# Characterizing the Performance of Parallel Applications on Multi-Socket Virtual Machines

Costin Iancu

Khaled Ibrahim, Steven Hofmeyr

*Lawrence Berkeley National Laboratory*

**Discovery 2015 – HPC and Cloud Computing Workshop**

***June 22 2011***

# Motivation

❖ **Virtualization is an enabling technology**

- Resource consolidation
- Fault tolerance & isolation

❖ **Virtualization Performance Expectations**

- Performance overhead is low (3-5% of raw)
- Current design and performance tuning techniques good enough!

❖ **HPC Workloads**

- Persistently use a large fraction of the system memory
- Data locality determines performance – NUMA support
- Sensitive to network bandwidth and latency – I/O support
- Use shared and/or distributed memory programming models – configuration/software support
- Most HPC studies are single socket or on dual core systems

# Virtualization Overhead
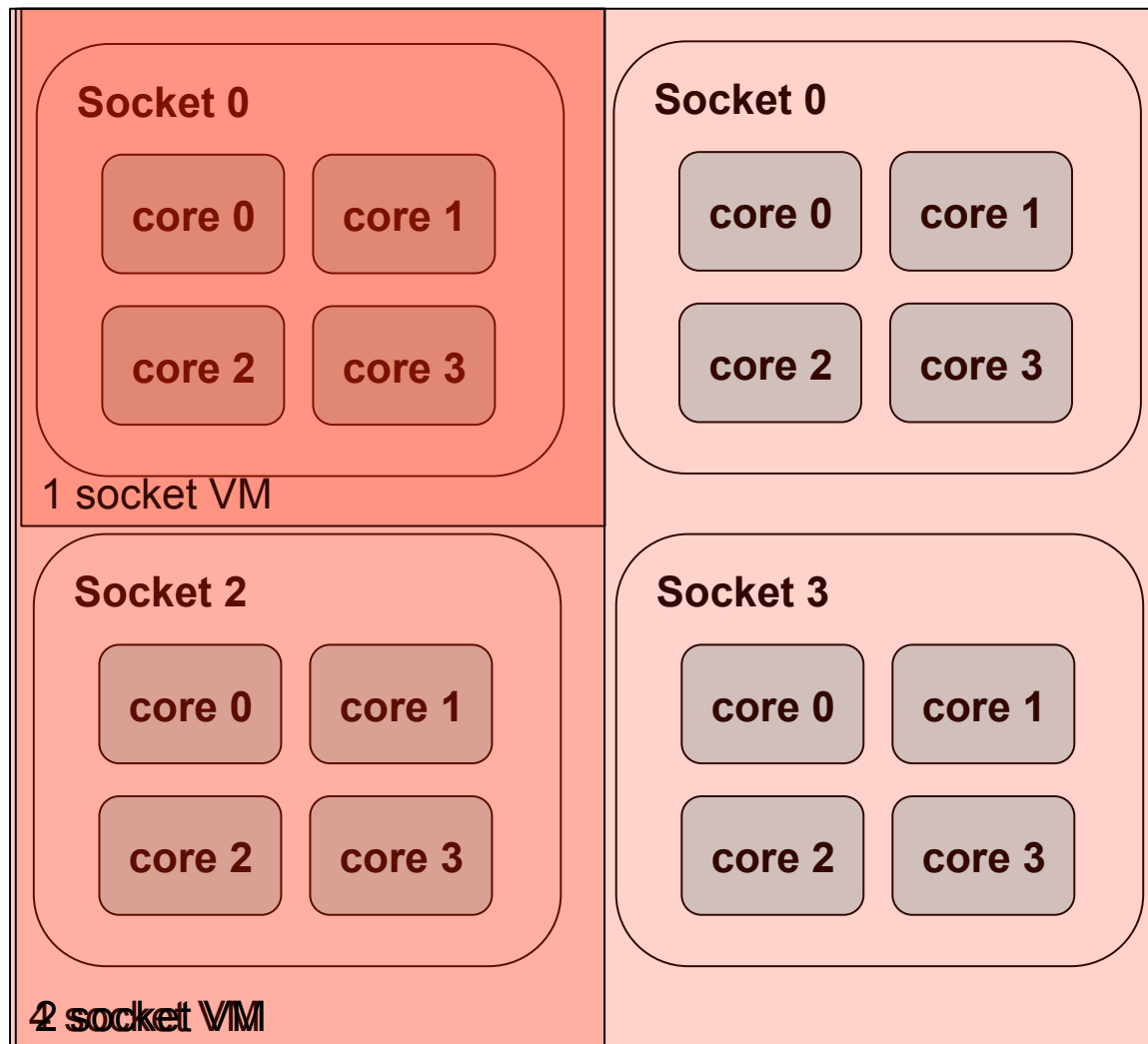
❖ **Three configurations**
  - 1 socket VM
  - 2 socket VM
  - 4 socket VM

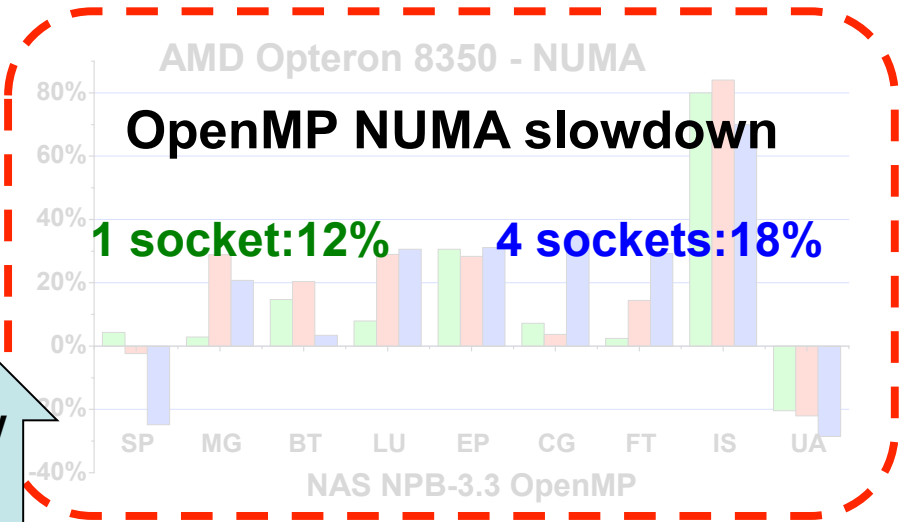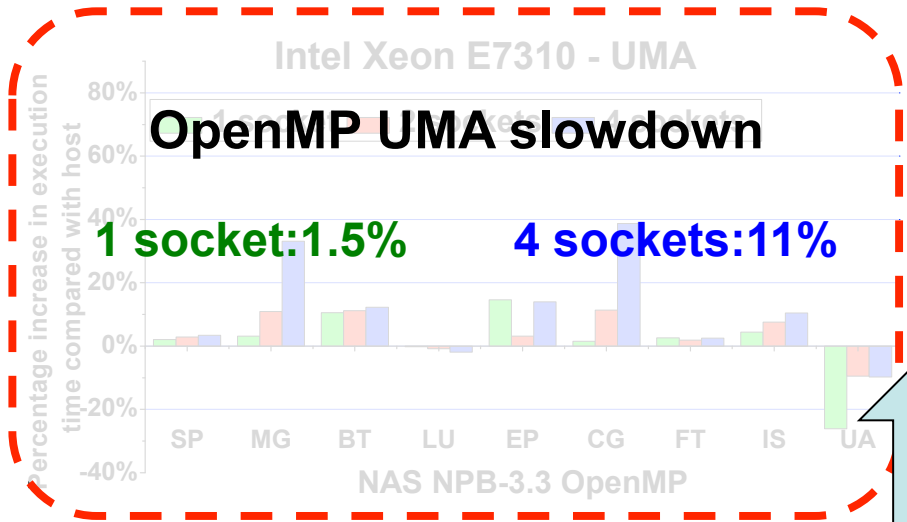❖ **Two architectures**
  - UMA
  - NUMA

❖ **Three programming models**
  - MPI
  - UPC
  - OpenMP

**Socket 0**
| | |
|---|---|
| core 0 | core 1 |
| core 2 | core 3 |

1 socket VM

**Socket 0**
| | |
|---|---|
| core 0 | core 1 |
| core 2 | core 3 |

**Socket 2**
| | |
|---|---|
| core 0 | core 1 |
| core 2 | core 3 |

**Socket 3**
| | |
|---|---|
| core 0 | core 1 |
| core 2 | core 3 |

4 socket VM    2 socket VM

**OpenMP UMA slowdown**

1 socket:1.5%        4 sockets:11%

**OpenMP NUMA slowdown**

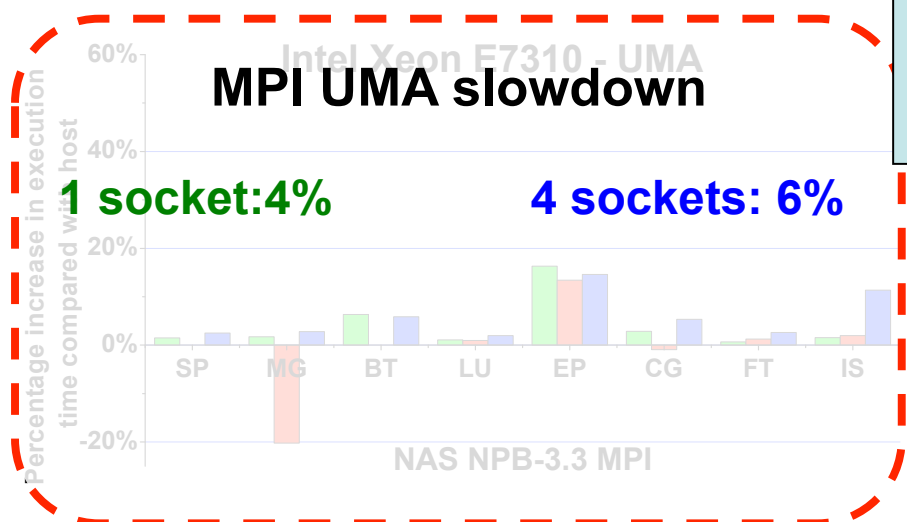1 socket:12%        4 sockets:18%

**MPI UMA slowdown**

1 socket:4%        4 sockets: 6%

**MPI NUMA slowdown**

1 socket:8%        4 sockets:40%

worse

❖ **Single socket performance is OK (KVM and Xen, matches performance expectations)**

❖ **Multi-socket UMA performance is OK ~ 10%**

❖ **High performance degradation when VMs span multiple NUMA domains:**

- KVM on average 40%
- Xen on average 233%

❖ **MPI seems to be slightly more affected than OpenMP**

1. **Reasons for performance degradation on multi-socket NUMA**

2. **Interaction between programming models and Virtualization**

3. **Techniques to improve NUMA support**

# Experimental Setup

- ❖ **Virtualization technology full H/W support for memory and I/O**
  - **KVM/QEMU 0.13.0**
  - **Xen 4.0**

- ❖ **NUMA support**
  - Xen 4.0 - NUMA support is the default setting for the hypervisor
  - Qemu-kvm allows NUMA emulation on the guest.

- ❖ **Benchmarks NAS Parallel benchmarks (3.3)**
  - MPI
  - OpenMP
  - UPC (Unified Parallel C)

- ❖ **Architectures-** Linux  (Kernel 2.6.32.8)
  - 4X4 UMA : Tigerton Xeon(R) CPU  E7310
  - 4X4 NUMA: AMD Opteron(tm) Processor 8350
  - 2X4 NUMA: Intel Xeon E5530 (Nehalem EP).

FUTURE TECHNOLOGIES GROUP

- ❖ **Vendor provided (Xen, KVM, VMWare, OpenBox, etc)**
  - ▪ Hypervisor manages NUMA locality of pages.
  - ▪ Guests are typically architecture neutral.

- ❖ **NUMA Page allocation**
  - ▪ On-demand: KVM, VMWare.
  - ▪ Pre-allocation: Xen (problematic for NUMA)
  - ▪ Two level translation (Xen, VMWare), three level (KVM)

- ❖ **Xen (The other open-source)**
  - ▪ 233% average slowdown (compared with 40% for KVM).

- ❖ **VMWare – restricted info**
  - ▪ Limited vcpus
  - ▪ Guest is not NUMA aware

- ❖ **Vendors advocate node confinement (1 VM per NUMA Domain).**

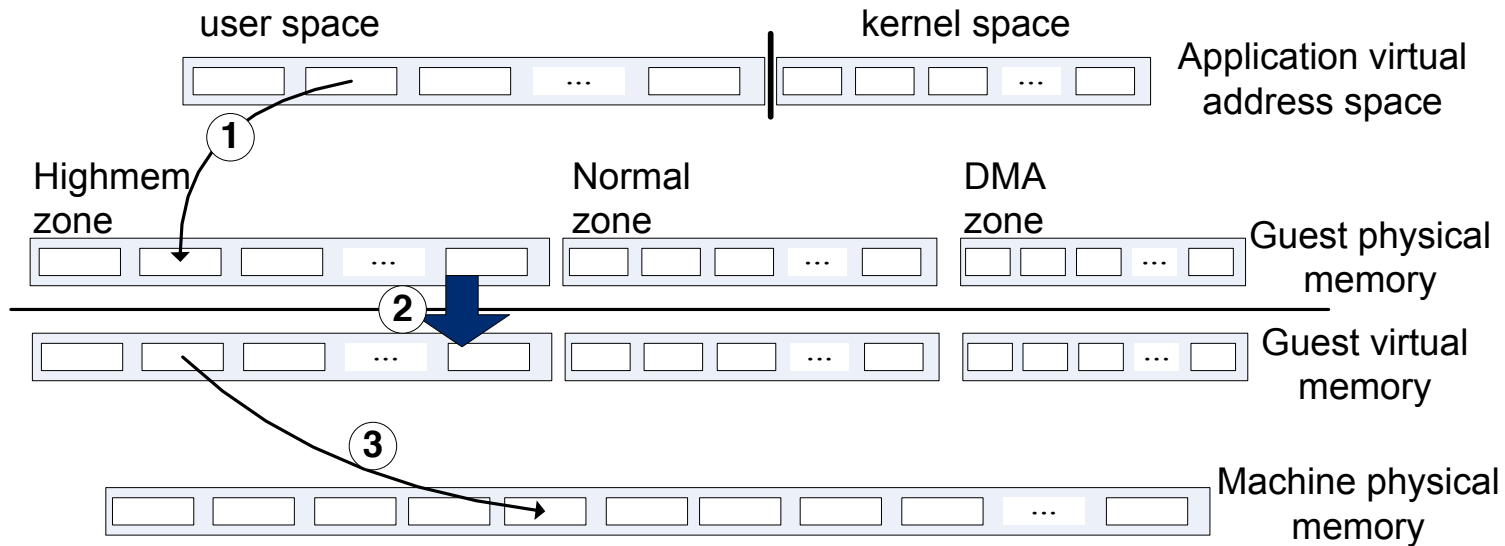❖ **Enabling NUMA, pinning and page granularity do not provide good multi-socket NUMA performance.**

❖ **Page granularity might affect performance**
  - Minor effect in our experiments.

❖ **Node confinement (1 VM per NUMA Domain).**
  - Implicitly assumes first-touch allocation
  - Requires pinning VMs and workloads, etc
  - Multi-socket?!

❖ **Is current support enough?**

user space

kernel space

Application virtual
address space

①

Highmem
zone

Normal
zone

DMA
zone

Guest physical
memory

②

Guest virtual
memory

③

Machine physical
memory
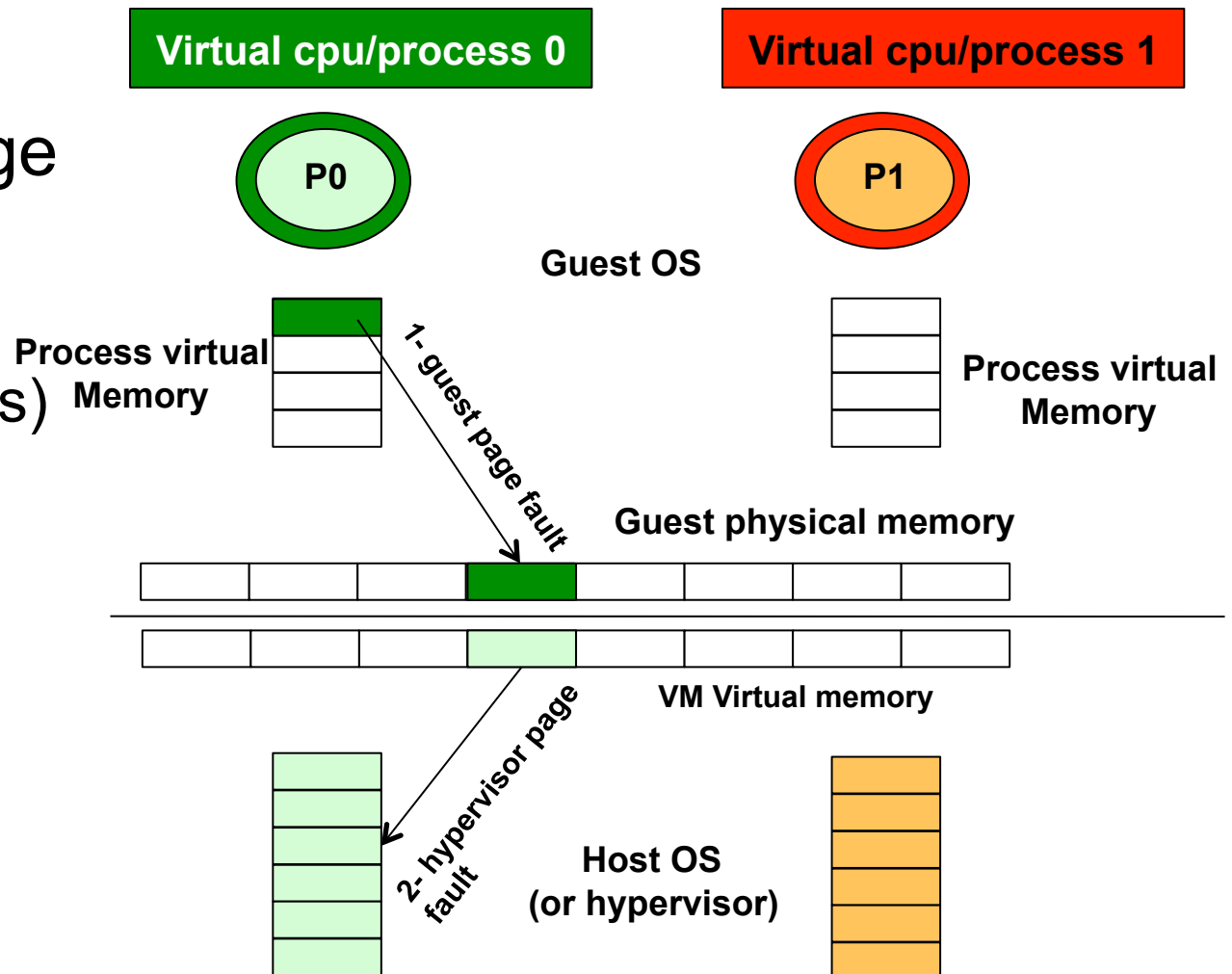
❖ **Three stage translation**
  ▪ 2 Dynamic (runtime) and one static (launch time)

❖ Cold touch involves two page faults

- Guest fault (NUMA oblivious)
- Hypervisor fault (NUMA aware)

**Virtual cpu/process 0**

**Virtual cpu/process 1**

P0

P1

Guest OS

Process virtual Memory

Process virtual Memory

1- guest page fault

Guest physical memory

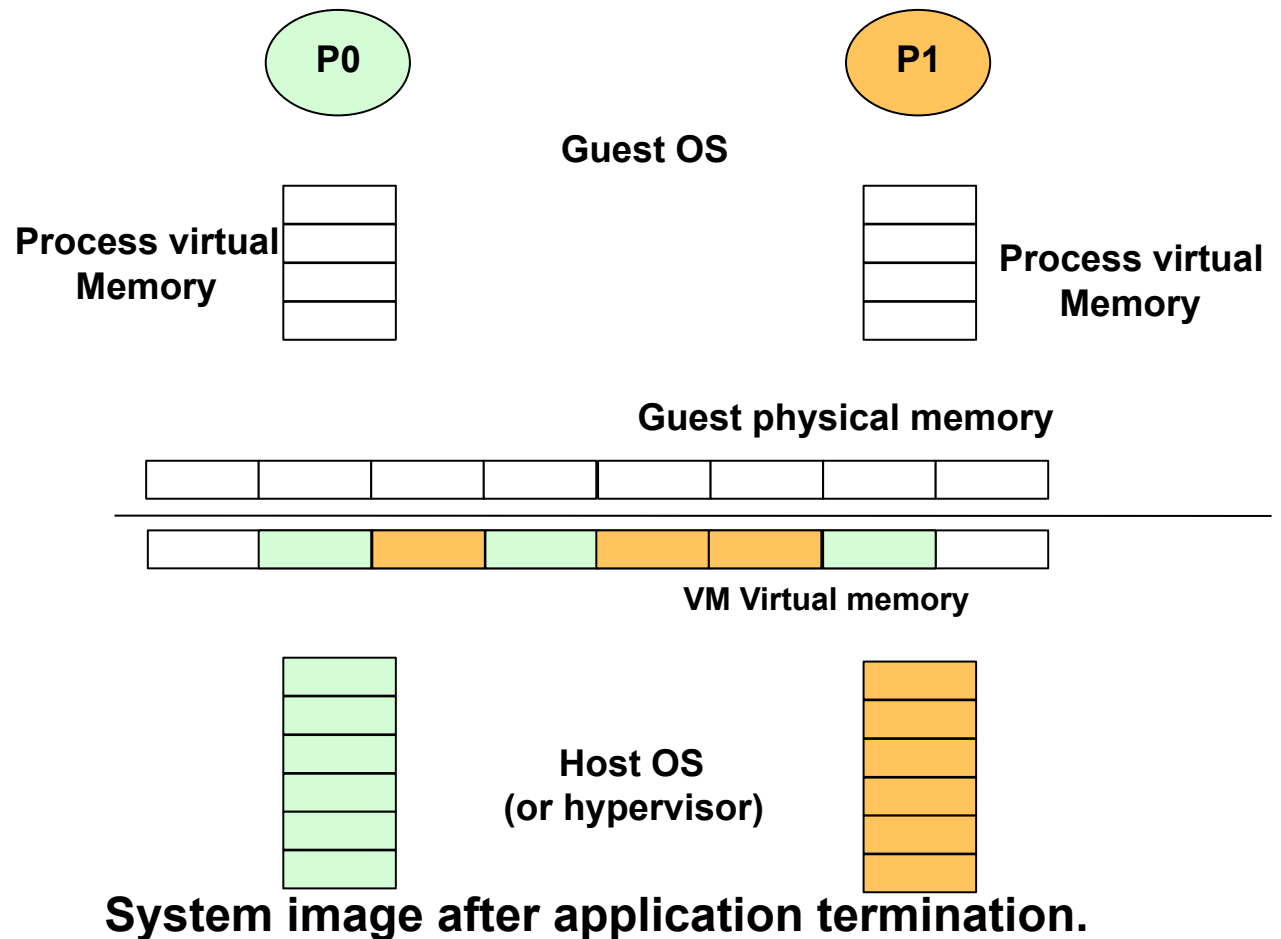VM Virtual memory

2- hypervisor page fault

Host OS (or hypervisor)

**Two phase translation mechanism for application for the first touch of a guest page**

❖ Correct NUMA affinity is managed by hypervisor.

**Virtual cpu/process 0**

**Virtual cpu/process 1**

P0

P1

**Guest OS**

**Process virtual Memory**

**Process virtual Memory**

**Guest physical memory**

**VM Virtual memory**

**Host OS (or hypervisor)**

**Two phase translation mechanism for application for the first touch of a page**

# Application Terminates

❖ Memory mappings in hypervisor are persistent.

P0

P1

Guest OS

Process virtual Memory

Process virtual Memory

Guest physical memory

VM Virtual memory

Host OS (or hypervisor)

**System image after application termination.**

**Virtual cpu/process 0**

**Virtual cpu/process 1**

P0

P1

Guest OS

❖ Hypervisor mapping is recycled and locality is not guaranteed.

Process virtual Memory

Process virtual Memory

Guest physical memory

VM Virtual memory

Host OS (or hypervisor)

**Page reuse results in host only page fault**

## Cold VM

Legend:
- Handled by guest
- Remote Multiple
- Remote Single
- Unmapped
- Local Node

By guest 18%
**Not mapped 75%**

Percentage of guest faults

NAS NBP-3.3 MPI

## Warm VM

Legend:
- Handled by guest
- Remote Multiple
- Remote Single
- Unmapped
- Local Node

By guest 70%
**Not mapped 0.6%**

Percentage of guest faults

NAS NBP-3.3 MPI

❖ **Hypervisor can provide locality**

❖ **Page faults are filtered by guests – do not reach hypervisor**

# Main Topics

1. Reasons for performance degradation on multi-socket NUMA

2. **Interaction between programming models and Virtualization**

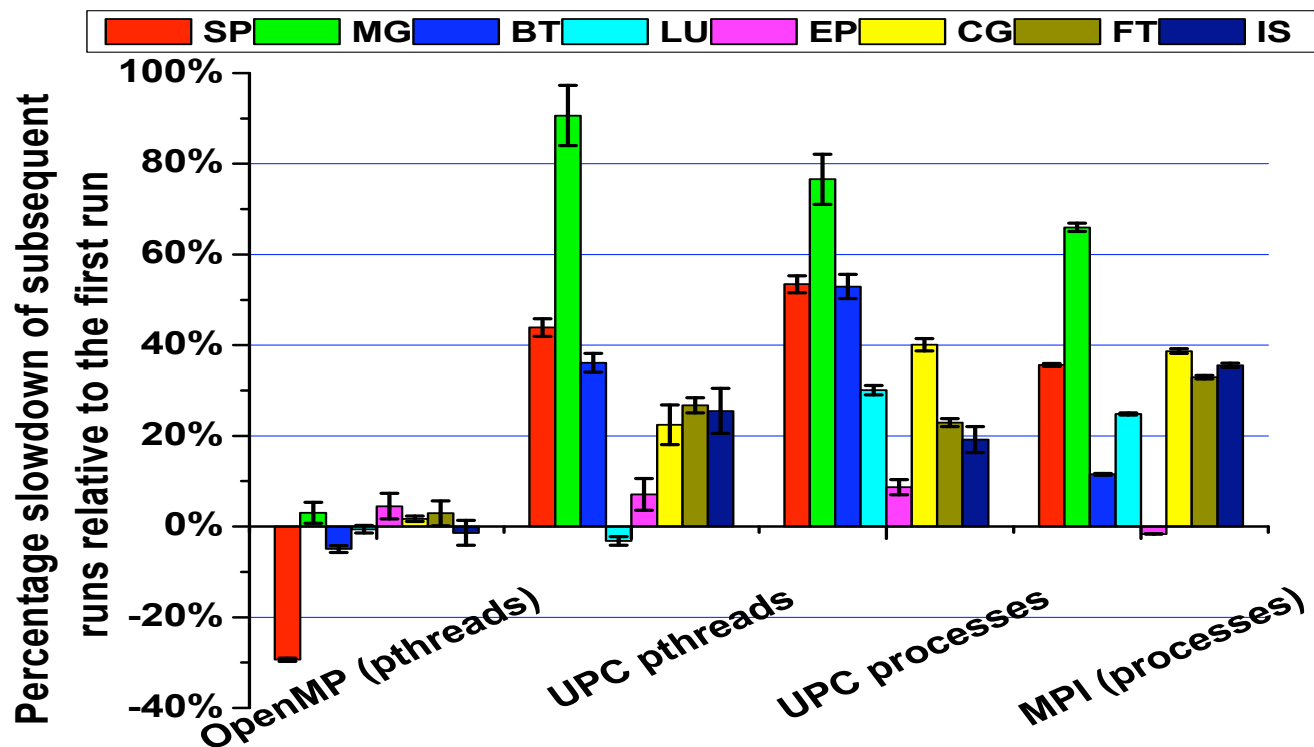3. Techniques to improve NUMA support

**Warm VMs provide lower performance!**



**First run avg. slowdown: 9%, second run avg. slowdown: 40%**

- **SP** ■ **MG** ■ **BT** ■ **LU** ■ **EP** ■ **CG** ■ **FT** ■ **IS**

*Percentage slowdown of subsequent runs relative to the first run* (y-axis, from -40% to 100%)

Categories: OpenMP (pthreads), UPC pthreads, UPC processes, MPI (processes)

❖ **Shared Memory (OpenMP)**

   ▪ No locality. Remote data are fetched each time they are needed.

❖ **Distributed Memory (MPI and UPC)**

   ▪ Implicit/explicit locality. Copy data locally before referencing them.

1. **Reasons for performance degradation on multi-socket NUMA**

2. **Interaction between programming models and Virtualization**

3. **Techniques to improve NUMA support**

- ❖ **How to improve locality?**
  - ▪ Hypervisor?
  - ▪ Guest?
  - ▪ Application? Shell? Runtime?

- ❖ **Expose NUMA architecture to the guest**
  - ▪ "Enlightenment" proposal for Xen

- ❖ **Modify memory management**
  - ▪ Page migration – hypervisor
  - ▪ Fault propagation – guests, hypervisor
  - ▪ Configuration/services

- ❖ **Use node confinement (partitioning)**
  - ▪ Transparent/configuration
  - ▪ With support – hypervisor, runtime

❖ **Expose NUMA architecture to the guest**

- How to over-commit memory?
- Can we handle non-contiguous NUMA nodes?
- How to flexibly manage memory of the VMs (reclamation, for instance)?
- How to resize memory?
- How to migrate VM to a non-compatible destination?
- Can the hypervisor commit to guarantee page node allocation?
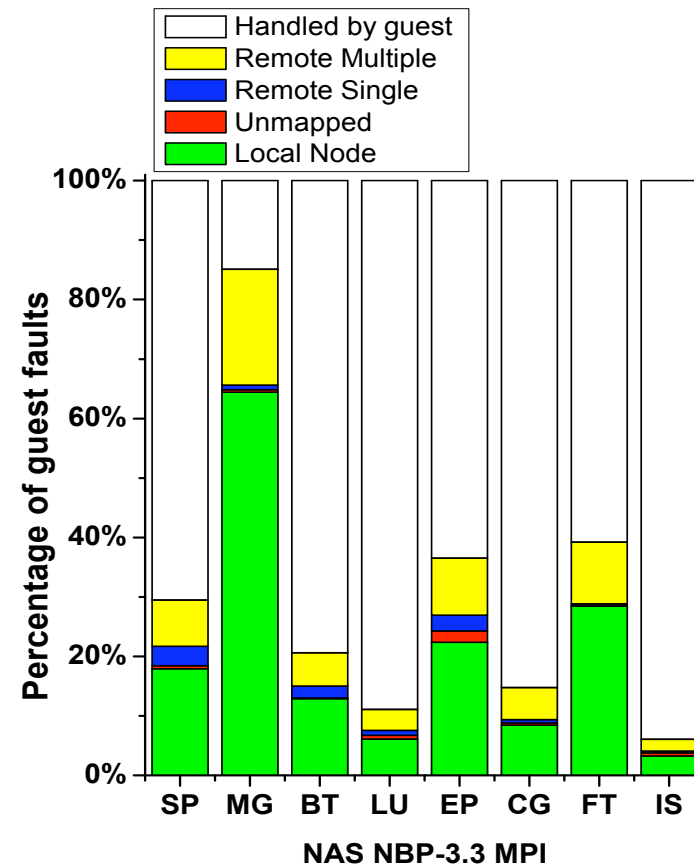
*Void virtualization benefits!*

❖ Page Migration: fix  locality for badly mapped pages.

  ▪ **1% remote single**

❖ Most faults handled  by the guest **(70%)**

❖ Propagating **faults** requires changes to all guest Oses

  ▪ Fast allocation

  ▪ Slow reclamation
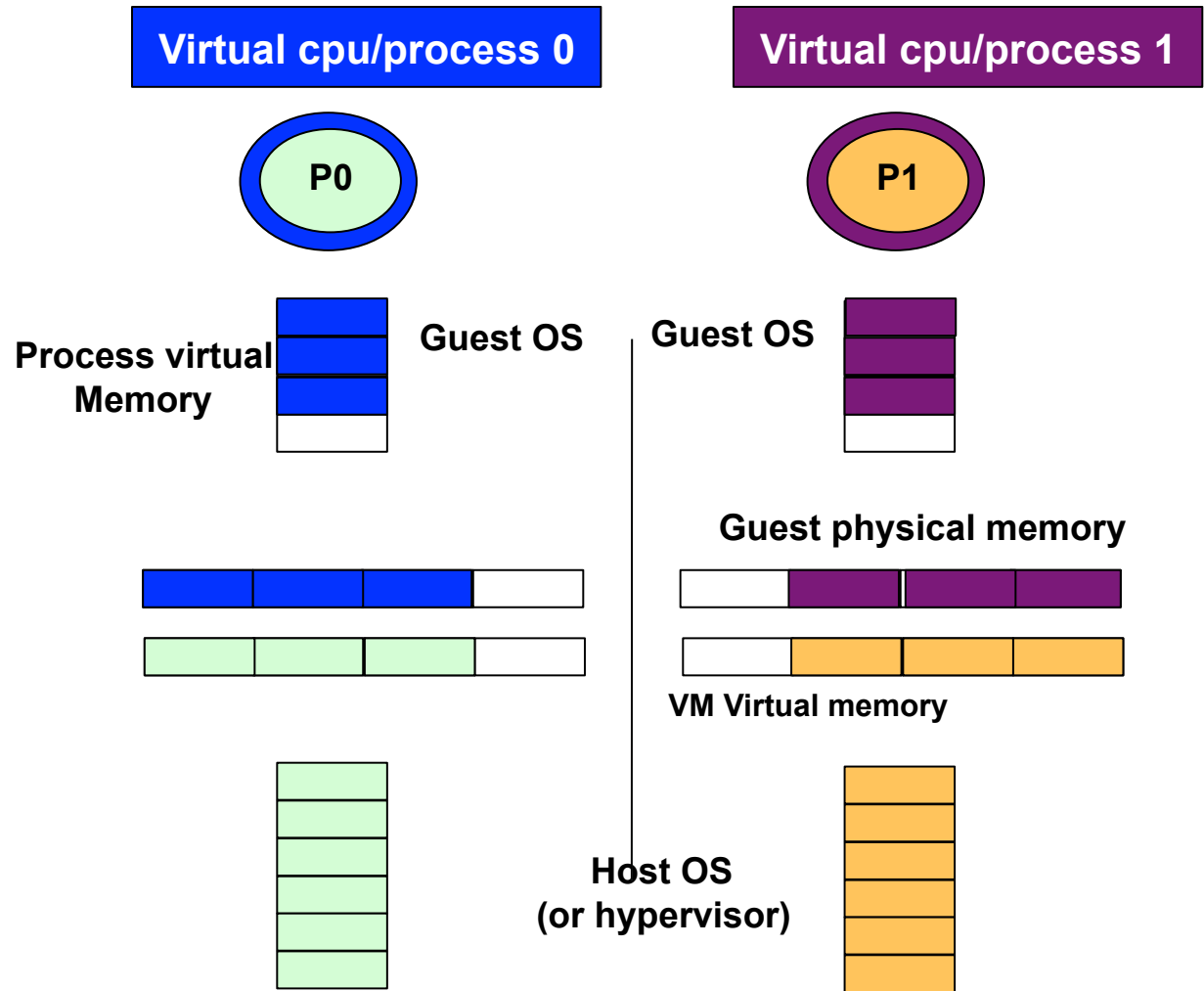
Page faults propagated to the hypervisor



Legend:
- ☐ Handled by guest
- ☐ Remote Multiple (yellow)
- ☐ Remote Single (blue)
- ☐ Unmapped (red)
- ☐ Local Node (green)

Y-axis: Percentage of guest faults (0% – 100%)
X-axis: SP MG BT LU EP CG FT IS

**NAS NBP-3.3 MPI**

❖ **Vendors advocate node confinement**

❖ **Performance:**
- Resource Contention
- Inter-VM communication

**Virtual cpu/process 0**

**Virtual cpu/process 1**

P0

P1

**Process virtual Memory**     Guest OS

Guest OS

**Guest physical memory**

**VM Virtual memory**

**Host OS (or hypervisor)**

**Page reuse results in host only page fault**

**1 VM per node 40% slowdown**

**1 VM per NUMA domain is 400% slowdown**



Legend:
- 16 VMs (4 VMs per socket)
- 8 VMs  (2 VMs per socket)
- 4 VMs  (1 VM per socket)
- 2 VMs  (2 sockets per VM)
- 1 VM   (4 sockets per VM)

y-axis: Percentage increase in execution time compared with host

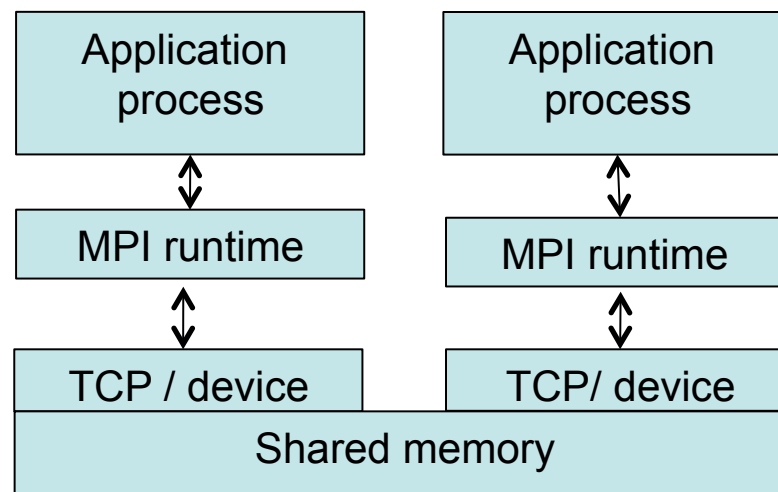x-axis categories: SP  MG  BT  LU  EP  CG  FT  IS

NAS NBP3.3 MPI

worse

- ❖ Up to **16x** performance degradation, mostly more than **2x**.
- ❖ HPC workloads depends on efficient inter-VM communication.

❖ Earlier proposals implement communication stack over- shared memory

- Zhang et al [Middleware'07] IP over shared memory.
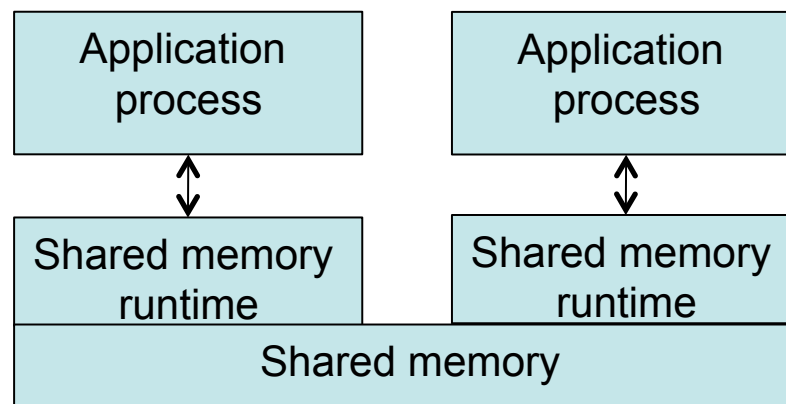- Huang et al. [SC'07] introduce IVC stack
- *virtio* essentially does the same.

❖ The bottleneck is in using the software stack.

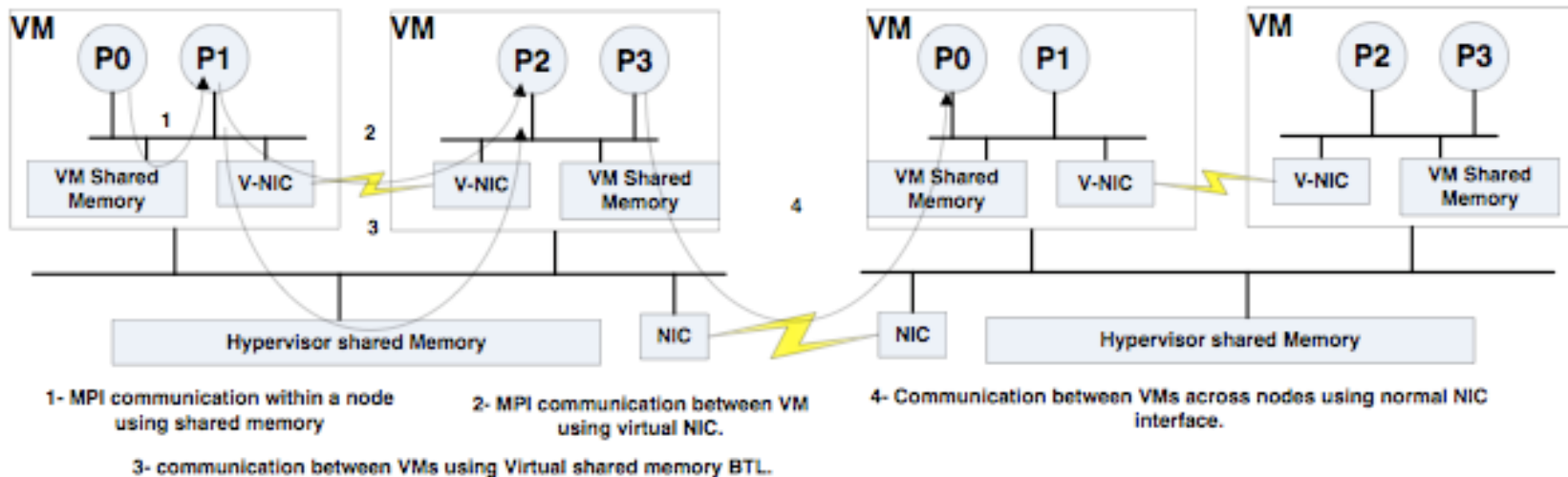❖ Instead, we implement inter-VM communication natively on top of shared memory.

Typical communication across nodes

| Application process | Application process |
|---|---|
| MPI runtime | MPI runtime |
| TCP / device | TCP/ device |
| Shared memory | |

Typical communication within a node

| Application process | Application process |
|---|---|
| Shared memory runtime | Shared memory runtime |
| Shared memory | |

1- MPI communication within a node using shared memory
2- MPI communication between VM using virtual NIC.
3- communication between VMs using Virtual shared memory BTL.
4- Communication between VMs across nodes using normal NIC interface.

- ❖ Shared memory exposed to guest as PCI device memory (ivshmem driver)
- ❖ Three new components handle the shared memory between different VMs
  - ▪ VM Shared memory communication component.
  - ▪ VM memory pool communication component.
  - ▪ VM collective communication component.
- ❖ New selection mechanism for communication component.
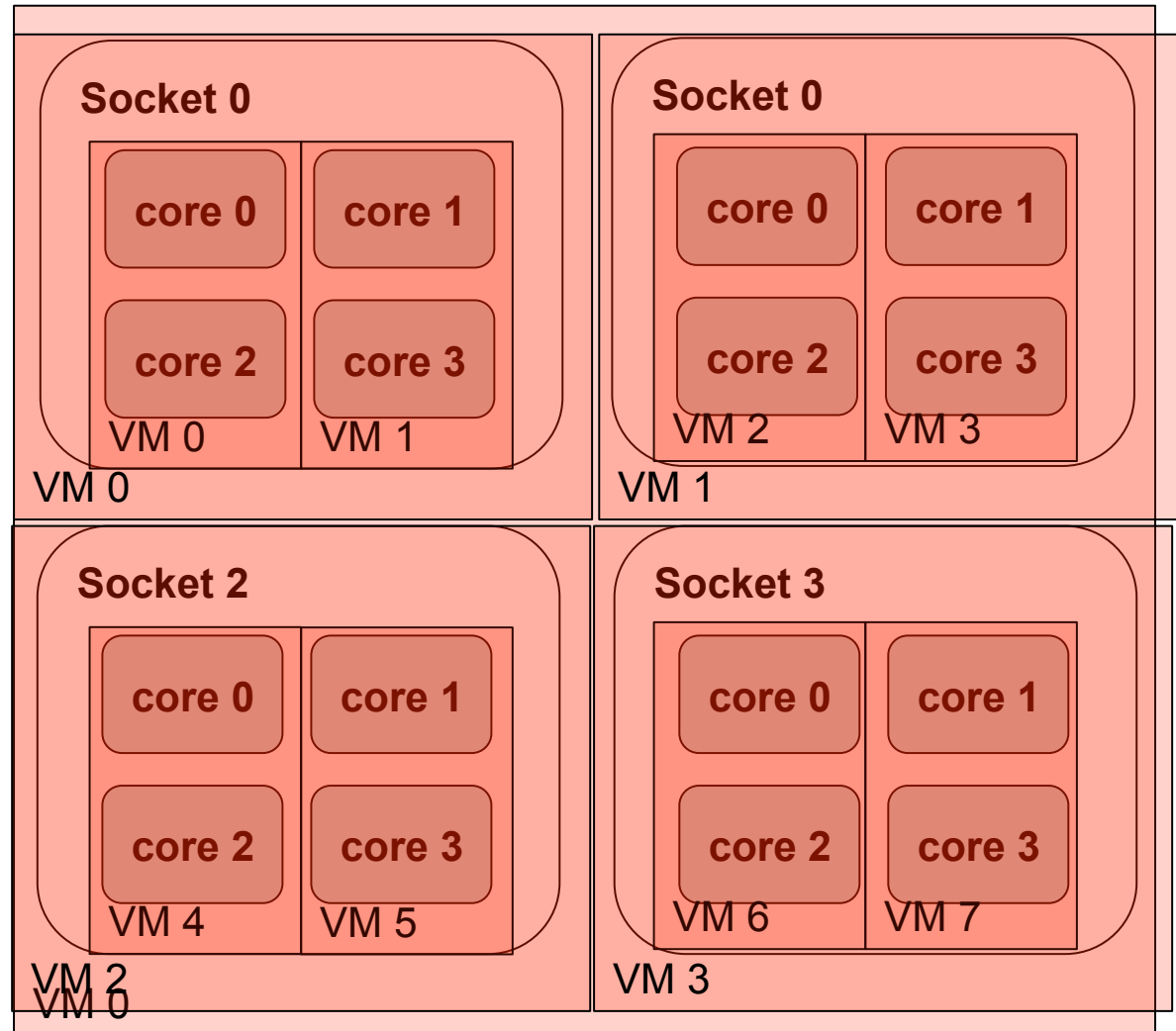- ❖ Similar mechanism is implemented for UPC

❖ Partitioning strategy

- 1 VM (4 socket per node)
- 4 VM (1 socket per node)
- 8 VM (2VM per socket)
- …

**One VM per node (1VM)**
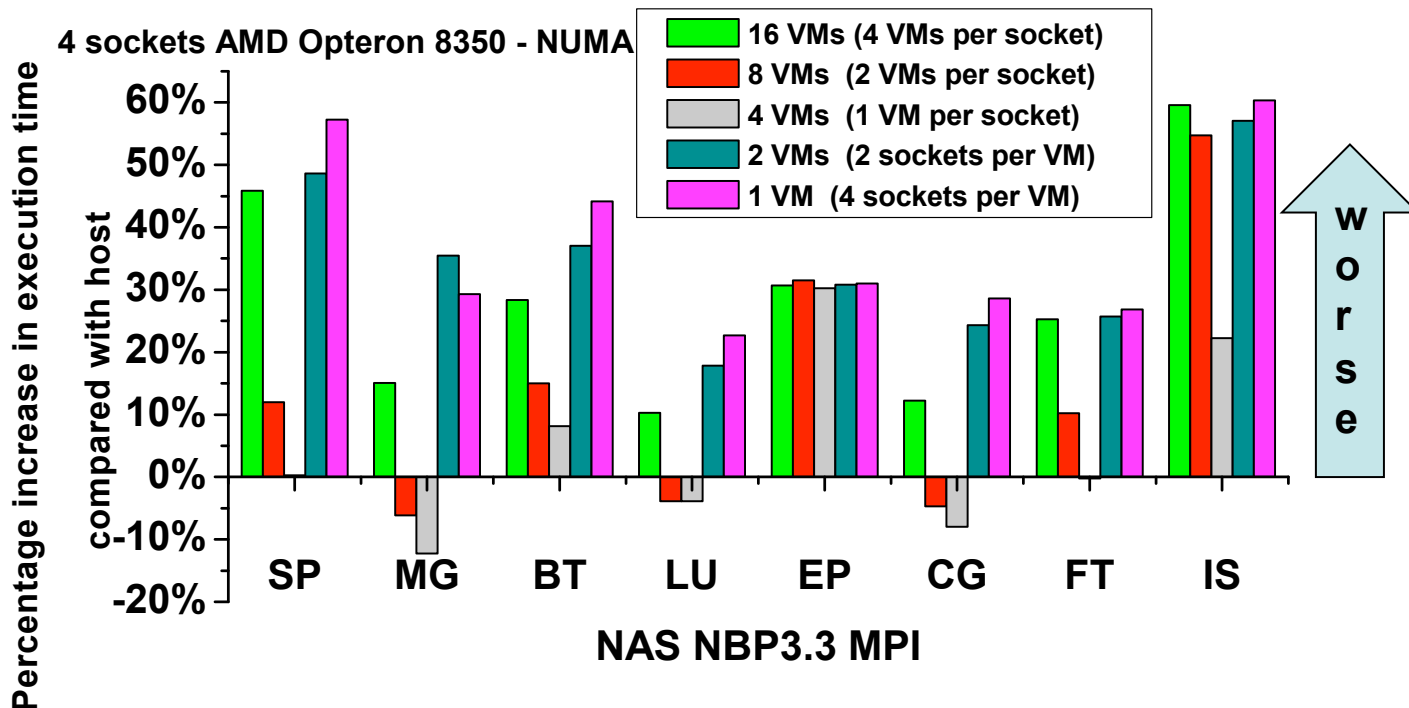**Slowdown: 40%**

**One VM per NUMA domain: (4VM)**
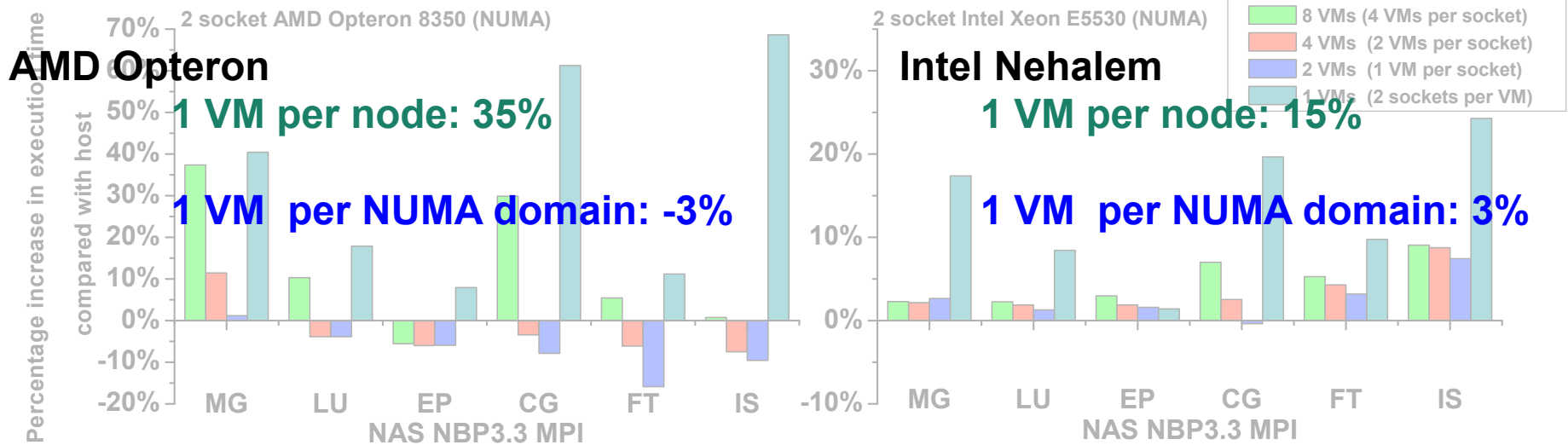**Slowdown: 3%**

4 sockets AMD Opteron 8350 - NUMA

Percentage increase in execution time compared with host

Legend:
- 16 VMs (4 VMs per socket)
- 8 VMs  (2 VMs per socket)
- 4 VMs  (1 VM per socket)
- 2 VMs  (2 sockets per VM)
- 1 VM   (4 sockets per VM)

NAS NBP3.3 MPI

SP  MG  BT  LU  EP  CG  FT  IS

worse

❖ One VM per socket is usually the best configuration.
❖ Efficient Inter-VM communication is key to performance.

❖ VM spanning sockets is always less efficient than multiple VMs with efficient inter-VM communication.



**AMD Opteron**

1 VM per node: 35%

1 VM per NUMA domain: -3%

2 socket AMD Opteron 8350 (NUMA)

**Intel Nehalem**

1 VM per node: 15%

1 VM per NUMA domain: 3%

2 socket Intel Xeon E5530 (NUMA)

Legend:
- 8 VMs (4 VMs per socket)
- 4 VMs (2 VMs per socket)
- 2 VMs (1 VM per socket)
- 1 VMs (2 sockets per VM)

# Other Benefits of Partitioning

❖ **Partitioning and resource contention**

- Introduces multi-level locking
- Reduces "system" overhead – e.g. MPI on UMA 6%->3%

❖ **Partitioning and I/O**

- KVM software driver – best is 1core per VM
  - MPI overhead: 17% on 32 VMs, 223% on 2 VMs
- Virtio – best is 8 cores per VM (12%)
  - MPI overhead: 34% on 32VMs, 63% on 2VMs

# Conclusions

❖ **The performance on NUMA machines is severely penalized if a VM span multiple sockets (avg. slowdown: 40% KVM, 223% Xen).**

❖ **NUMA cannot be handled by hypervisor alone**
- Lacking (Xen), or hindered by guest (KVM locality leakage).

❖ **VM partitioning requires efficient inter-VM communication**
- Better than virtualized IO or communication stack on top of shared memory.
- Our implementation reduces slowdown to 3% on average.

❖ **Other solutions may be needed for shared memory programming models, for instance OpenMP.**

# Questions

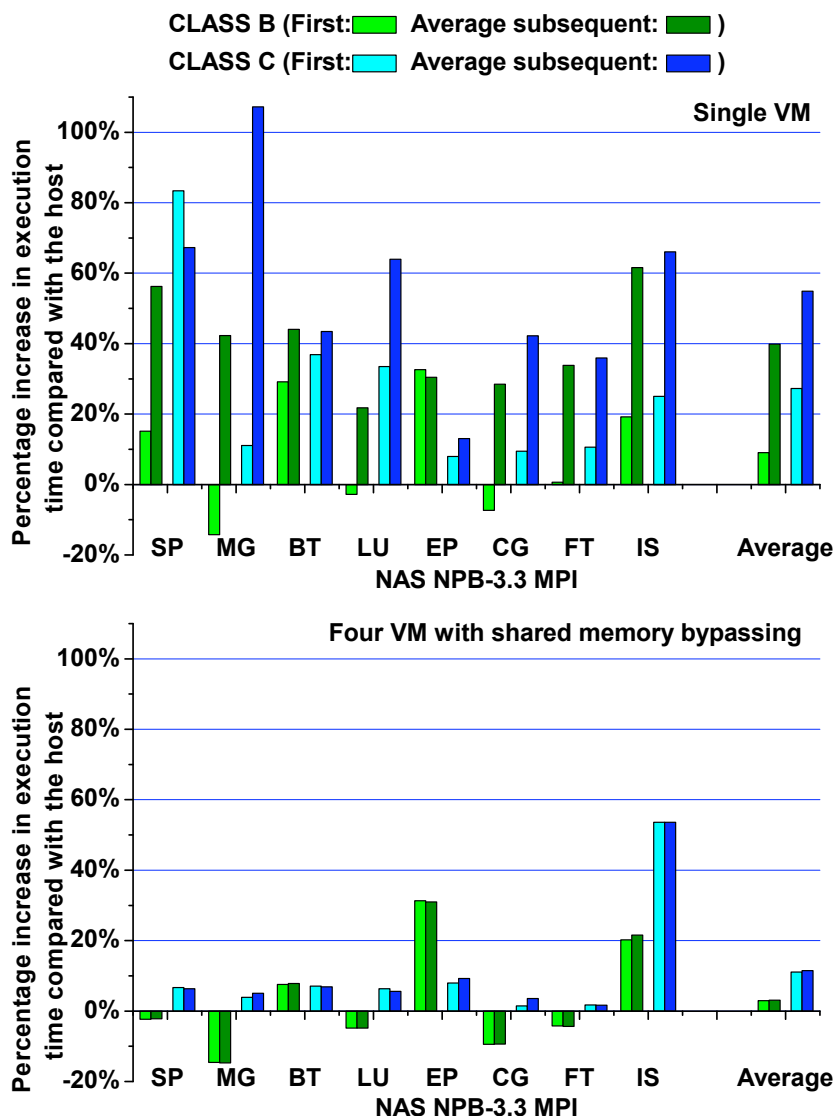**Thanks for attending!**

# Impact of Dataset

- ❖ First run performance becomes less optimal for the large dataset.
- ❖ Less data are cached so bad locality is associated with higher cost.
  - Class B: avg. 40% slowdown (up to 61%)
  - Class C: 57% in average (up to 105%)
- ❖ With partitioning and efficient communication
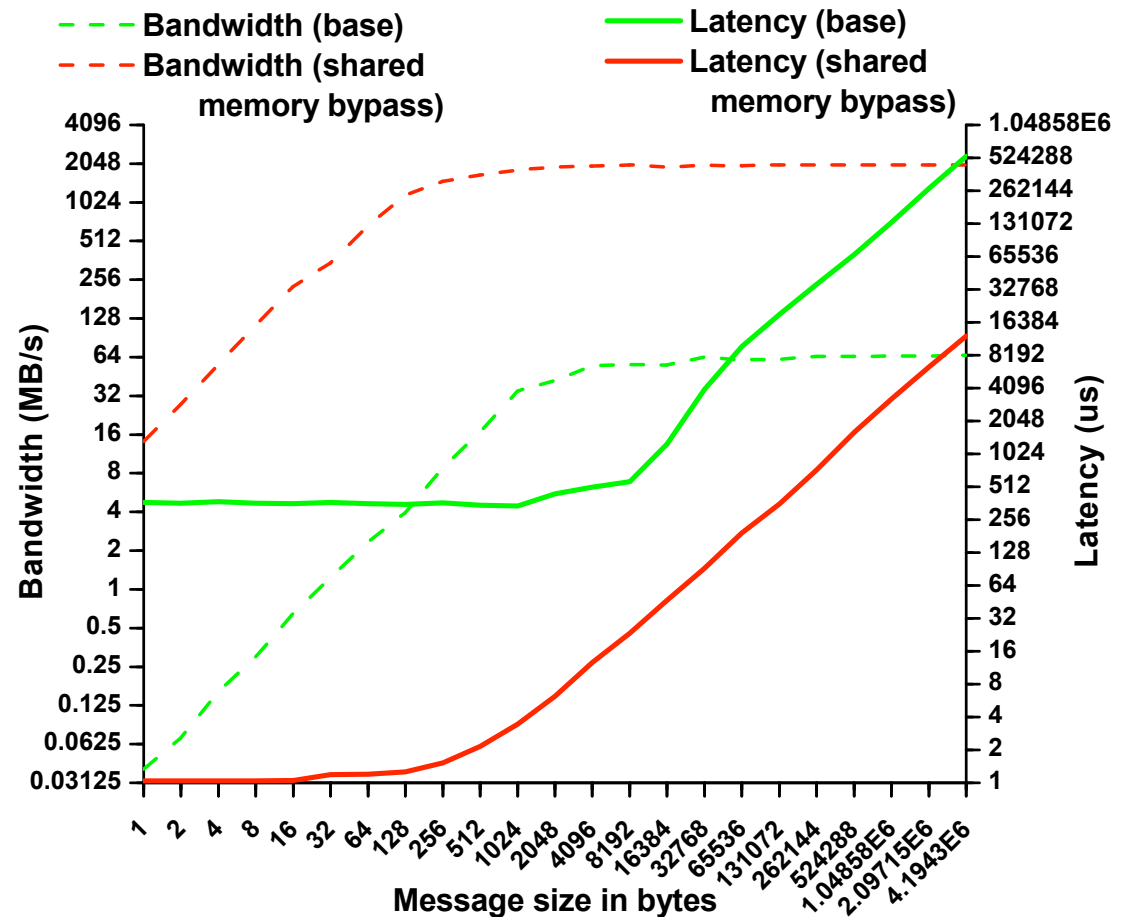  - Class B: avg. 3%
  - Class C: avg. 11%



CLASS B (First: ▮ Average subsequent: ▮ )
CLASS C (First: ▮ Average subsequent: ▮ )

Single VM

Percentage increase in execution time compared with the host

NAS NPB-3.3 MPI

Four VM with shared memory bypassing

- ❖ MPI network performance for TCP network vs shared memory bypassing.

- ❖ Benchmark implementations have similar NUMA domain distribution ( have well balanced page fault distribution across domains)
- ❖ The implementation of the programming model  affects behavior:
  - ❖ pthread model vs. processes (Higher percentage of faults exposed in the first run for OpenMP.)

- ❖ NUMA distribution + implementation of runtime do not explain perforamance differences
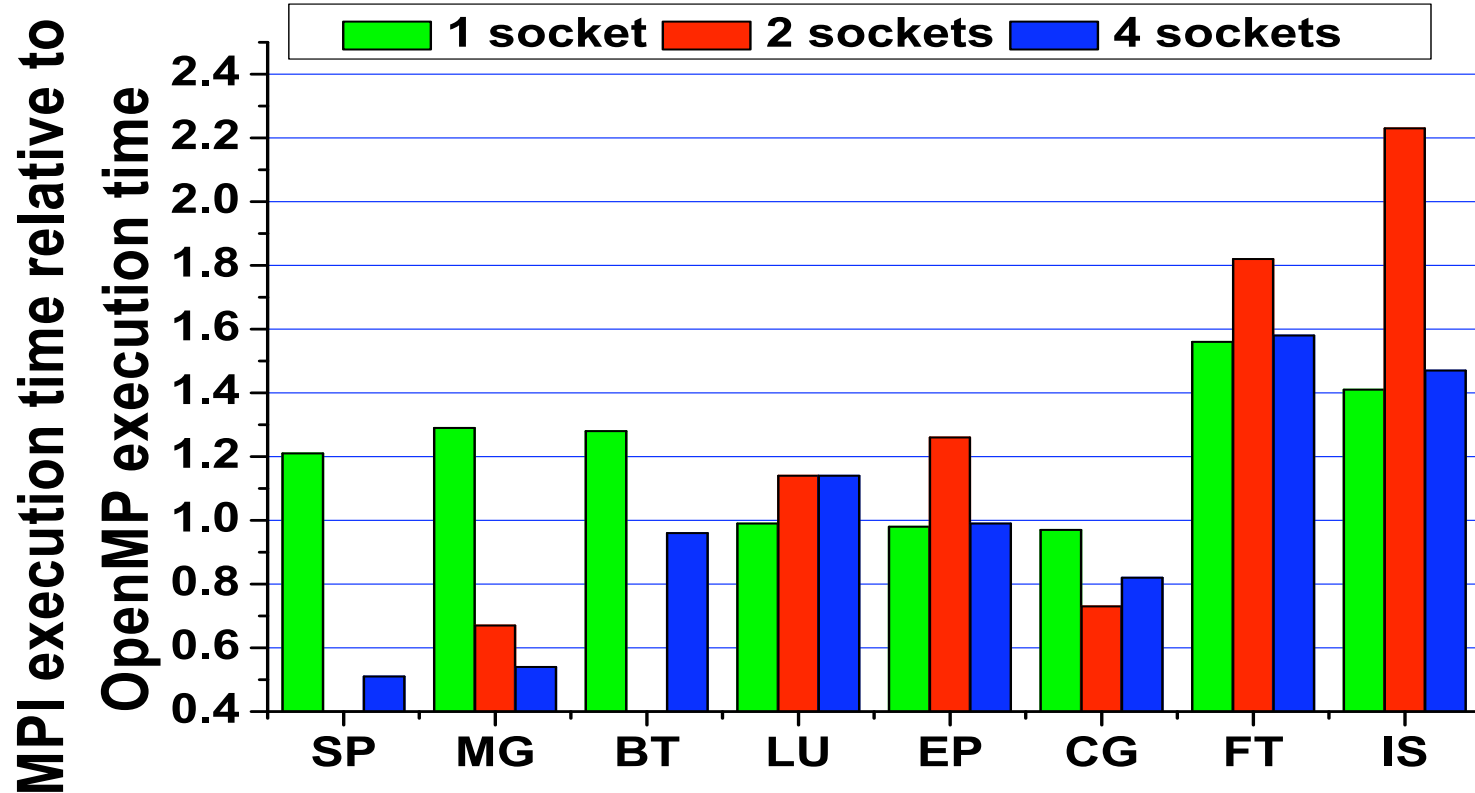
❖ < 1 ➜ MPI is better

- ❖ Single socket performance is OK (KVM and Xen, matches performance expectations)
- ❖ Multi-socket UMA performance is:
- ❖ High performance degradation when VMs span multiple NUMA domains:
  - KVM on average 40%
  - Xen on average 233%

- ❖ VMWare and HyperV
  - Limited number of vcpu per guest – node confinement
  - Restrictions in reporting performance in addition to lack of source code.



Intel Xeon E5530 - NUMA

Percentage increase in execution time compared with host

1 socket    2 sockets

MG    LU    EP    CG    FT    IS

NAS NPB-3.3 MPI

worse



Xen 4.0 – 4 sockets

Percentage increase in execution time compared with host

AMD Opteron 8350 - NUMA
Intel Xeon E7310 - UMA

SP    MG    BT    LU    EP    CG    FT    IS

NAS NBP 3.3 MPI

❖ Xen – GrantTables

❖ KVM - Base shared memory is PCI-based IOMEM driver (an extended version of ivshmem) driver.


❖ Severe restrictions on sizes – MPI works, UPC not

❖ Breaks migration ? What else?

❖ Does not work for OpenMP

❖ OpenMPI is based on Open Component architecture.

❖ Communication is done through communication components that are chosen based on runtime condition.

❖ Shared memory BTL is higher priority (higher exclusivity) transport layer than all other network (only less than self).

❖ Each processor tries to find all transport modules (BTLs) that it can use to reach each destination processors. The highest exclusivity BTL win the registration competition.



1- MPI communication within a node uses shared memory (using sm BTL)

2- MPI communication across nodes uses the fastest available network card (using one of the tcp, IB, ... BTLs).