

NUMA Aware I/O in Virtualized Systems

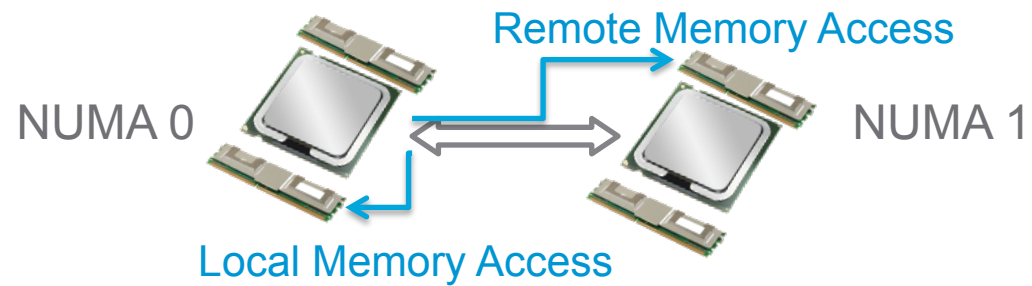
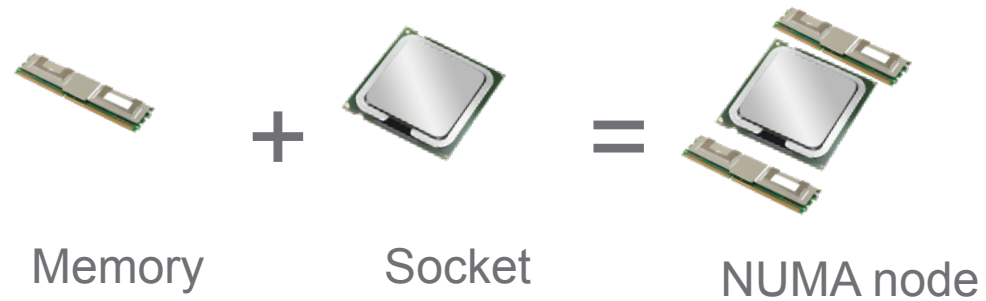
Rishi Mehta, Zach Shen,
Amitabha Banerjee



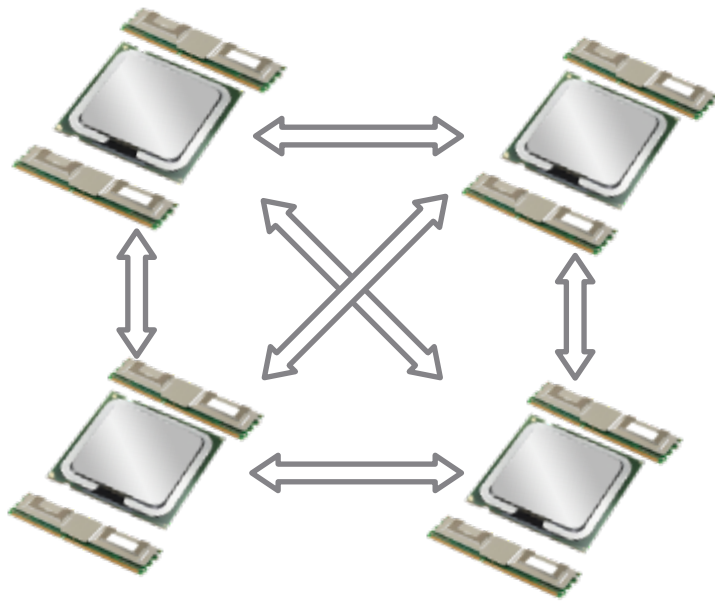
vmware®

© 2014 VMware Inc. All rights reserved.

Non-Uniform Memory Access



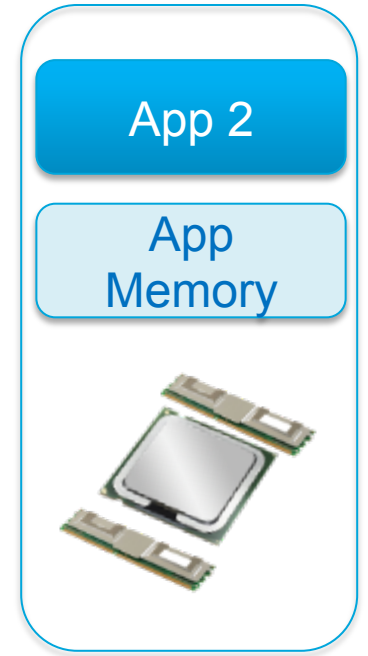
Non-Uniform Memory Access



Modern Server Architecture



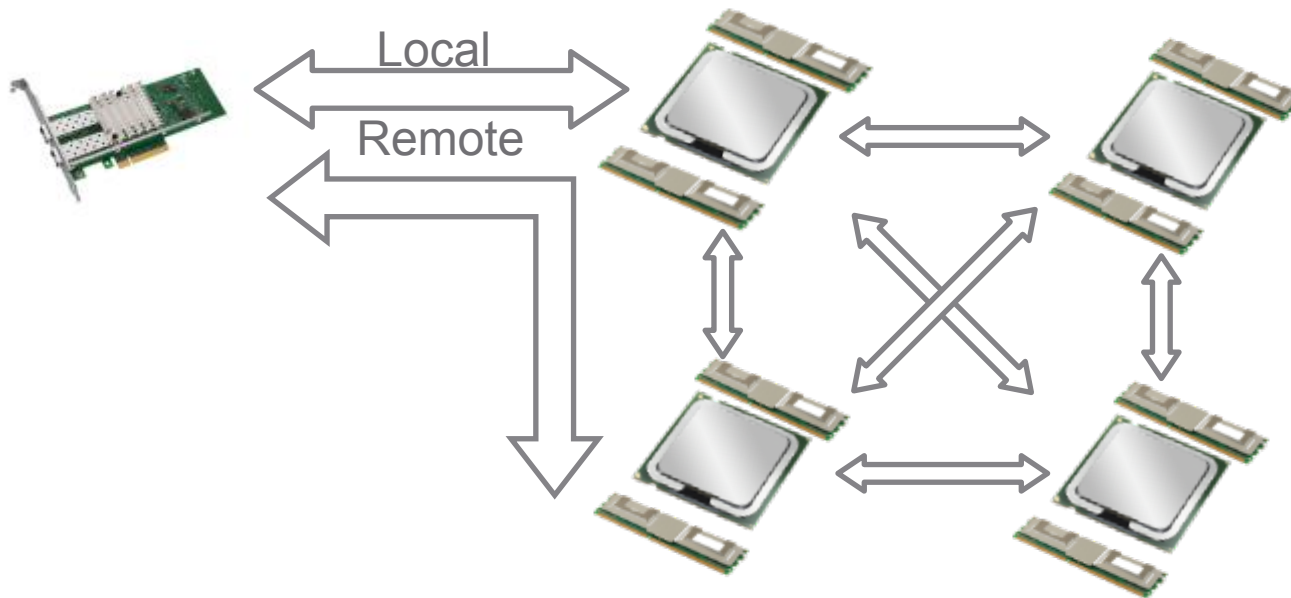
NUMA 1



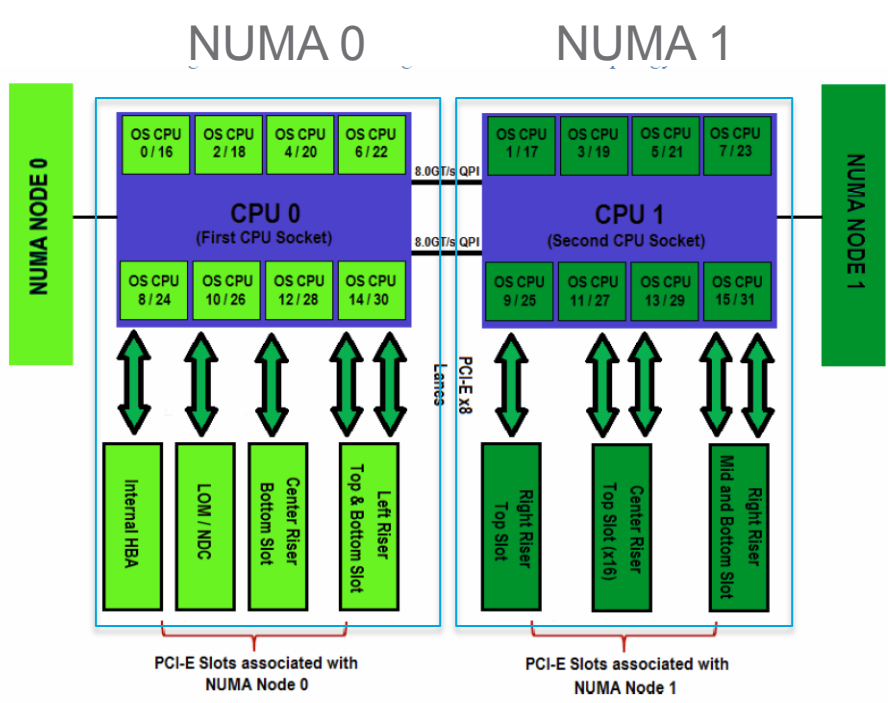
NUMA 2



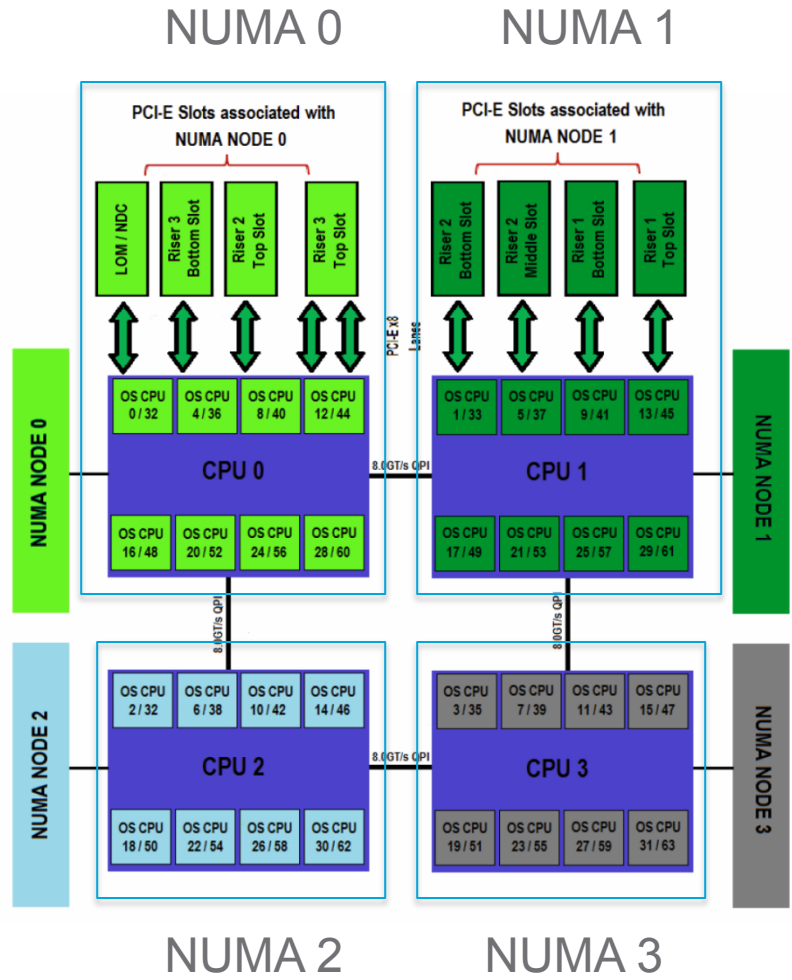
Non-Uniform I/O Access



Non Uniform I/O Access

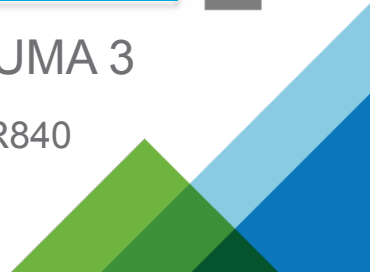


2 –socket Dell PE R720

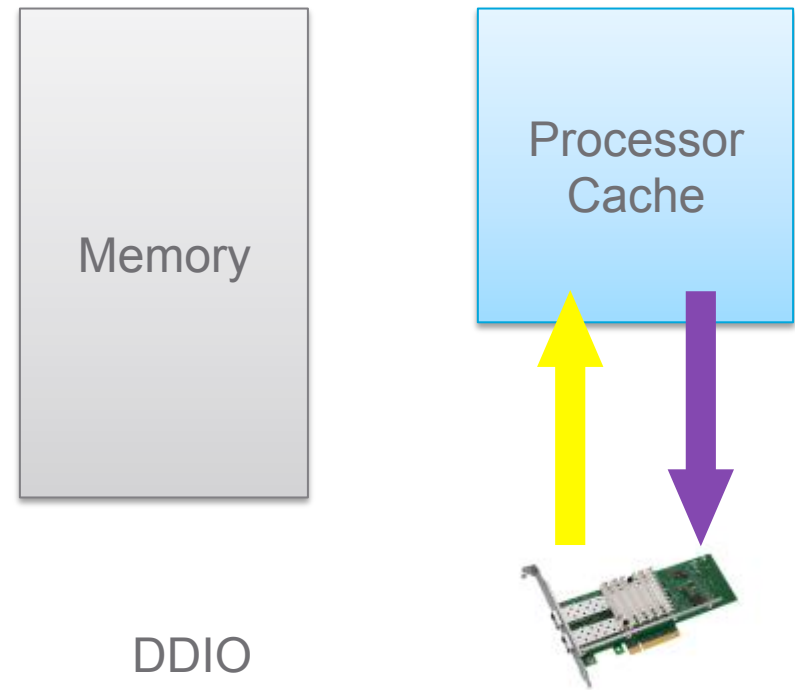
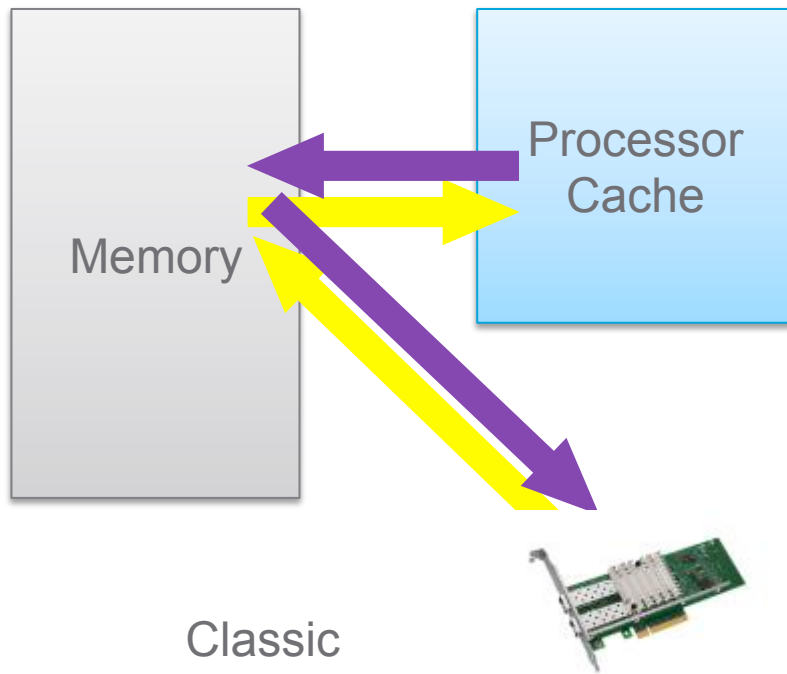


NUMA 2 NUMA 3

4 –socket Dell PE R840



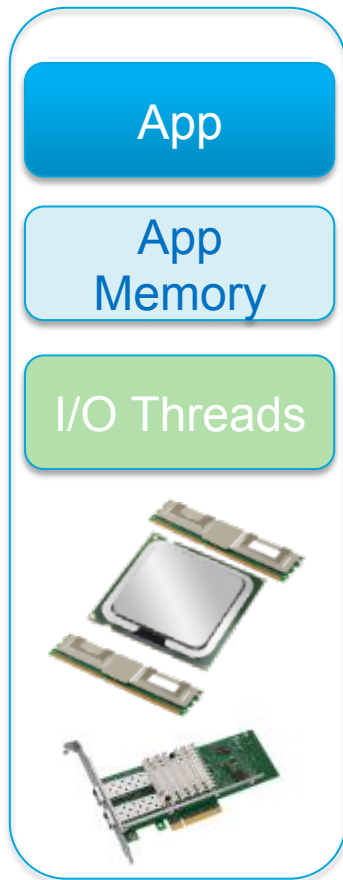
Intel DDIO



*Currently works with Local PCI Device



Non-Uniform I/O Access



Ideal Situation

Intel DDIO Benefits

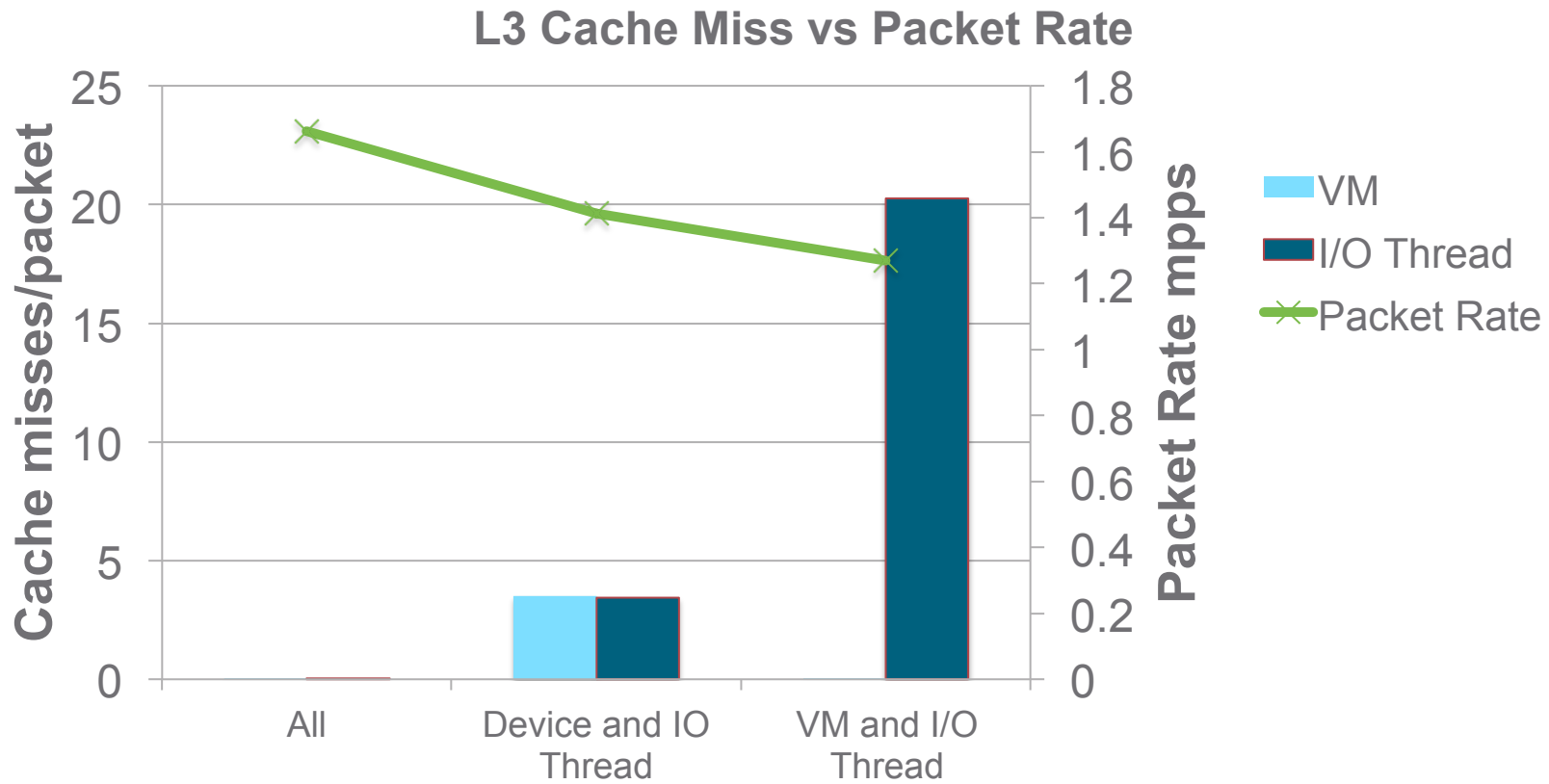


Real Situation

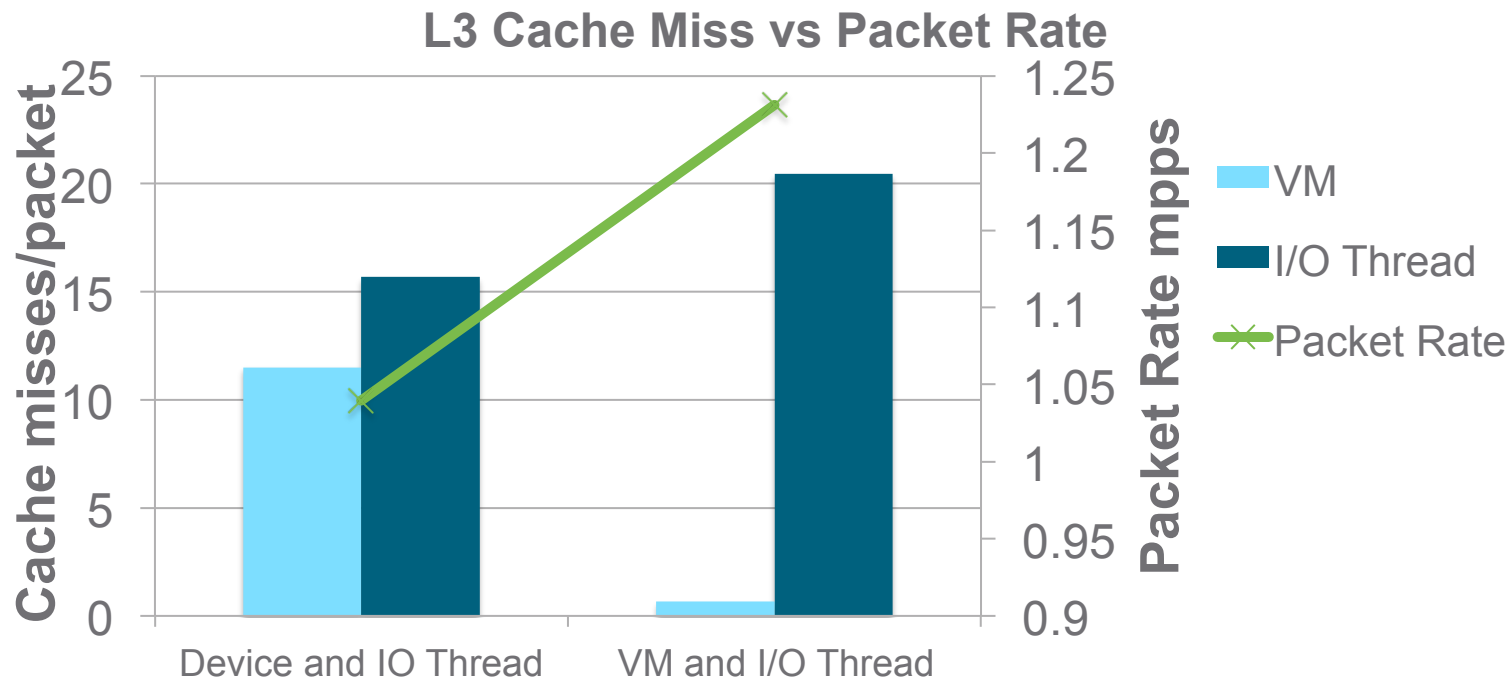
Local Memory Access



DDIO benefits vs Local I/O : 1VM



DDIO benefits vs Local I/O : 4VM

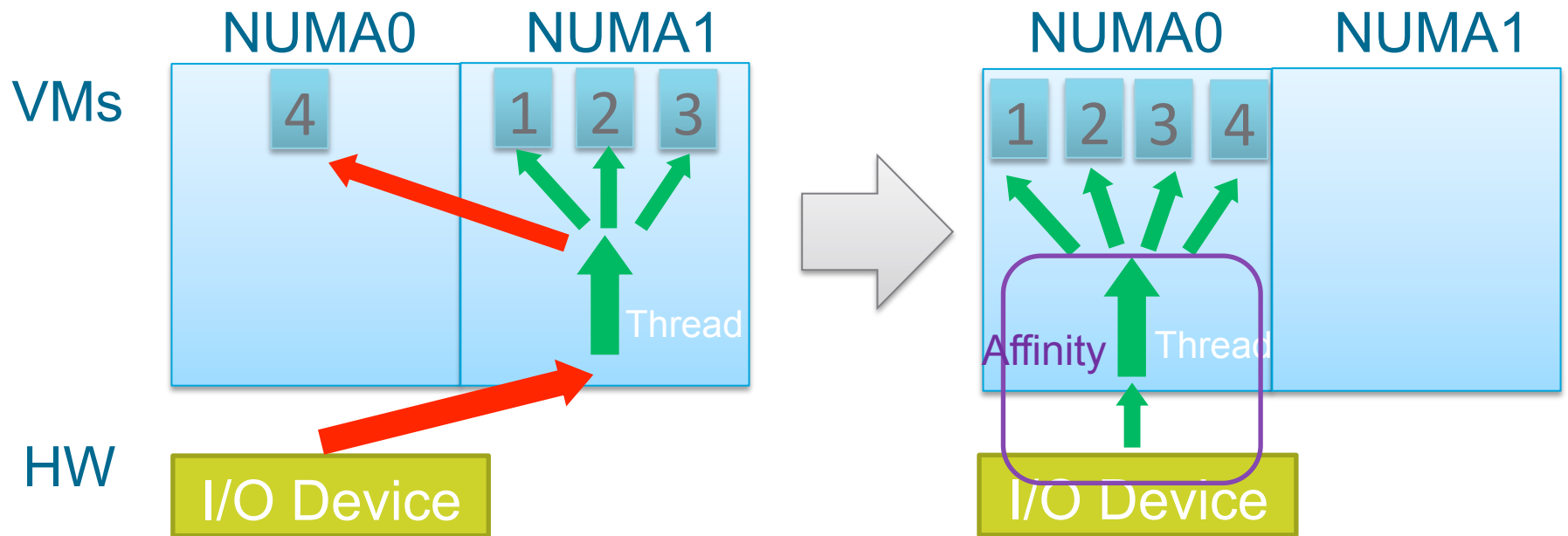


NUMA I/O Scheduler

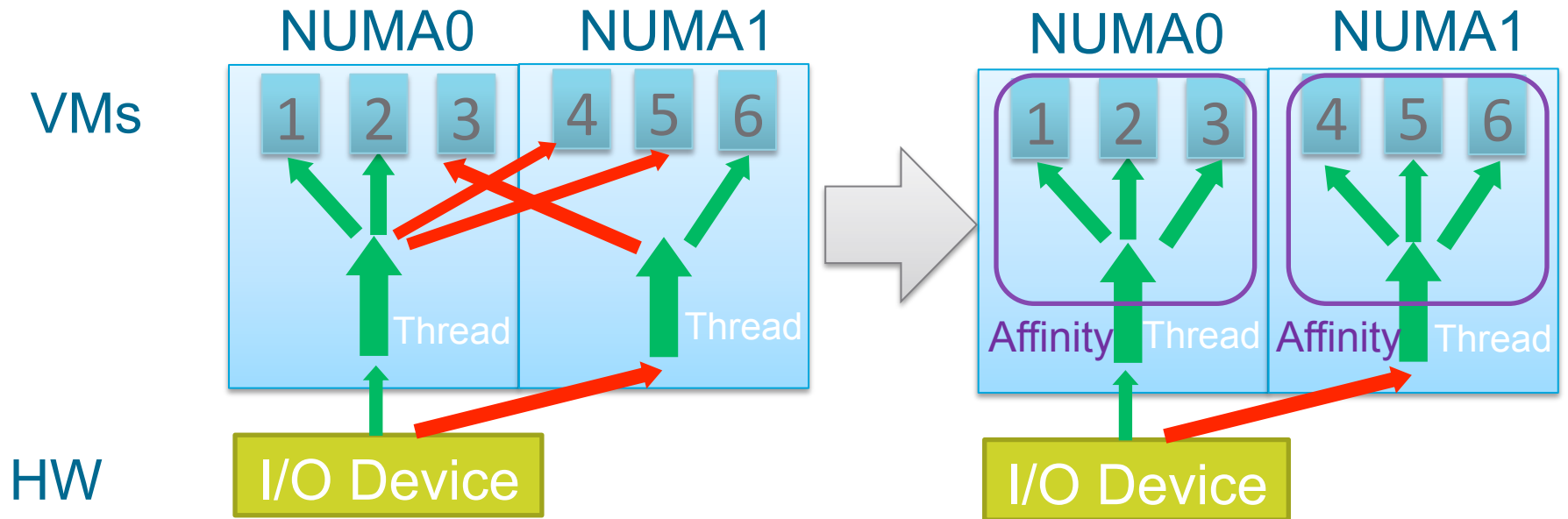
- Hybrid Mode
 - Low Load :
 - One I/O thread is sufficient for networking traffic
 - Pin I/O Thread to device NUMA Node
 - Let the scheduler migrate I/O intensive VM to device NUMA Node
 - High Load:
 - Sufficient load for multiple I/O Threads.
 - Create I/O Thread per NUMA Node.
 - Use I/O Thread based on Virtual Machines home node.



NUMA I/O: Light load

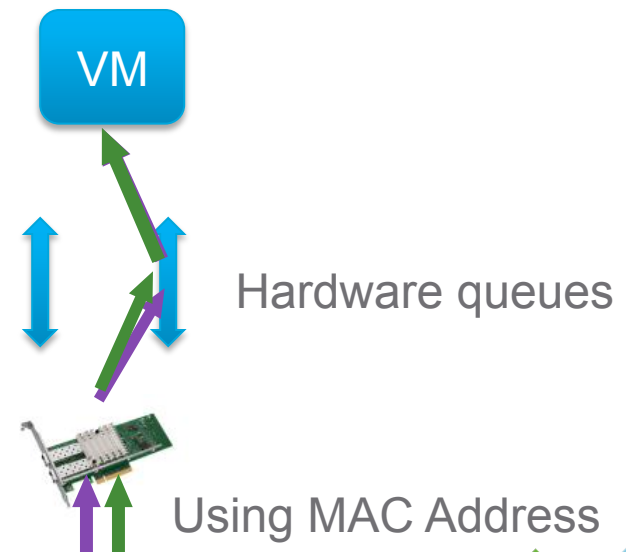
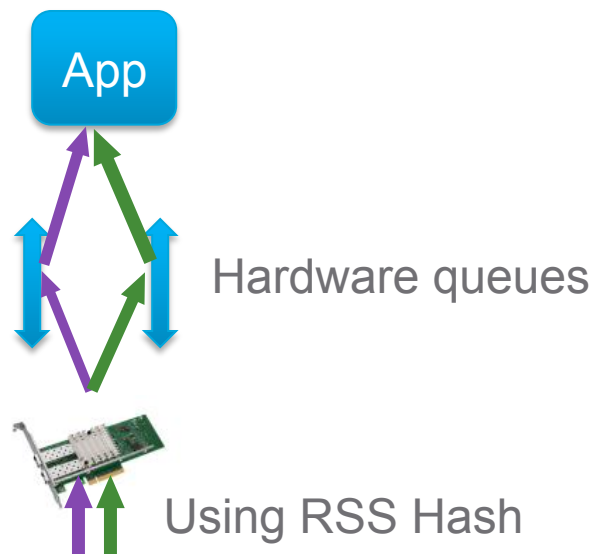


NUMA I/O: Heavy load

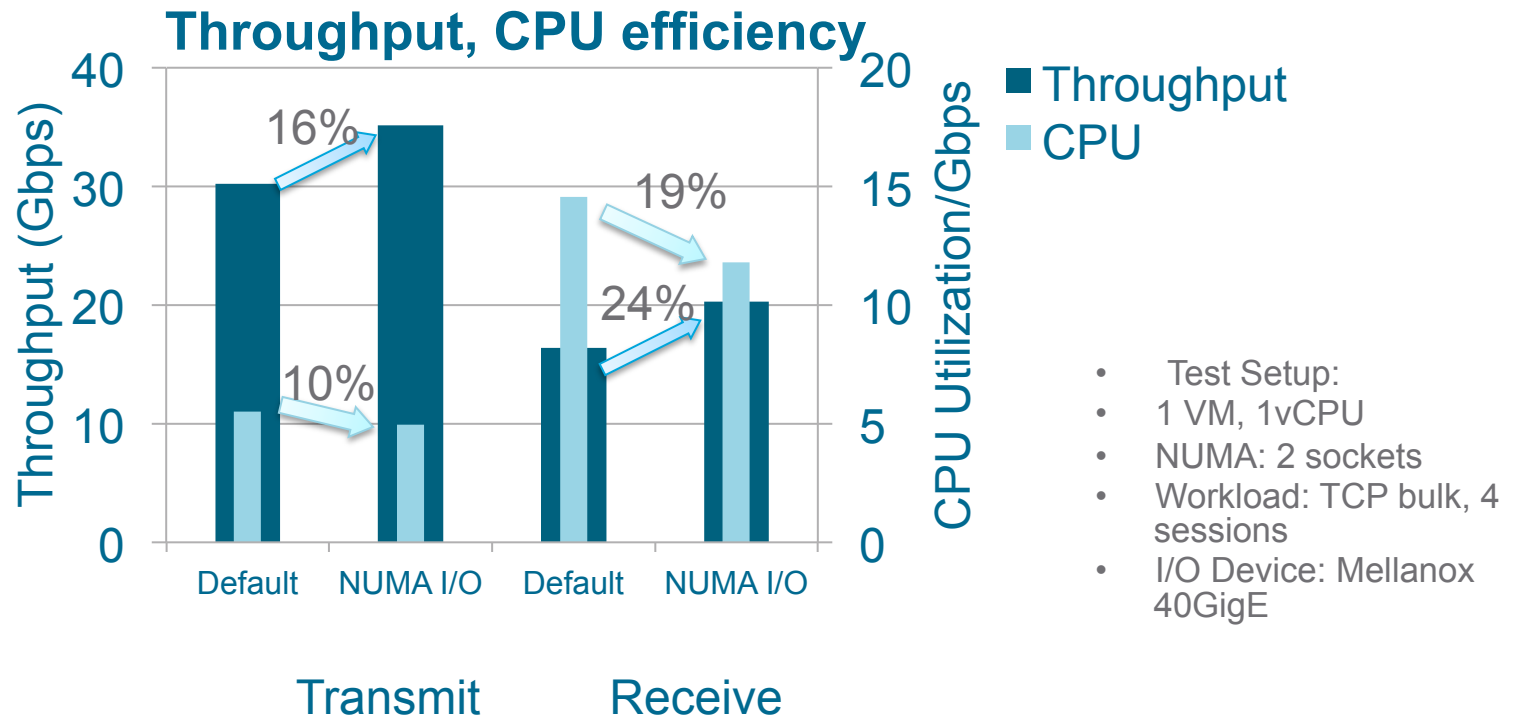


Hardware Support

- Most NICs support multiple queues.
 - Multiple contexts to process I/O requests.
 - More scheduling choices
- Virtual machine provides abstraction at an application layer
 - Easier to Isolate Apps and I/O traffic using MAC Address as traffic classification instead of 5-tuple hash.

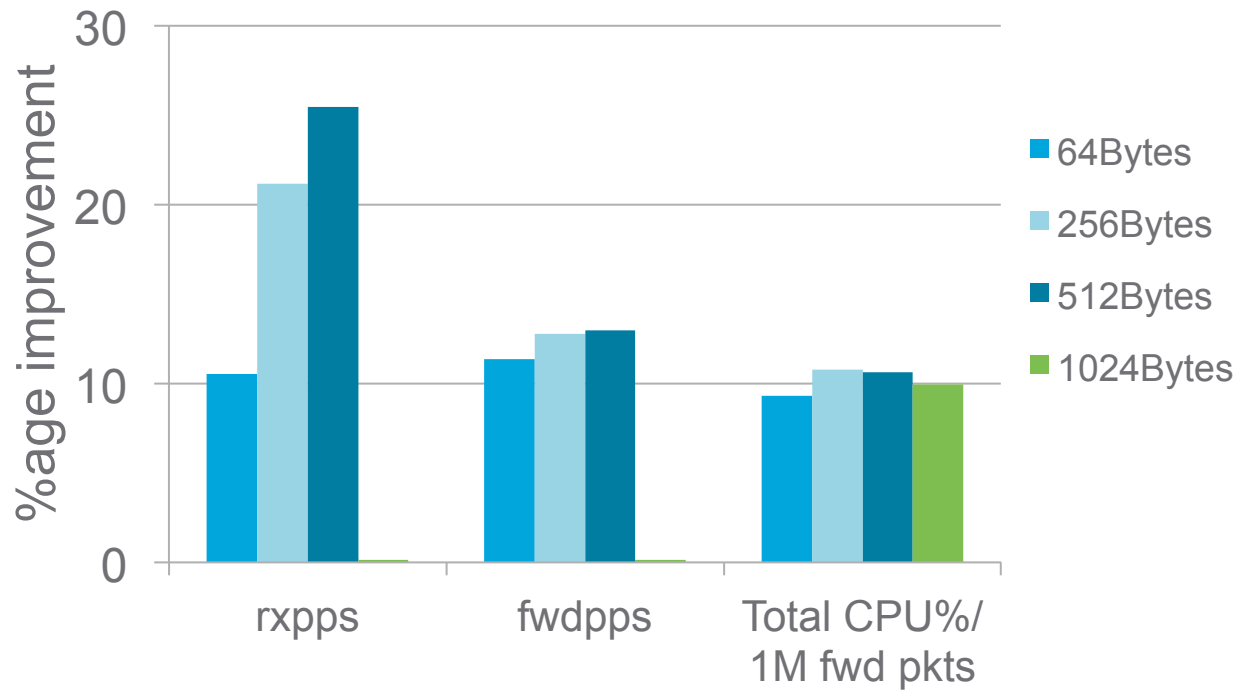


Improvement: Light load



Improvement: Light Load

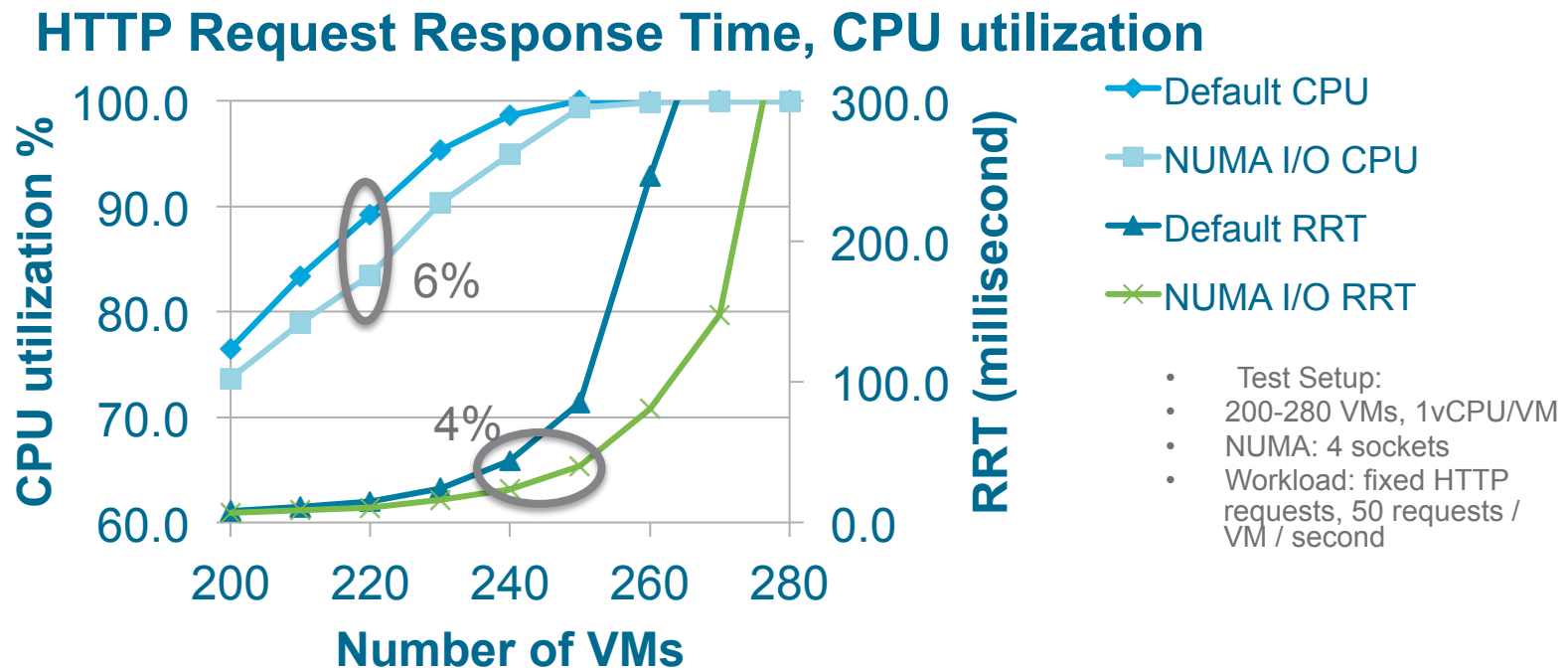
DPDK L2 Forwarding



Test Setup:
1VM, 1vCPU
NUMA: 2 sockets
Workload: UDP packet forwarding
I/O Device: Intel 10G Device



Improvement: Heavy load



Recap: NUMA aware I/O

- Hybrid design: low and high load
- Improvement:
 - higher throughput (Increased by 25%)
 - lower CPU (Reduced by 20%)
 - better VM consolidation (Up to 5% more VMs per host)
- Future work: GPU, storage, RDMA



Questions ?

vmware®

