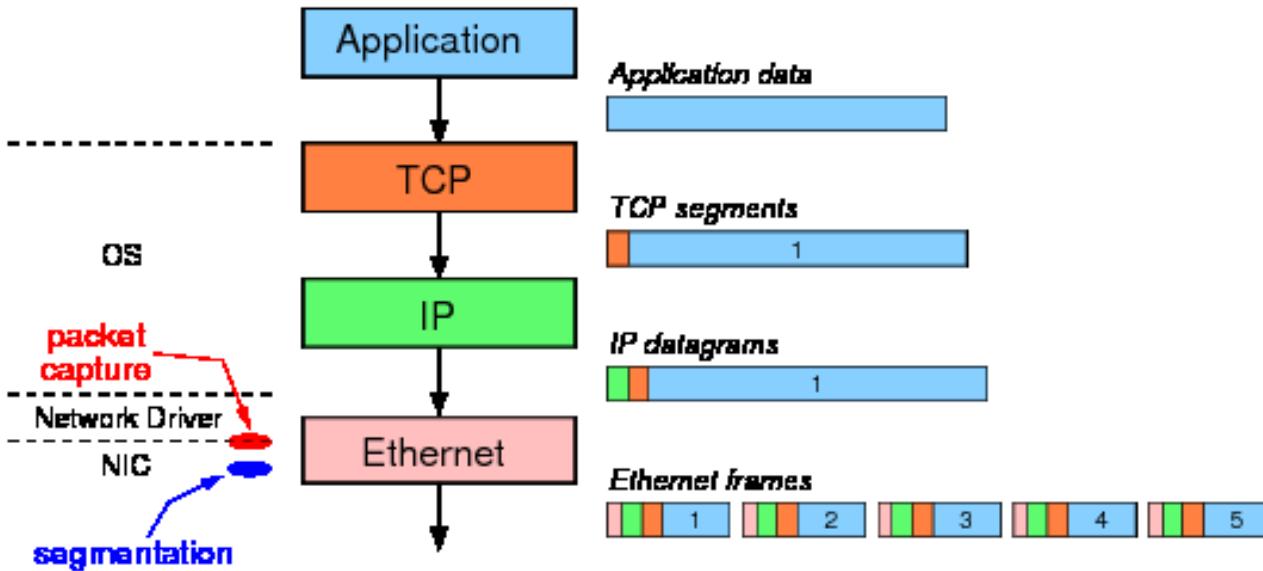
Segmentation Offload – Receive Segment Coalescing (RSC)



Generic Segmentation Offload (GSO)
Large Send Offload (LSO)
Generic Receive Offload (GRO)
Large Receive Offload (LRO)
TCP Segmentation Offload (TSO)

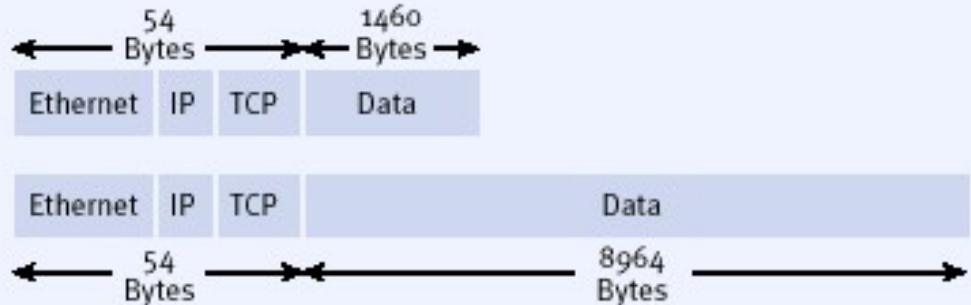
- Receive Segmentation Coalescing
 - » 일반적으로 Packet 당 Interrupt 를 처리 하지만 RSC 가 Enable 되어 있을 경우 Flow 단위로 한꺼번에 Packet 을 처리
 - » 대용량 Buffer 를 지원할 수 있어야 함
 - » 10G NIC 은 기본적으로 RSC를 지원
 - » 최대 64KB Payload 를 한꺼번에 처리할 수 있음 (v2는 256KB)
 - » Windows 2012에 소개
 - » Linux Kernel 2.6 부터 정식 지원 (2.4.x 소개)

Segmentation Offload - Large Send Offload v2 (LSOv2)

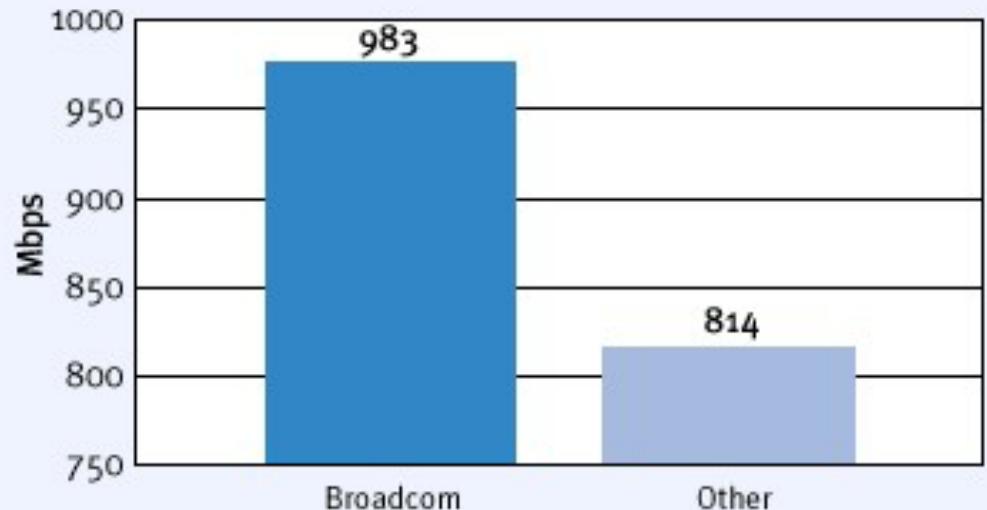


TCP	66	ndl-aas > 37295 [ACK]	TCP	66	ndl-aas > 37293 [ACK]
TCP	2762	[TCP segment of a reas	TCP	1414	[TCP segment of a reas
TCP	66	37295 > ndl-aas [ACK]	TCP	66	37293 > ndl-aas [ACK]
TCP	4017	[TCP segment of a reas	TCP	1414	[TCP segment of a reas
TCP	66	37295 > ndl-aas [ACK]	TCP	66	37293 > ndl-aas [ACK]
TCP	3538	[TCP	TCP	1414	[TCP
TCP	66	3729	TCP	66	3729
TCP	3538	[TCP	TCP	1414	[TCP
TCP	66	3729	TCP	66	3729
TCP	98	[TCP segment of a reas	TCP	1414	[TCP segment of a reas
TCP	66	37295 > ndl-aas [ACK]	TCP	66	37293 > ndl-aas [ACK]
TCP	16242	[TCP segment of a reas	TCP	1414	[TCP segment of a reas

Jumbo Frame



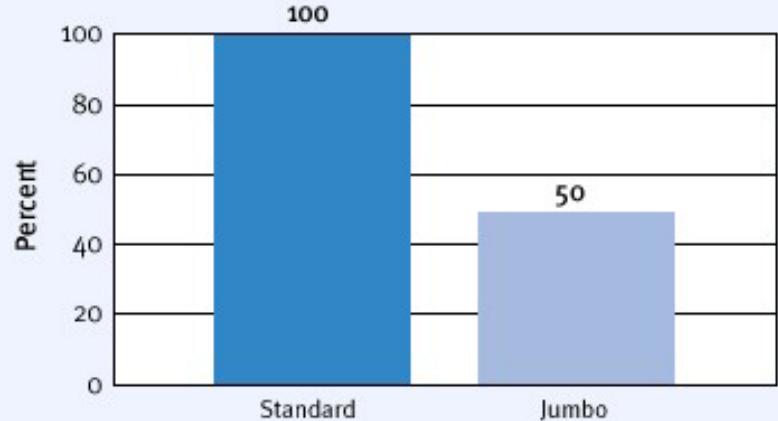
Performance



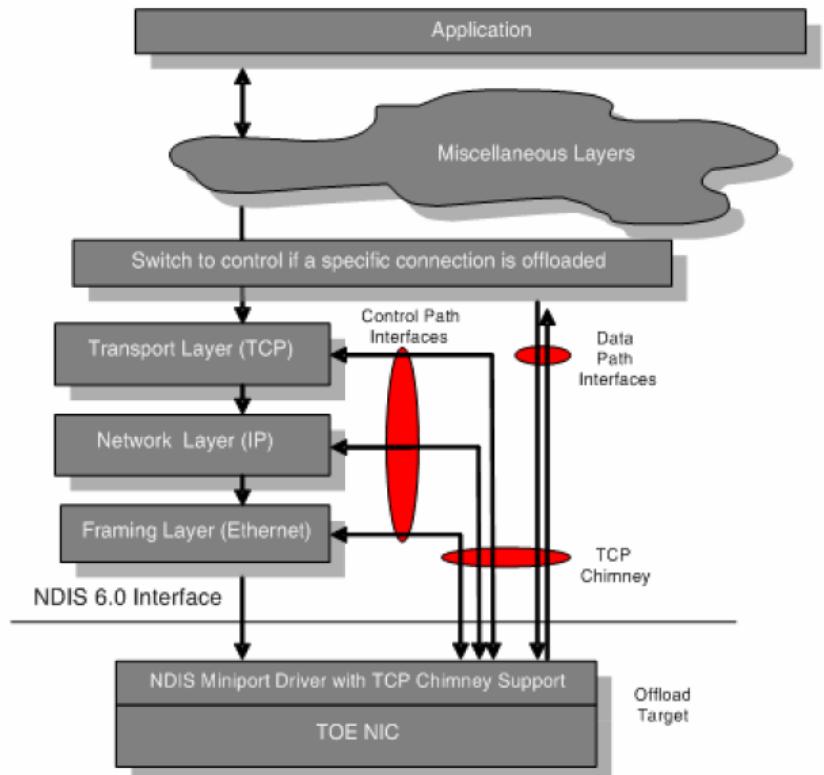
Throughput



CPU utilization



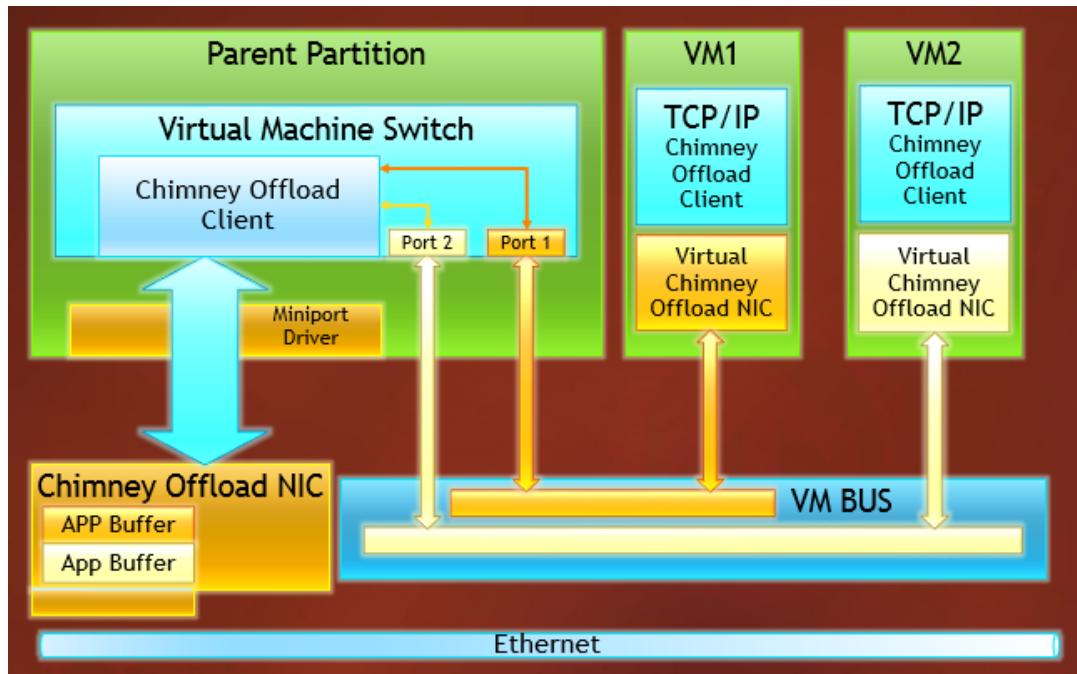
TCP Chimney Offload



- Application에서 OS의 TCP/IP Stack을 거치지 않고 Direct로 TOC NIC에 직접 Connection
- 초기 NIC Driver bug로 인하여 강제로 기능을 사용하지 않는 것을 추천하였으나 최근 안정화됨(Driver 영향이 큼)
- Alacritech사 특허

Merry!

TCP Chimney Offload on Hypervisor



- Hypervisor에서도 Chimney Offload 기능을 사용할 수 있음
- Virtual Chimney Offload NIC을 통해 서 복제됨 (VM BUS를 이용하므로 copy 부하는 적용)

Considerations

Performance Metric	Loopback Fast Path	Registered I/O (RIO)	Large Send Offload (LSO)	Receive Segmentation Offload (RSC)	Receive Side Scaling (RSS)	Virtual Machine Queues (VMQ)	Remote DMA (RDMA)	Single Root I/O Virtualization (SR-IOV)
Lower End-to-End Latency	X	X					X	X
Higher Scalability		X			X	X		
Higher Throughput	X	X	X	X	X	X	X	X
Lower Path Length	X	X	X	X			X	X
Lower Variability		X						

I/O Virtualization Technologies

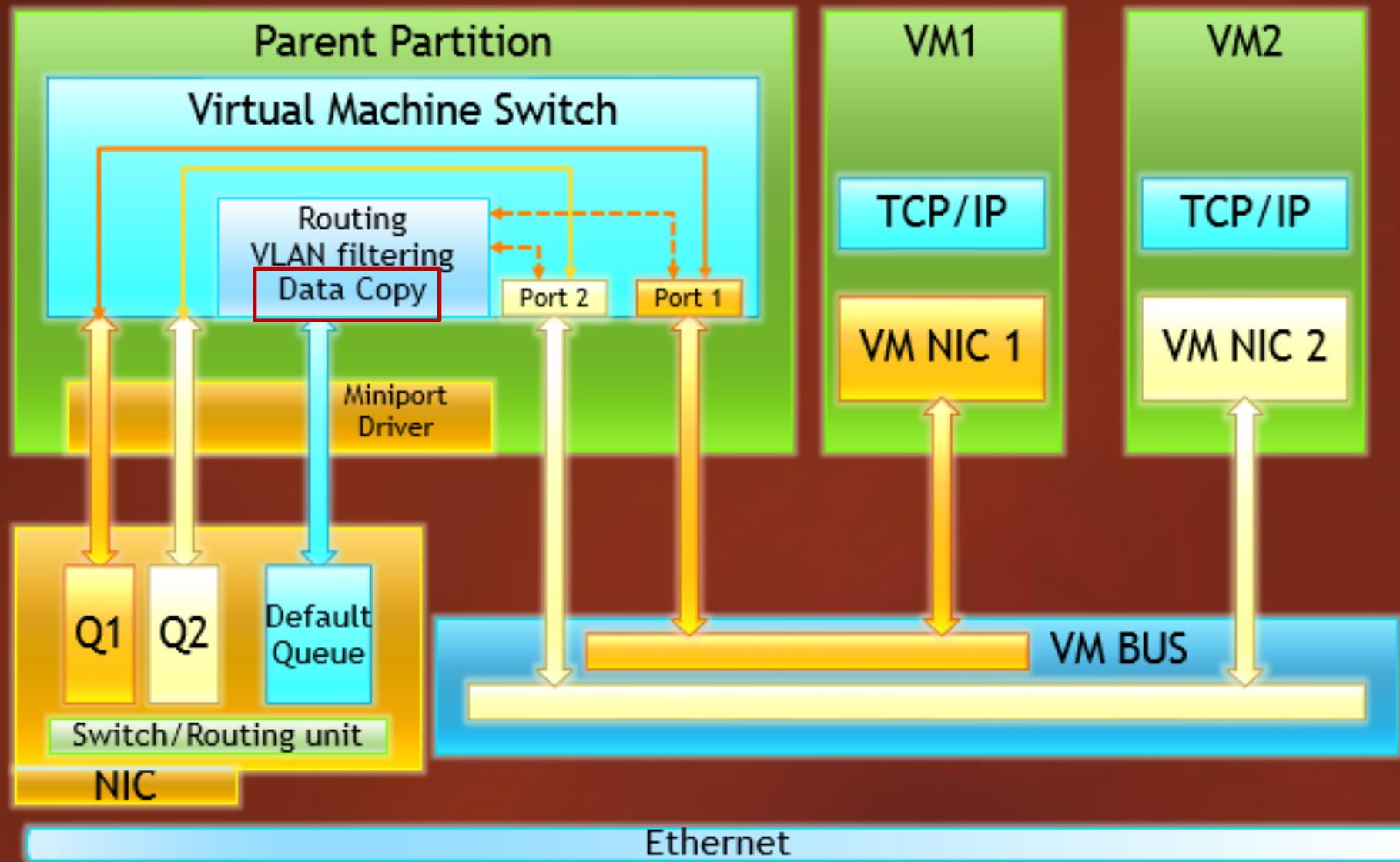


Virtual Machine Queue (VMQ)

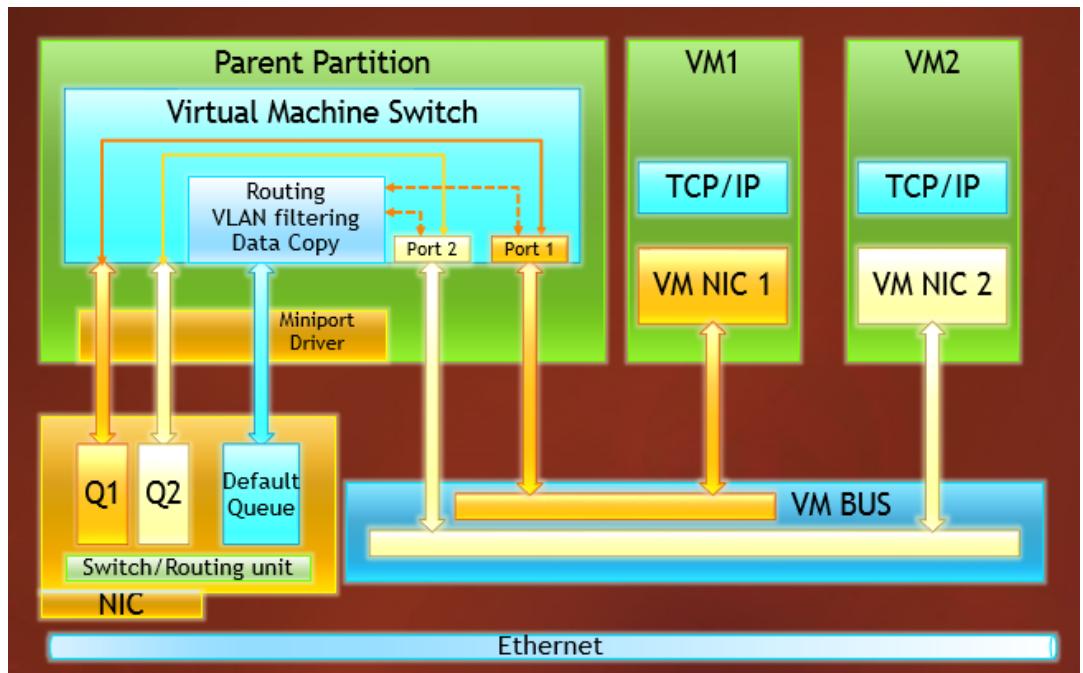
SR-IOV

MR-IOV

NIC Partitioning

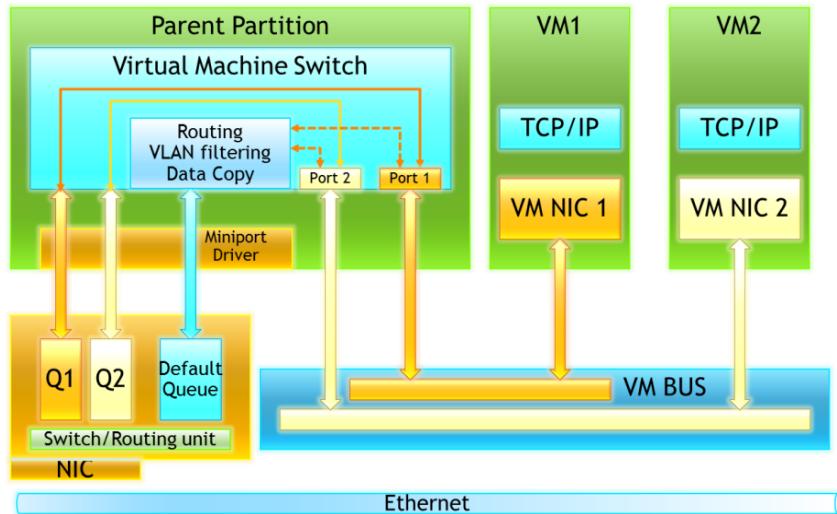
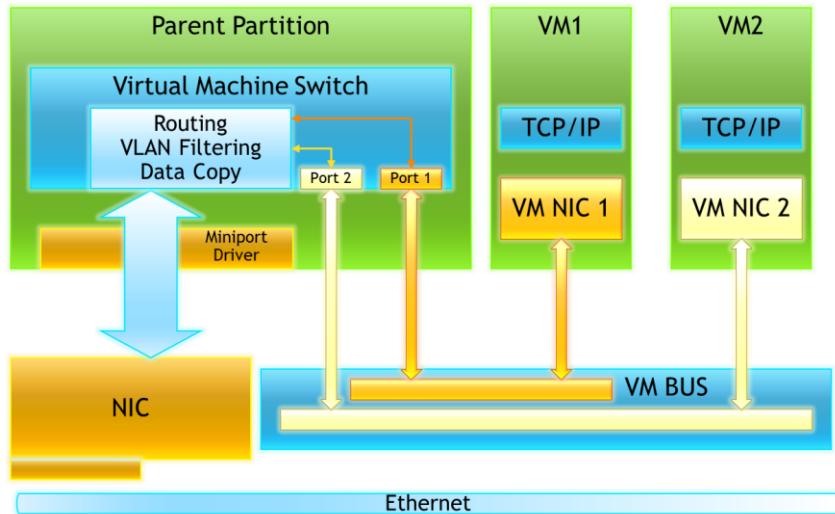


Virtual Machine Queue (VMQ)



- NIC에 Virtual Machine Queue를 할당하여 GuestVM과 구분되어 통신하도록 지원하는 기능
- 각종 Offload 기능을 그대로 사용 가능
- Parent Partition을 거쳐야 하므로 (Emulation) Direct I/O 기능을 지원하는 Virtual NIC 선택 필요
- Parent Partition 통과 여부가 SR-IOV, MR-IOV와 차별화 되는 점

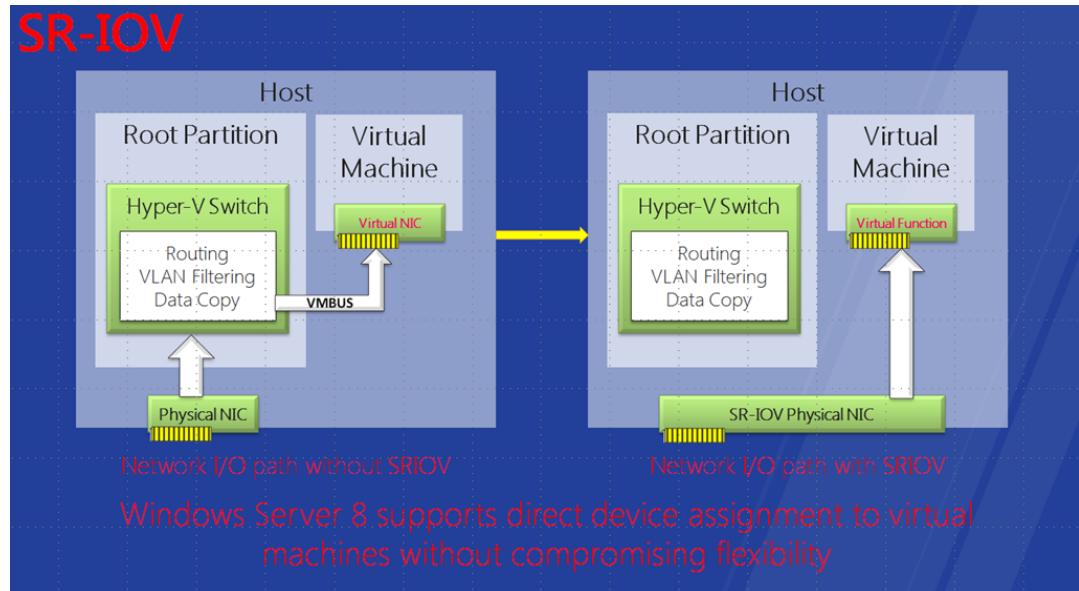
Virtual Machine Queue (VMQ)



without VMQ

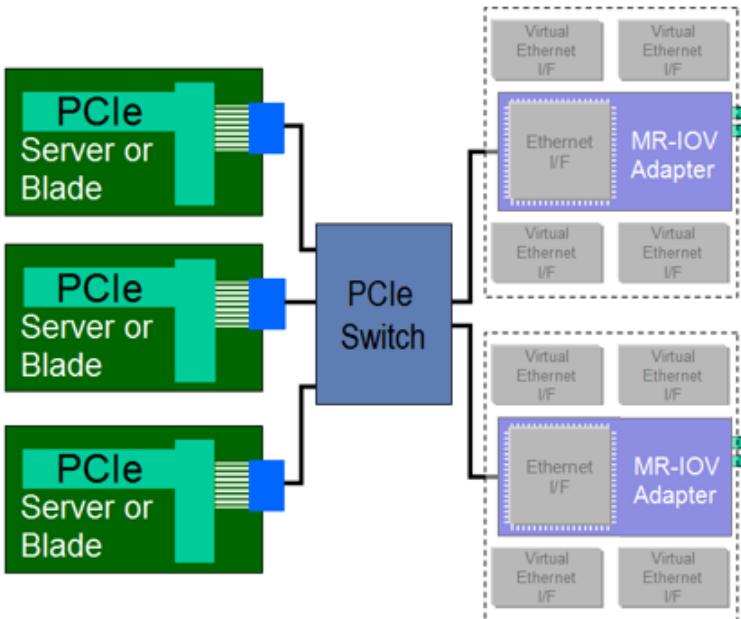
with VMQ

SR-IOV (Single Root IO Virtualization)



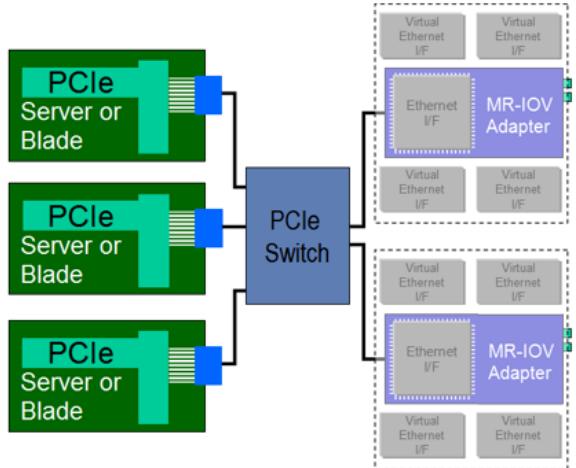
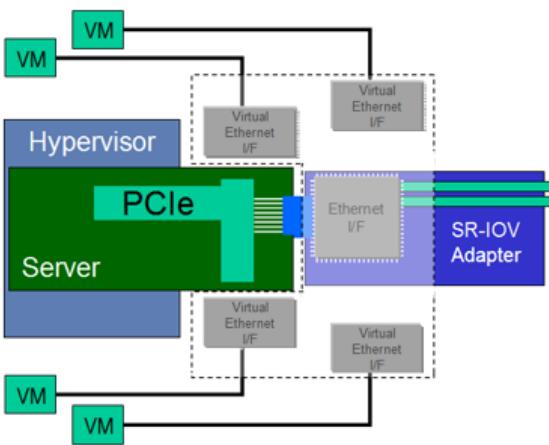
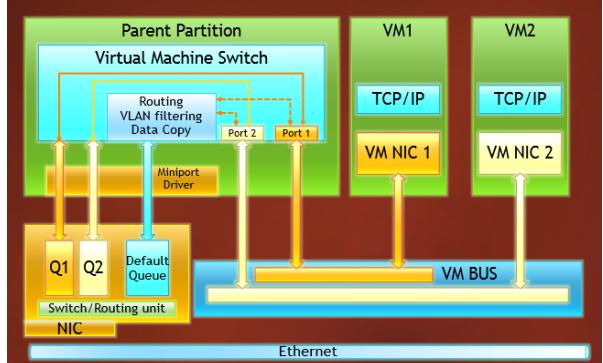
- Hypervisor 를 통하지 않고 I/O 장치에 Guest OS 가 Direct Access
- 기존 Hypervisor 에서 I/O 는 Emulation 으로 인한 속도 저하 문제가 있음
- Guest VM 에서는 SR-IOV Virtual NIC 으로 인식 됨

MR-IOV (Multi Root IO Virtualization)



- Hypervisor 를 통하지 않고 I/O 장치에 Guest OS 가 Direct Access (SR-IOV 와 동일)
- 여러개의 Server 혹은 Blade 에서 PCIe Switch 를 이용하여 NIC에 접근
- PCIe 자체를 가상화 하므로 Guest VM 에서는 별도의 NIC 으로 보임

VMQ vs. SR-IOV vs. MR-IOV

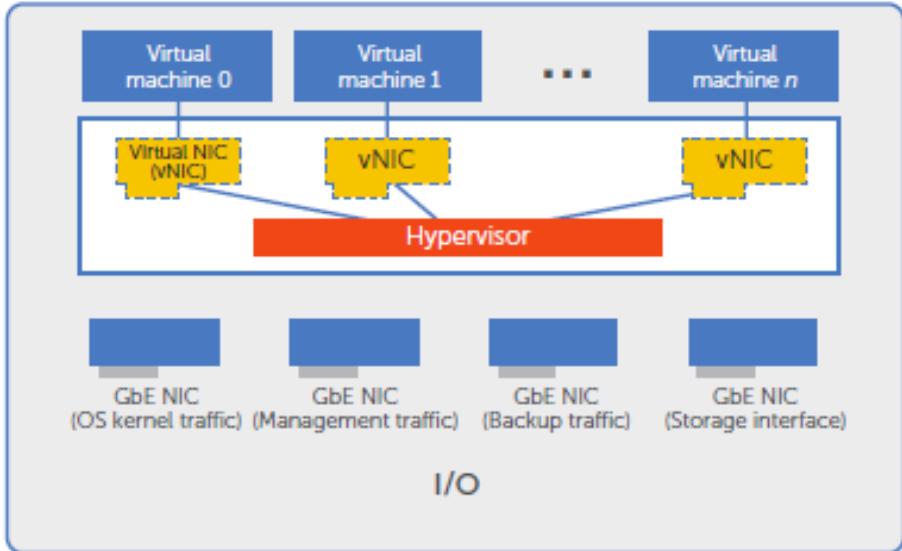


- **VMQ & Non-I/O Virtualization**
 - 장치메모리 ↔ PCI 버스(R) ↔ 하이퍼바이저 ↔ Guest VM 메모리
- **SR-IOV**
 - 장치메모리 ↔ 하나의 PCI 버스(R) ↔ Guest VM 메모리
- **MR-IOV**
 - 장치메모리 ↔ 여러 개의 PCI 버스(R) ↔ Guest VM 메모리

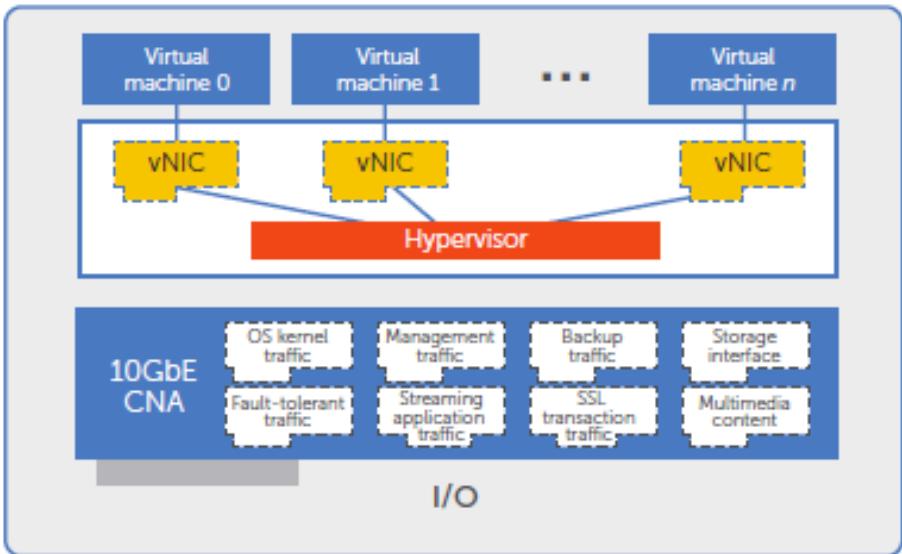


NIC Partitioning and etc. (Qlogic VMFlex, Broadcom)

Without NPAR



With NPAR



마
국
으
리



**CPU, Memory
Hypervisor**

기타 장치 Emulation

달라요

다르다

Different

다르다니까

다를껄

Onloading
Offloading

가속 기술 들들

VMQ
SR-IOV
MR-IOV

