# KVM as The NFV Hypervisor

Jun Nakajima

Contributors:

Mesut Ergin, Yunhong Jiang, Krishna Murthy, James Tsai, Wei Wang, Huawei Xie, Yang Zhang

(intel) Software

Intel OpenSource TECHNOLOGY CENTER

# Legal Disclaimer

# Agenda

- KVM Enhancements for NFV at OPNFV
- Deterministic Execution and Minimal Latency
- Inter-VM communication: vhost-user-shmem

Classical Network Appliance Approach

Message Router
CDN
Session Border Controller
WAN Acceleration

DPI
Firewall
Carrier Grade NAT
Tester/QoE monitor

SGSN/GGSN
PE Router
BRAS
Radio Access Network Nodes

Independent Software Vendors

Virtual Appliance

Orchestrated, automatic & remote install.

Standard High Volume Servers

Standard High Volume Storage

Standard High Volume Ethernet Switches

Network Virtualisation Approach

**Software**

**IT**

**Virtualization**

**Standard High Volume Servers**

**ETSI's Vision**
**European Telecommunications Standards Institute**

4

# Architecture Framework



Management and Orchestration

OSS/BSS

Os-Ma

NFV Orchestrator

Or-Vnfm

Virtual Network Function

Ve-Vnfm

VNF 1    VNF 2    VNF 3

VNF Manager(s)

Service, VNF and Infrastructure Description

Vi-Vnfm

NFV Infrastructure

NFVI

Virtual Computing    Virtual Storage    Virtual Network

Virtualisation Layer

Vl-Ha

Nf-Vi

Virtualised Infrastructure Manager(s)

Or-Vi

Hardware resources

Computing Hardware    Storage Hardware    Network Hardware

KVM, Containers, ...

openstack

Execution reference points    Other reference points    Main NFV reference points

OPEN DAYLIGHT

# KVM is Crucial to OPNFV

**OPNFV**

AN OPEN PLATFORM
TO ACCELERATE NFV
A Linux Foundation Collaborative Project

Upstream Projects:

Linux

KVM

openstack™

OPEN DAYLIGHT

…

## PLATINUM MEMBERS

at&t · BROCADE · 中国移动 China Mobile · CISCO · DELL · EMC² · ERICSSON

hp · HUAWEI · IBM · intel · JUNIPER NETWORKS · NEC · NOKIA

NTT docomo · redhat · TELECOM ITALIA · vodafone · ZTE

## SILVER MEMBERS

6WIND · ADVA Optical Networking · Alcatel·Lucent · ALTERA · ARM · Array Networks · Brain4Net · BROADCOM

CableLabs · ubuntu Supported by Canonical · CAVIUM · CenturyLink · CertusNet · ciena · CITRIX · ClearPath networks

ConteXtream · Coriant · cyan · Dialogic · dorado SOFTWARE · ENEA · H3C · ixia

kt · Metaswitch Networks · midokura · MIRANTIS · ooredoo · orange · OVERTURE · QOSMOS

sandvine · SK telecom · Sonus · SPIRENT · Sprint · Stratus Technologies · vmware · WIND RIVER

XILINX ALL PROGRAMMABLE.

# Project: NFV Hypervisors-KVM  OPNFV

1. Minimal Interrupt latency variation for data plane VNFs (Virtual Network Function)
2. Inter-VM Communication
3. Fast Live Migration

**Developers from:**

CAVIUM  HUAWEI  NOKIA  WIND RIVER

中国移动 China Mobile  intel  redhat  ZTE

https://wiki.opnfv.org/nfv-kvm

# Deterministic Execution and Minimal Latency

# Causes of Latency Variation
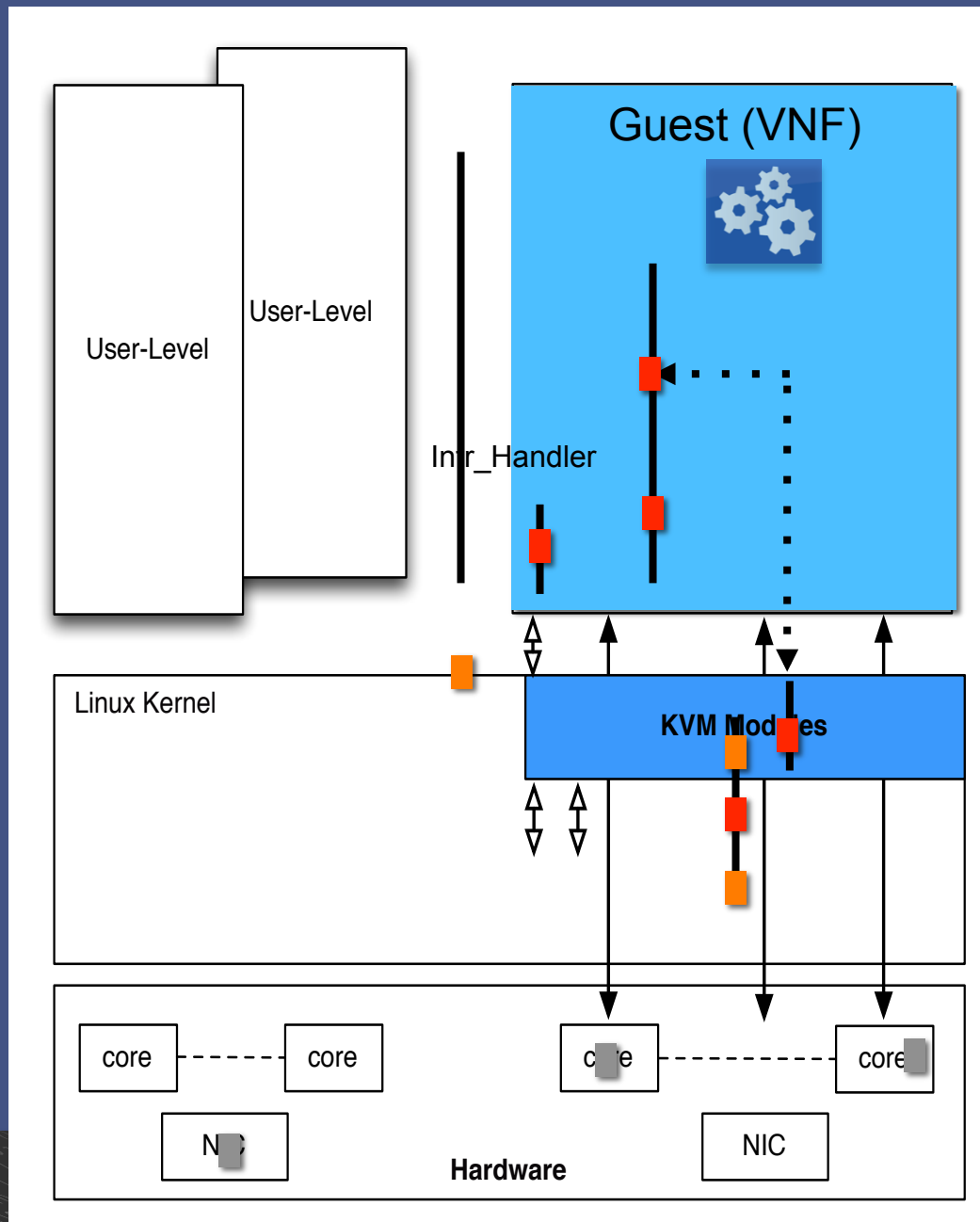
**Asynchronous Events**

Interrupts, VM Exits, Cache/TLB Misses

**Software**

Spin Locks, Loops, Scheduling, Exit to user-level

**Hardware/Firmware**

SMI, Power Management, NIC



Guest (VNF)

User-Level

User-Level

User-Level

Intr_Handler

Linux Kernel

KVM Modules

core — core   core — core

NIC    NIC

**Hardware**

# Solutions ✅

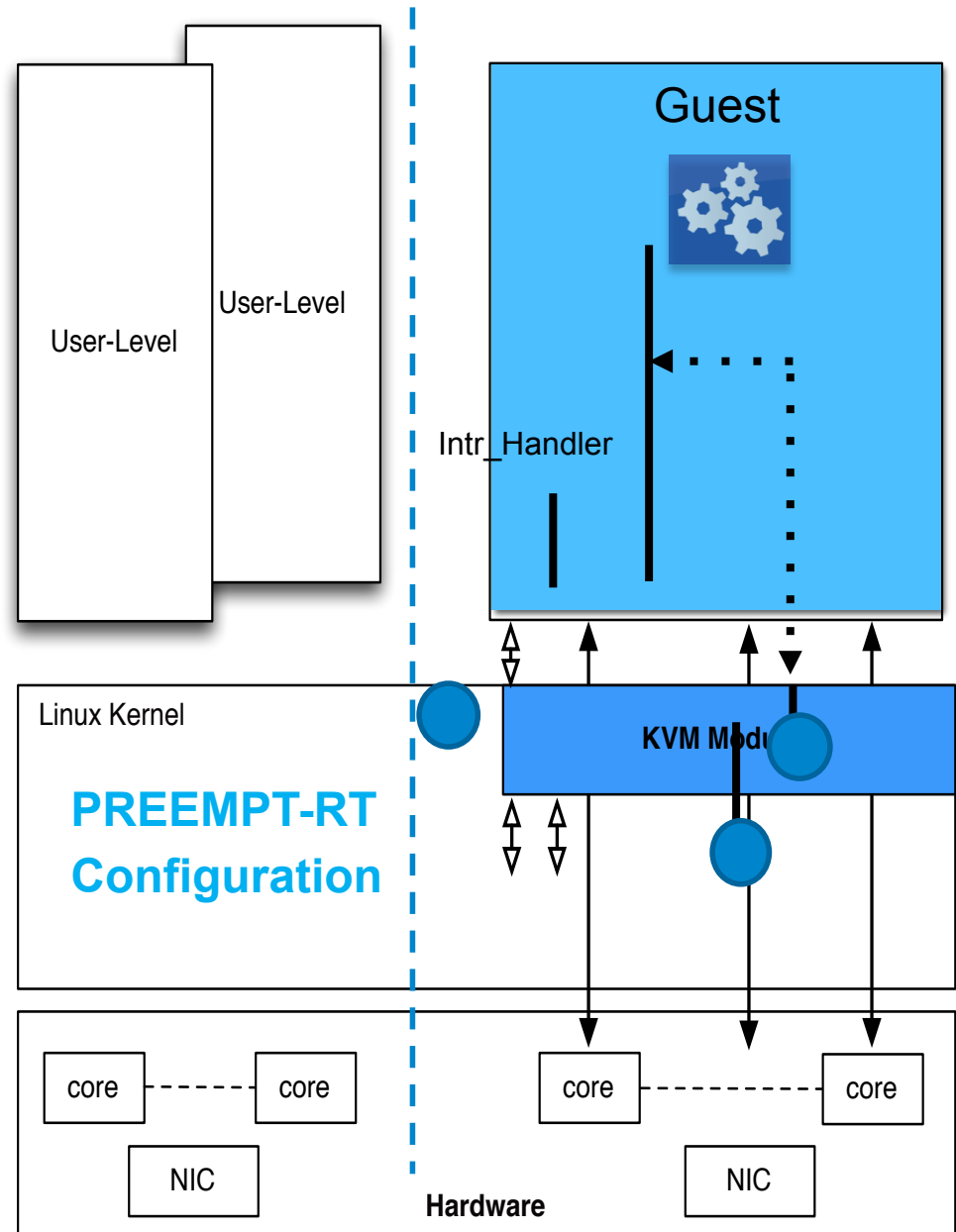**Excusive/Static Allocation**

Soft "Partitioning", CPU Binding, Huge Pages

**Software**

PREEMPT-RT Linux, Code inspection, testing/measurements
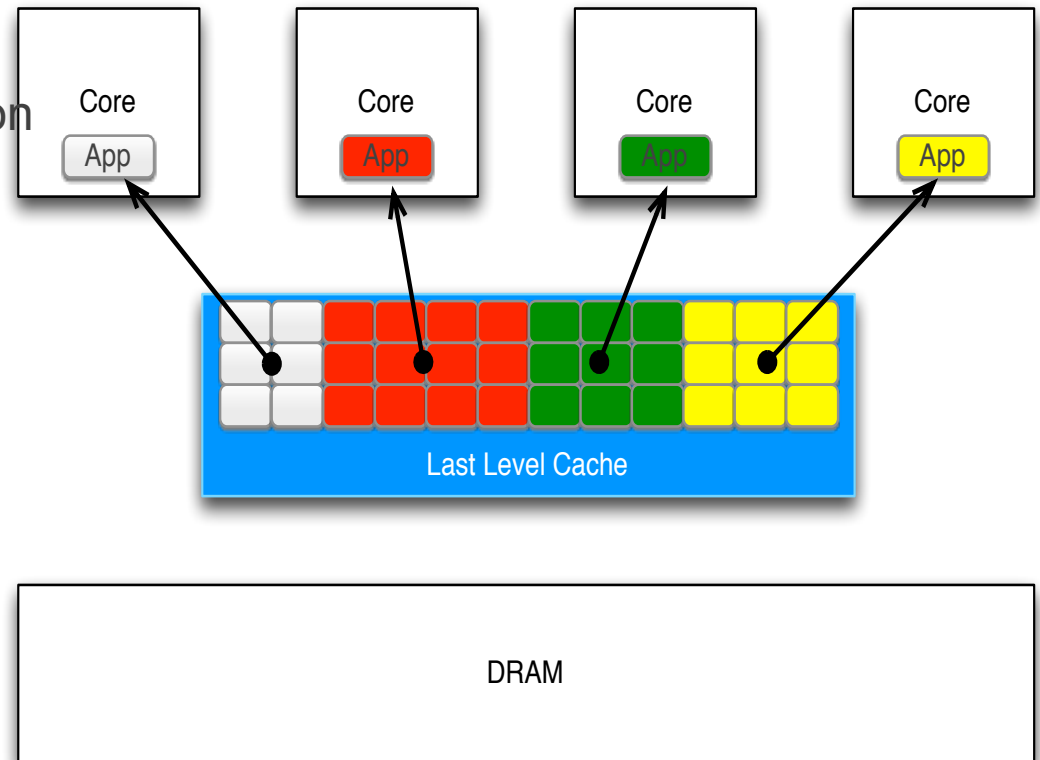
**Hardware Technologies**

Cache Allocation Technology, Advanced VT features

User-Level

User-Level

Guest

Intr_Handler

Linux Kernel

KVM Module

**PREEMPT-RT Configuration**

core – – – core

NIC

core – – – core

NIC

**Hardware**

# Cache Allocation Technology

- Last Level Cache partitioning mechanism enabling the separation of an application
- VMs can be isolated to increase determinism
- Having limited cache is still better than "unlimited cache and noisy neighbors"

| Core | Core | Core | Core |
|------|------|------|------|
| App  | App  | App  | App  |

Last Level Cache

DRAM

CAT is supported on the following 6 SKUs for Intel Xeon processor E5 v3 family: E5-2658 v3, E5-2658A v3, E5-2648L v3, E5-2628L v3, E5-2618L v3, and E5-2608L v3 and Intel(R) Xeon(R) processor D family.
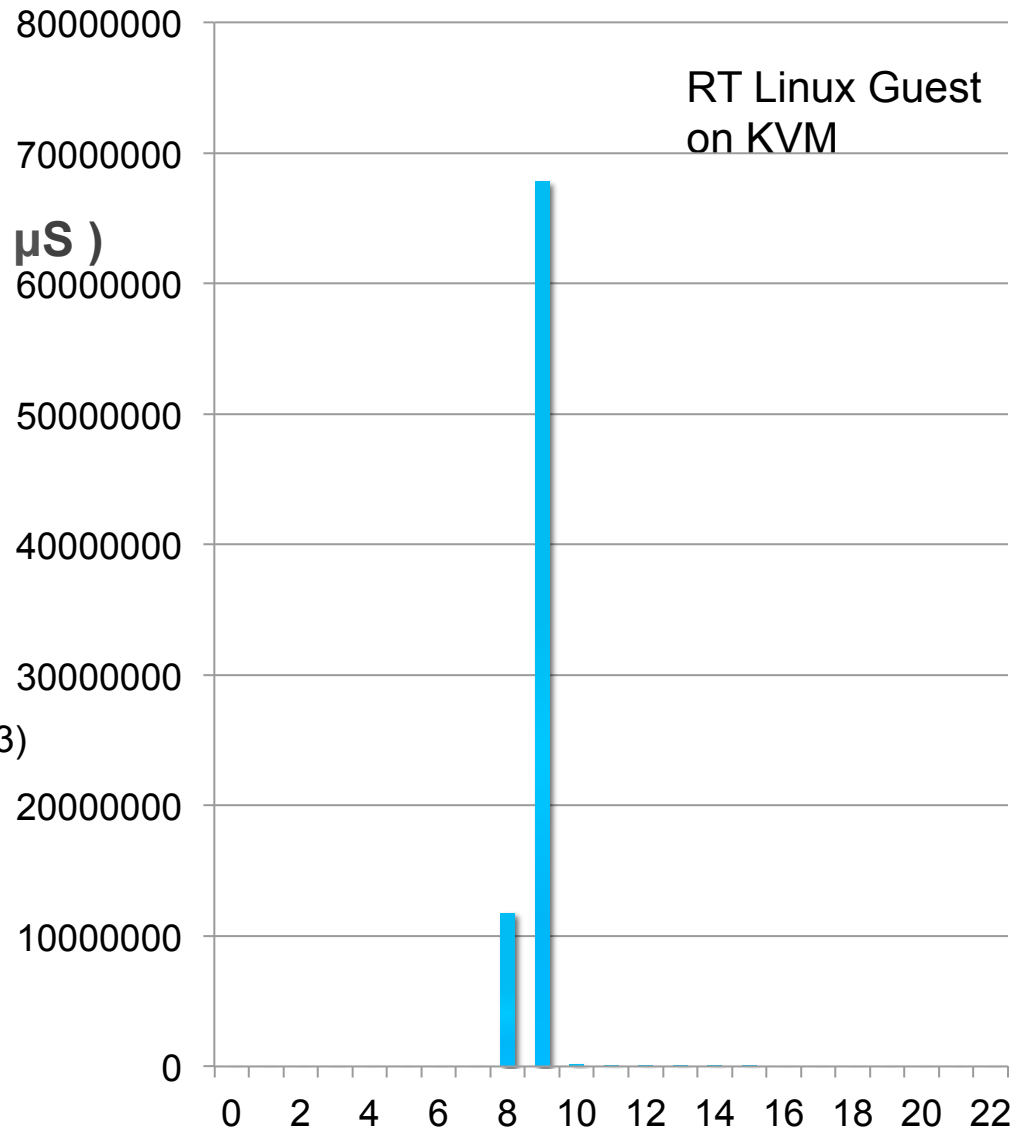
# Latency Data 1: Cyclictest

## Histogram

**Cyclic Test in Guest: Latency (in µS )**

- Min: 7
- Avg: 9
- Max: 16

Latency  Occurrences

```
000007 000003
000008 11757562
000009 67812652
000010 159222
000011 069100
000012 011004
000013 000379
000014 000207
000015 000049
000016 000005
```

99.69%
(Total #: 79,810,183)

Host: Linux with RT patches

RT Linux Guest
on KVM

Latency (µs)

Intel Software

Intel OpenSource TECHNOLOGY CENTER

# Latency Data 2: Latency from Periodic External Interrupts

**Latency from periodic external interrupt:**

- Time delta from interrupt occurrence to invocation of interrupt handler in guest (unit: in µS)

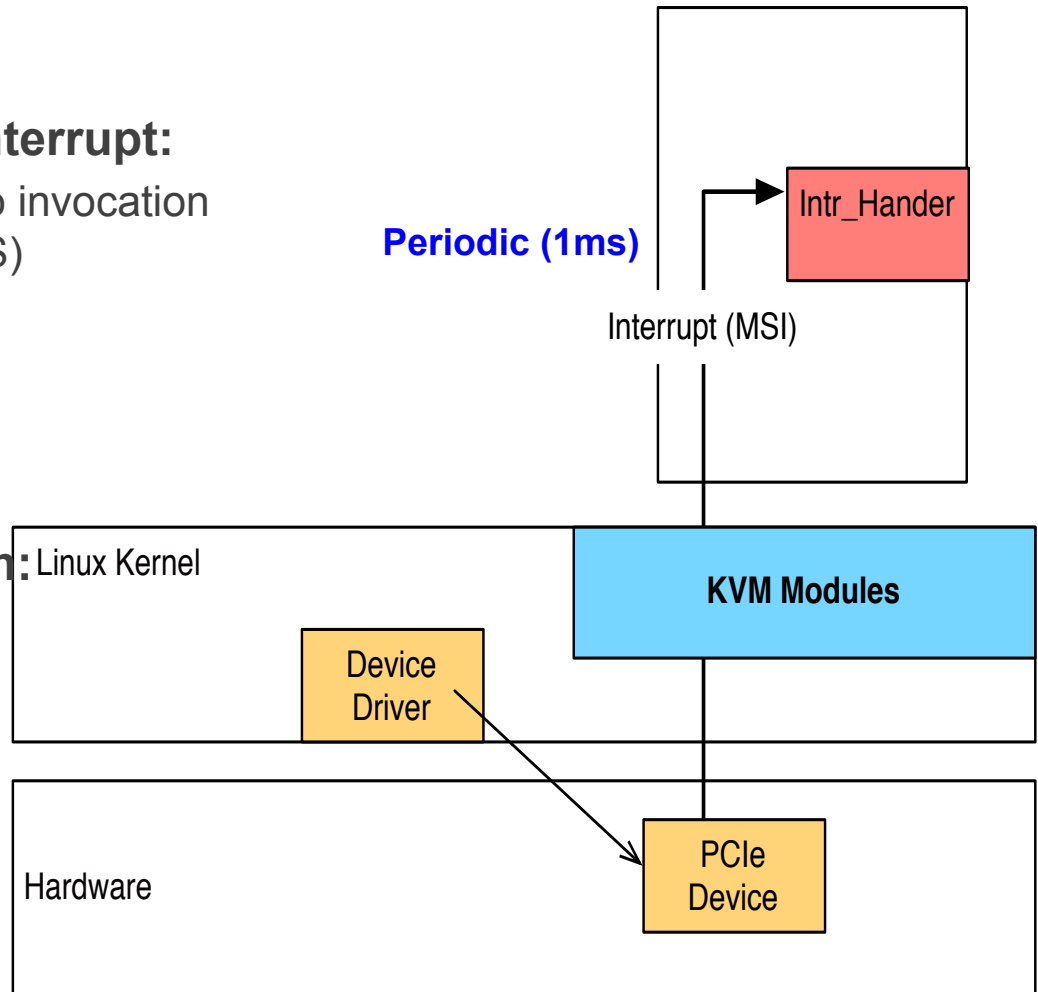  Min: 3.98

  Avg: 4.42

  Max: 9.10
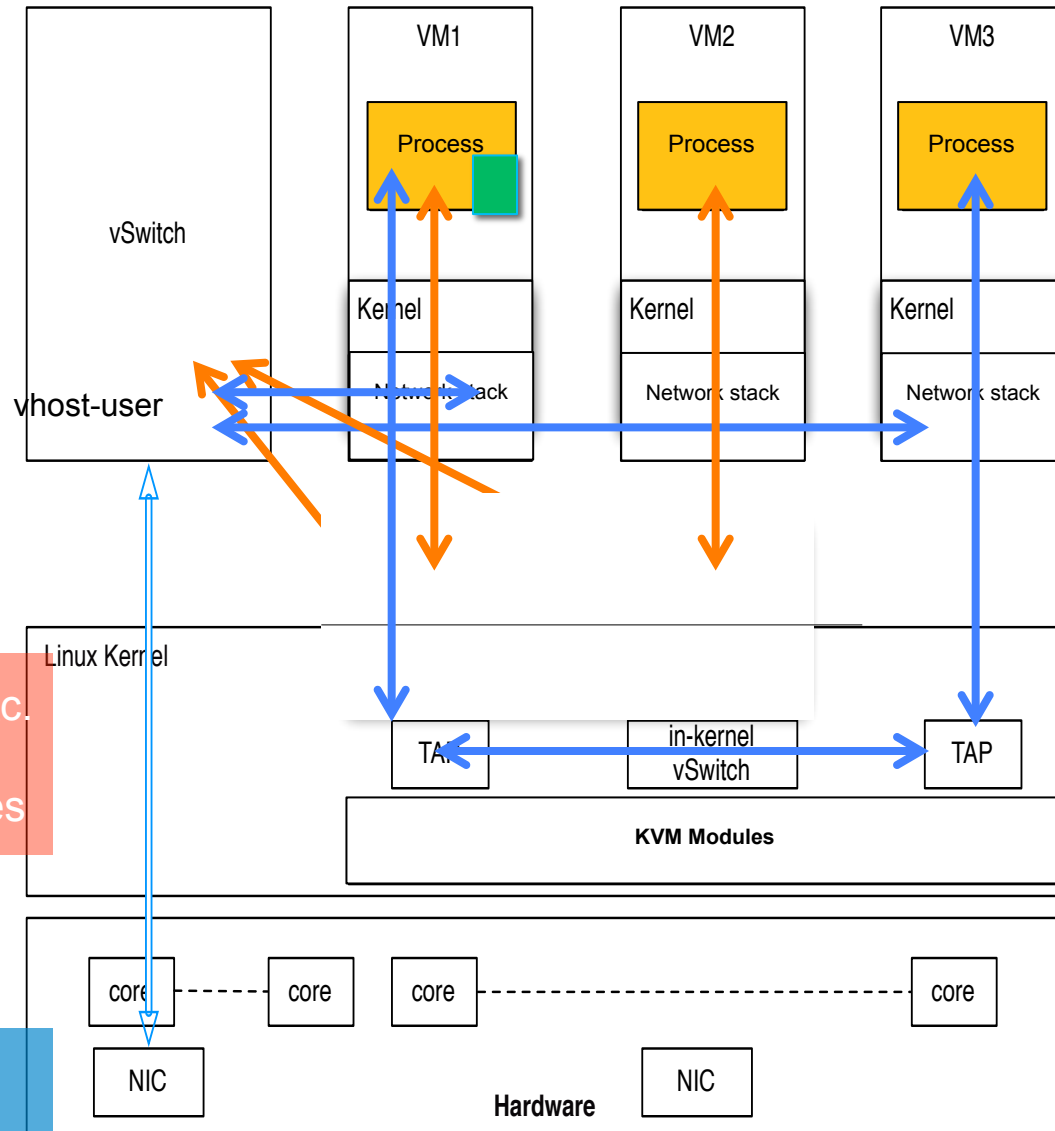
**Expecting even better results with:**

- Posted Interrupts and
- CAT (Cache Allocation Technology)

**Periodic (1ms)**

Intr_Hander

Interrupt (MSI)

Linux Kernel

KVM Modules

Device Driver

Hardware

PCIe Device

# Inter-VM Communication

# Communication Models

vSwitch

vhost-user

**Specific API**

- Shared memory, etc.
- Directly used by particular Processes

**Networking API**

- Use vSwitch in hypervisor
- Generic

VM1
Process
Kernel
Network stack

VM2
Process
Kernel
Network stack

VM3
Process
Kernel
Network stack

Linux Kernel

TAP

in-kernel vSwitch

TAP

**KVM Modules**

core --- core    core -------- core
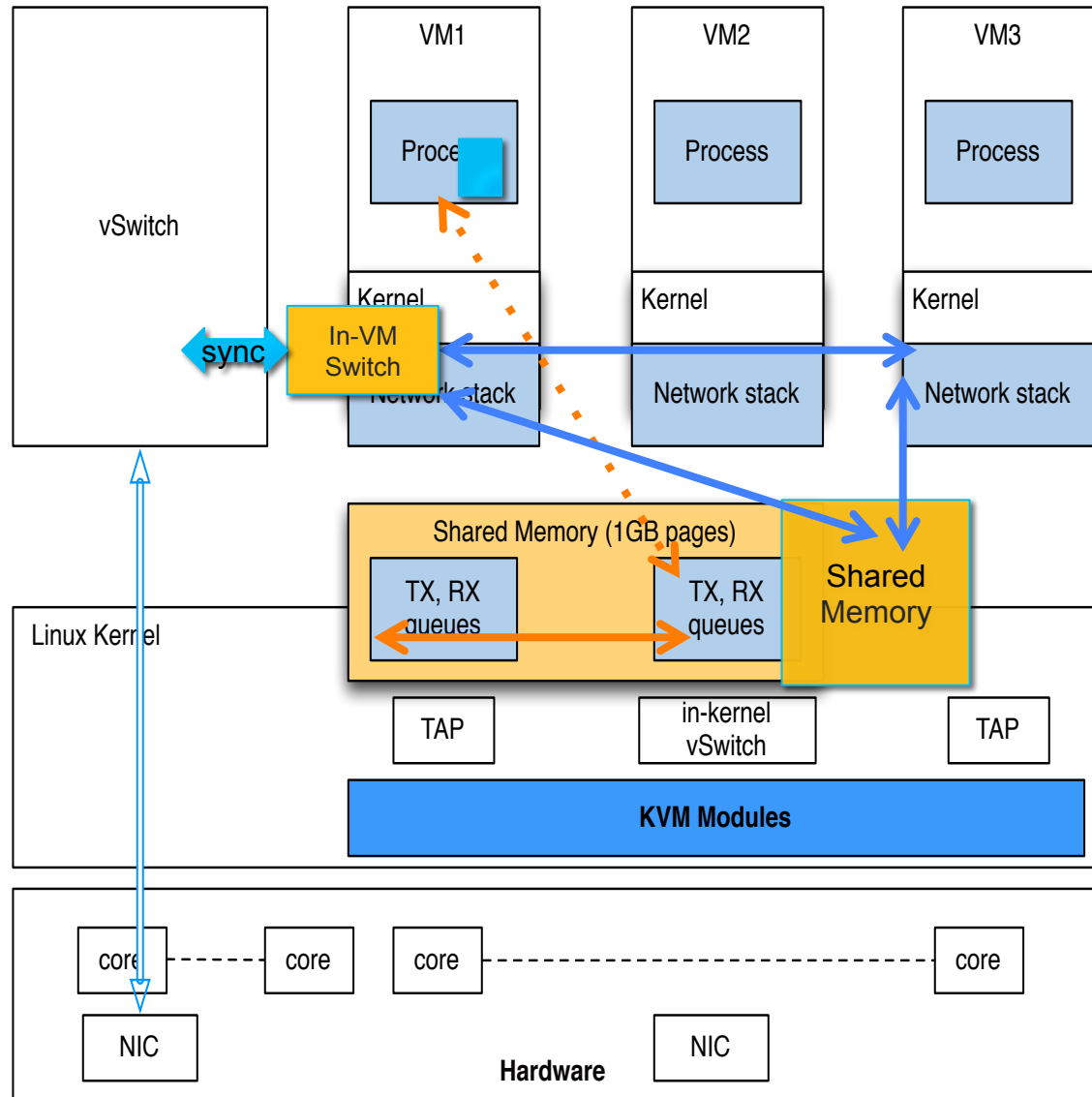
NIC

NIC

**Hardware**

# Fast Paths: Inter-VM Communication

## Specific API

- Need to improve security when using shared memory

## Networking API

- Access Destination VM memory
- Use In-VM Switch



vSwitch

VM1

Process

Kernel

In-VM Switch

sync

Network stack

VM2

Process

Kernel

Network stack

VM3

Process

Kernel

Network stack

Shared Memory (1GB pages)

TX, RX queues

TX, RX queues

Shared Memory

Linux Kernel

TAP

in-kernel vSwitch

TAP

**KVM Modules**

core

core

core

core

NIC

NIC

**Hardware**

# Implementing Inter-VM Communication: vhost-user-shmem

# Goals

- Add fast-paths in VMs as optimized inter-VM communication
- Maintain consistent flow table entries in VMs
- Enable protected access to the destination VM or shared memory
  - Open the Window when needed
  - Close it immediately when done



vSwitch (e.g. OVS-dpdk, Snabb Switch)

VM1

Process

Process

API Libraries

virtio-net

If (in-VM Switch is present) Go to Fast Path

In-VM Switch

Flow Table Entries

In-VM Switch

Inter-VM data transfer

vhost-user

VM2

virtio-net

R X | T X | In-VM Switch

VM3

virtio-net

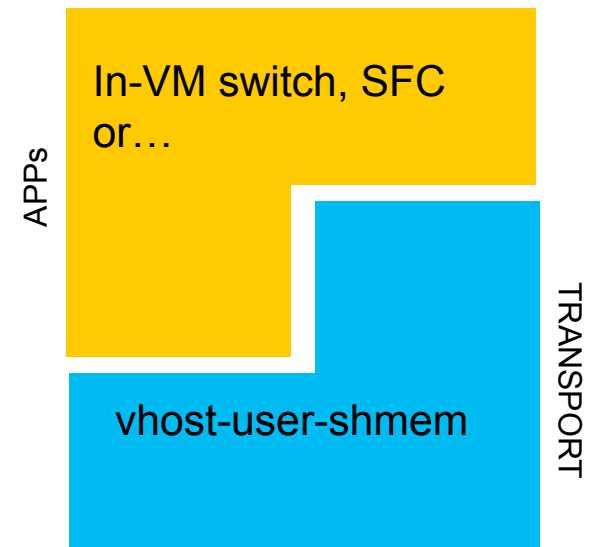R X | T X | In-VM Switch

Shared Memory

NIC

# Clean Design Objectives

**Extend vhost-user as transport mechanism over shared memory/virtqueues:**

- Deliver packets to another guest's virtio device/virtqueue directly
- Provide memory mapping (GPAs), protected access, destination addressing

**Build innovative high-performance networking applications, e.g:**
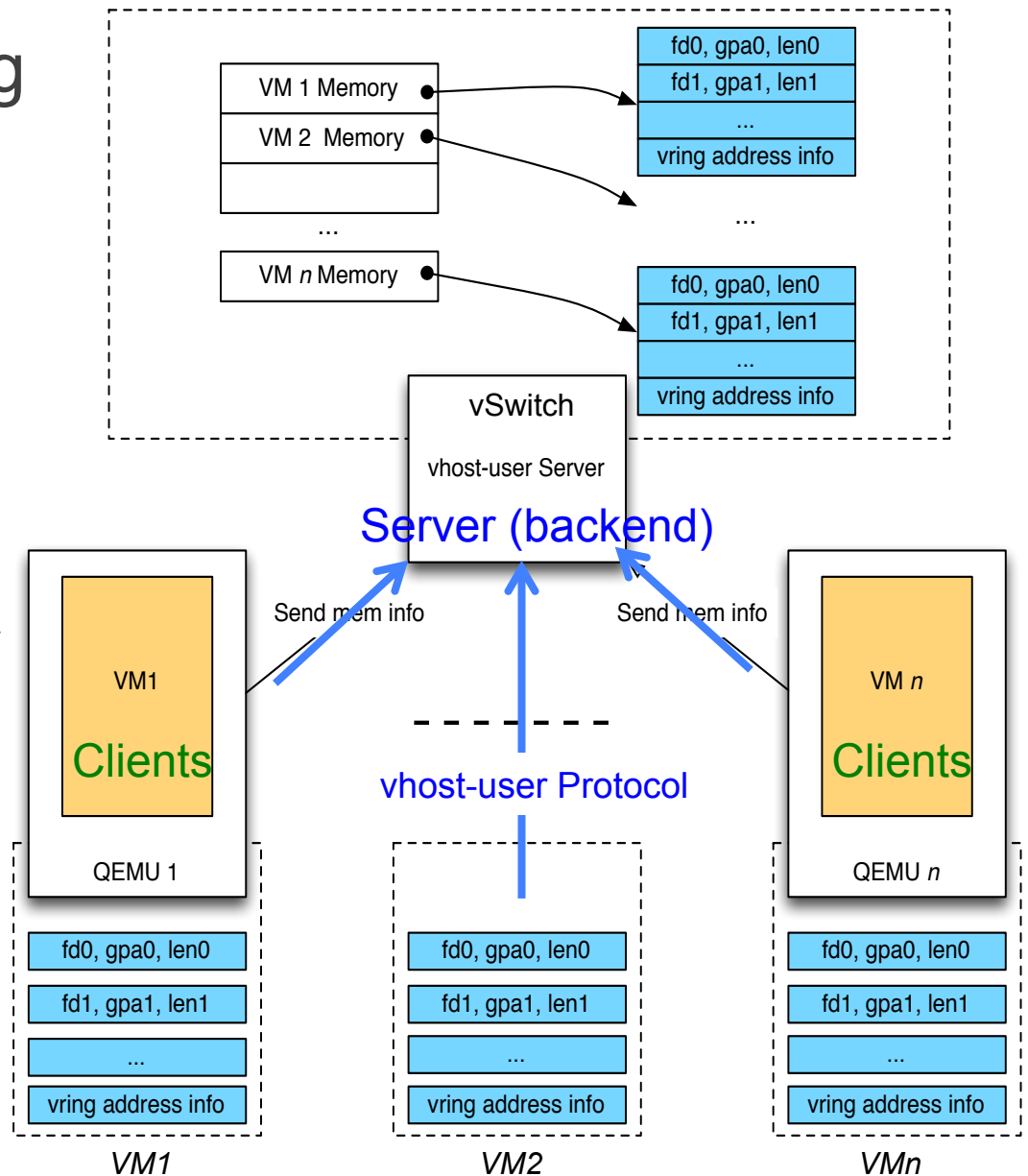
1. In-VM switch as a fast cached-datapath for the full-blown virtual switch
2. Lightweight and fast Service Function Chaining
3. Next big NFV app you are developing

APPs

In-VM switch, SFC or…

vhost-user-shmem

TRANSPORT

# Shared Memory Using vhost-user Server

**vhost-user server (backend) has sufficient info and capability to host shared memory:**

- Gather mem info to access virtuques from vhost-user clients (QEMUs)
- It can allocate its own memory for sharing purposes
  - E.g. large pages shared by guests (like ivshmem)

# Extending it for Inter-VM Communication

- **vhost-user server (backend)** becomes a client
  - Send mem info to QEMUs
  - QEMU extends memory regions
- Allows **vhost-user clients** to access their virtqueues each other
- Provides **vhost-user clients** with shared memory

# Simple Example: VM1 and VM2



vSwitch

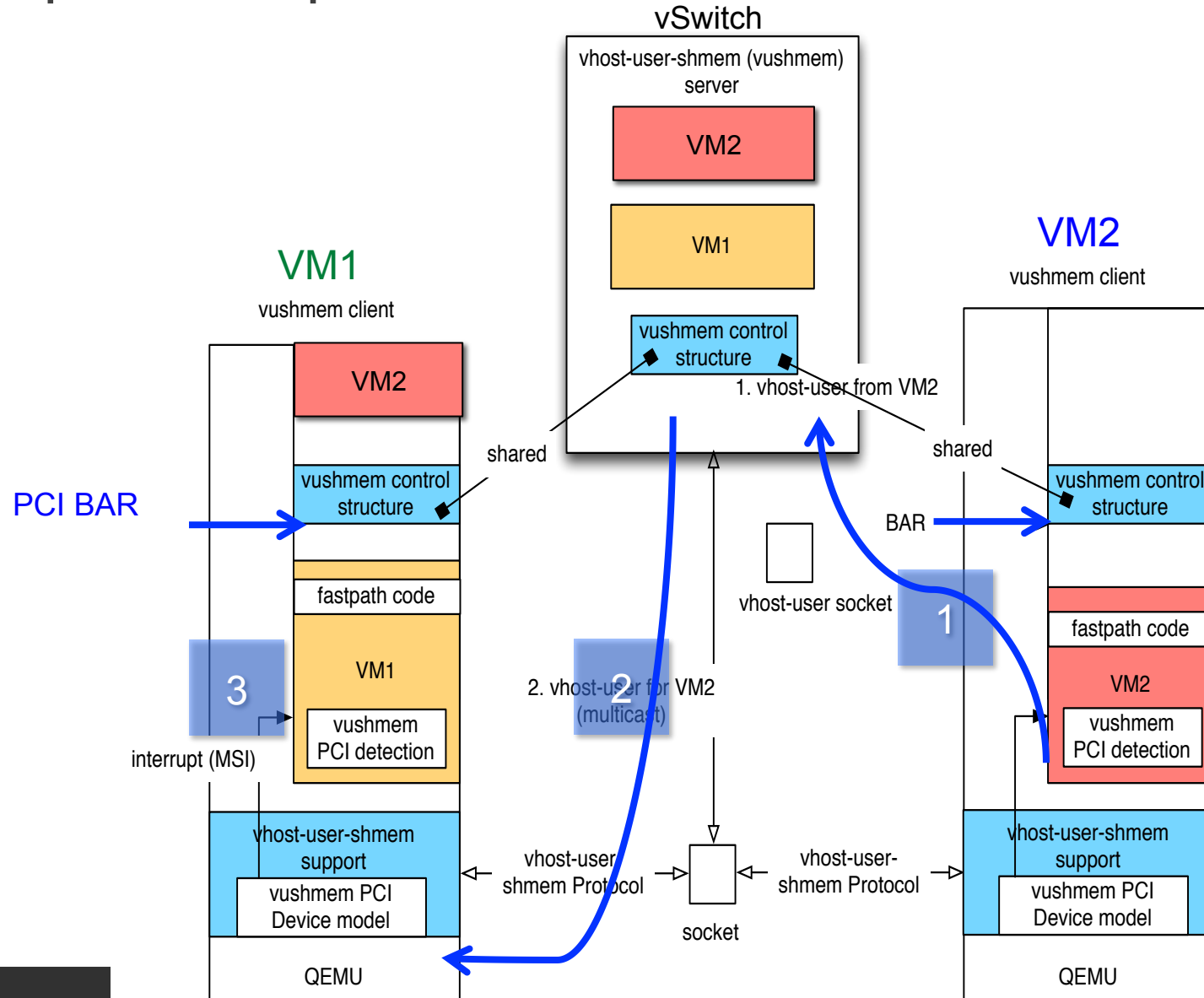vhost-user-shmem (vushmem) server

VM2

VM1

vushmem control structure

1. vhost-user from VM2

**VM1**
vushmem client

VM2

vushmem control structure

PCI BAR

fastpath code

VM1

vushmem PCI detection

interrupt (MSI)

vhost-user-shmem support

vushmem PCI Device model

QEMU

shared

vhost-user socket

2. vhost-user for VM2 (multicast)

vhost-user-shmem Protocol

socket

vhost-user-shmem Protocol

**VM2**
vushmem client

BAR

vushmem control structure

shared

fastpath code

VM2

vushmem PCI detection

vhost-user-shmem support

vushmem PCI Device model
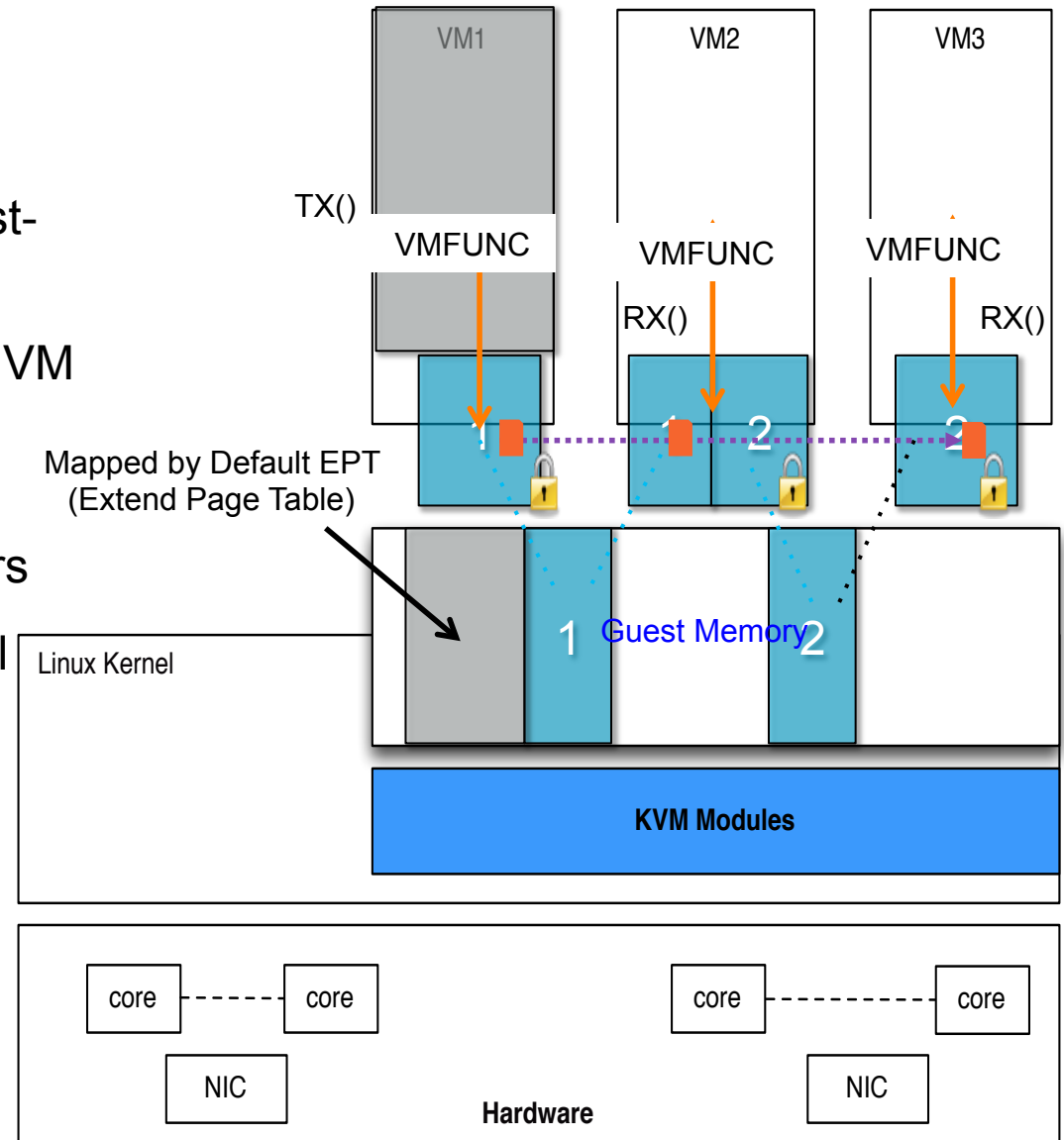
QEMU

1

2

3

# Adding Protected Access

- **Extends memory** to access fast-path channel or destination VM

- **VMFUNC** instruction in VM w/o VM exit

  - #0 (EAX): Switches EPT (Extend Page Table) Pointers

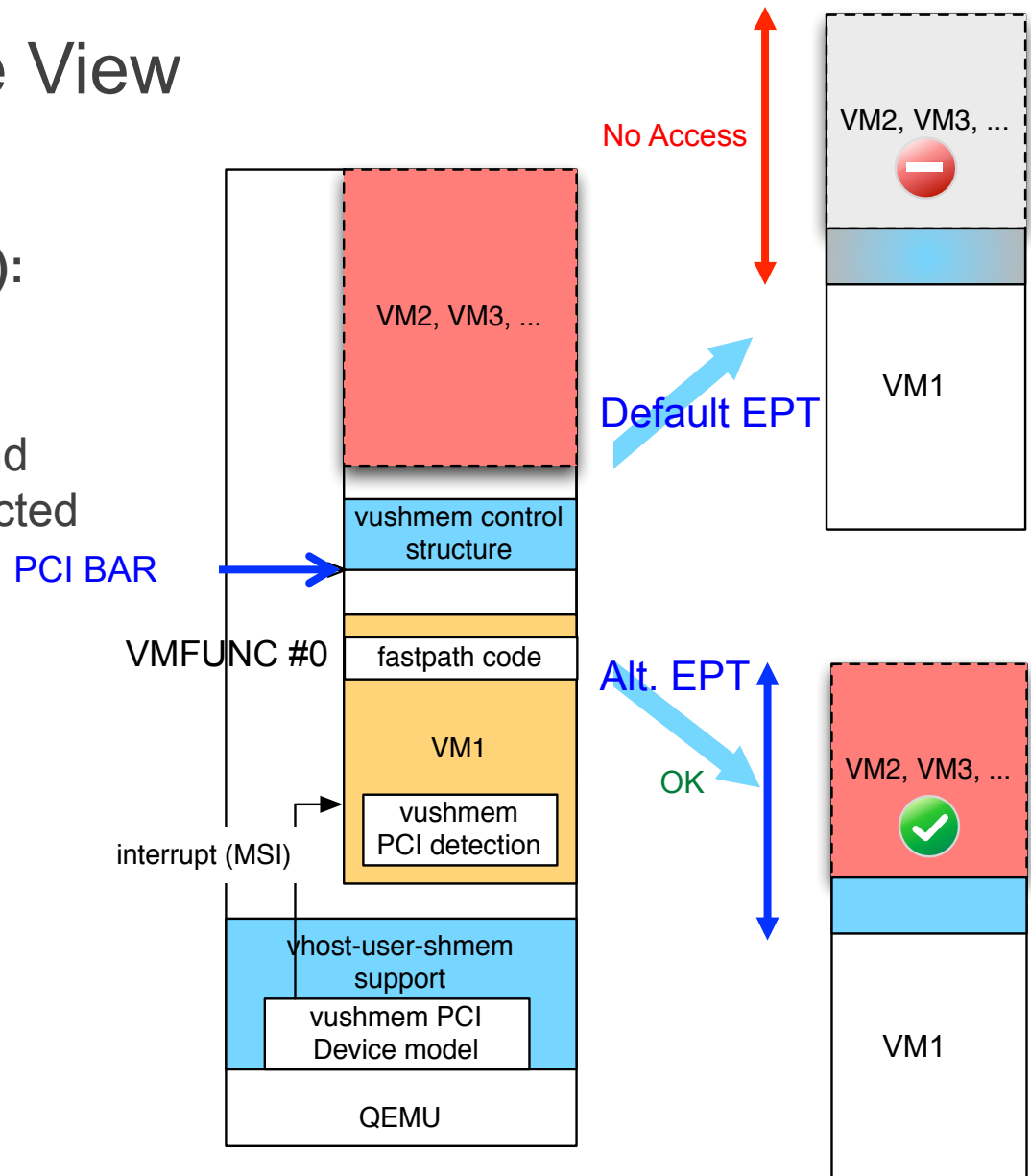  - Alternate EPT has additional translation



VM1   VM2   VM3

TX()

VMFUNC   VMFUNC   VMFUNC

RX()   RX()

Mapped by Default EPT (Extend Page Table)

Linux Kernel

1   Guest Memory   2

KVM Modules

core ---- core      core ---- core

NIC      NIC

**Hardware**

# Adding EPT Alternate View

**Fast Pass Code (Protected Code):**

- Upon VMFUNC #0, EPT View is changed

- Access other shared memory and virtqueues of other VMs in protected fashion

**KVM ioctl options for QEMU to extend Guest Memory:**

1. W/O protection, or

2. W/ protection
   - Extend only in alternate EPT view

No Access

VM2, VM3, ...

VM2, VM3, ...

Default EPT

VM1

vushmem control structure

PCI BAR

VMFUNC #0

fastpath code

Alt. EPT

VM1

vushmem PCI detection

OK

VM2, VM3, ...

interrupt (MSI)

vhost-user-shmem support

vushmem PCI Device model

VM1

QEMU

intel Software

Intel OpenSource TECHNOLOGY CENTER

# Implementing Protected Access



Default EPT

EPT Permission

Full (X, W, R)

```
start_xmit(*skb, *dev) {

...
    send(packets);
}
```

Page Boundary

**Trampoline Page**

- Registered by a trusted entity
- Entry/Exit to/from Trusted Code

Write-Protected (X, -, R)

```
send(*packet) {
...
    VMFUNC #0, EPTP;
    Tx(packets);
    VMFUNC #0, 0
}
```

No Access (-, -, -)

```
Tx(*packet) {
    move_data();
    notifify()
}
```

**Protected Code**

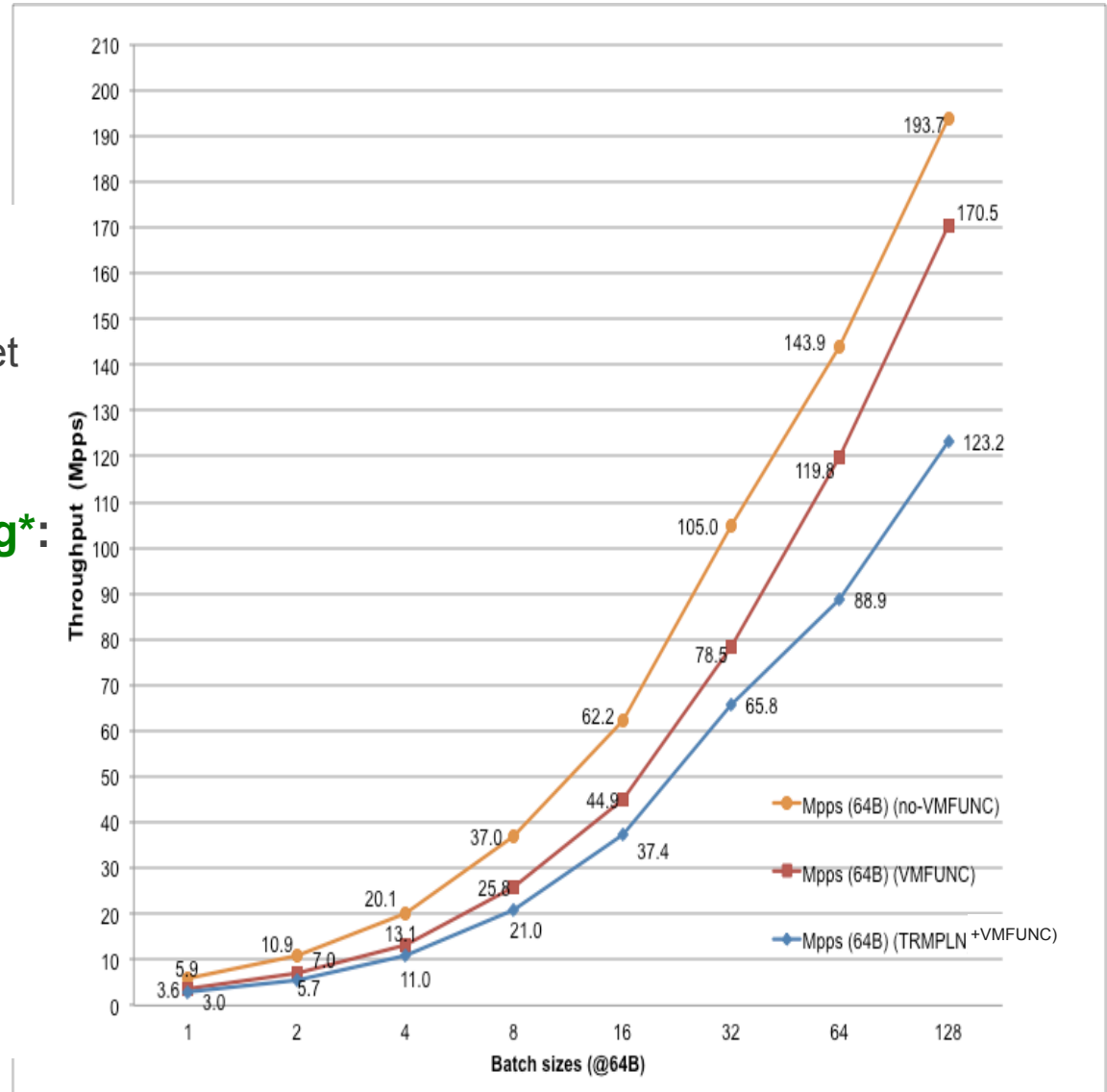- Specified at registration time

Additional pages

# Performance Estimate from PoC

**Measure cost of VMFUNC and Trampoline Code:**

- Transfer 64B packets from virtio-net to another VM (fast path)

**65Mpps with 32-packet batching\*:**

- Same batching size as DPDK

\*Intel internal estimation



Chart: Throughput (Mpps) vs Batch sizes (@64B)

- Mpps (64B) (no-VMFUNC): 5.9, 10.9, 20.1, 37.0, 62.2, 105.0, 143.9, 193.7
- Mpps (64B) (VMFUNC): 3.6, 7.0, 13.1, 25.8, 44.9, 78.5, 119.8, 170.5
- Mpps (64B) (TRMPLN +VMFUNC): 3.0, 5.7, 11.0, 21.0, 37.4, 65.8, 88.9, 123.2

Batch sizes: 1, 2, 4, 8, 16, 32, 64, 128

# Summary

1. Minimal Interrupt latency variation for data plane VNFs (Virtual Network Function)
   - On Track
2. Inter-VM Communication
   - Preliminary performance data from PoC with trampoline code
   - Implementation proposal (vhost-user-shmem) based on vhost-user
3. Fast Live Migration
   - Next presentation

## Join OPNFV Projects!

OPNFV
AN OPEN PLATFORM
TO ACCELERATE NFV
A Linux Foundation Collaborative Project

Q & A