# Host Hypervisor Trace Mining for Virtual Machine Workload Characterization

Hani Nemati, Vahid Azhari,

and Michel Dagenais

Dec 6, 2018

# Outline

## 1

### Introduction

- Agent-less
- Restricted Access
- Old Kernel
- Limited Resources

## 2

### Previously on VM Analysis

- New Tracepoint
- Nested VM
- Wait Analysis
- Critical Path

## 3

### Feature Extraction

- Virtual Interrupt
- CPU Metrics
- Disk Metrics
- Network Metrics

## 4

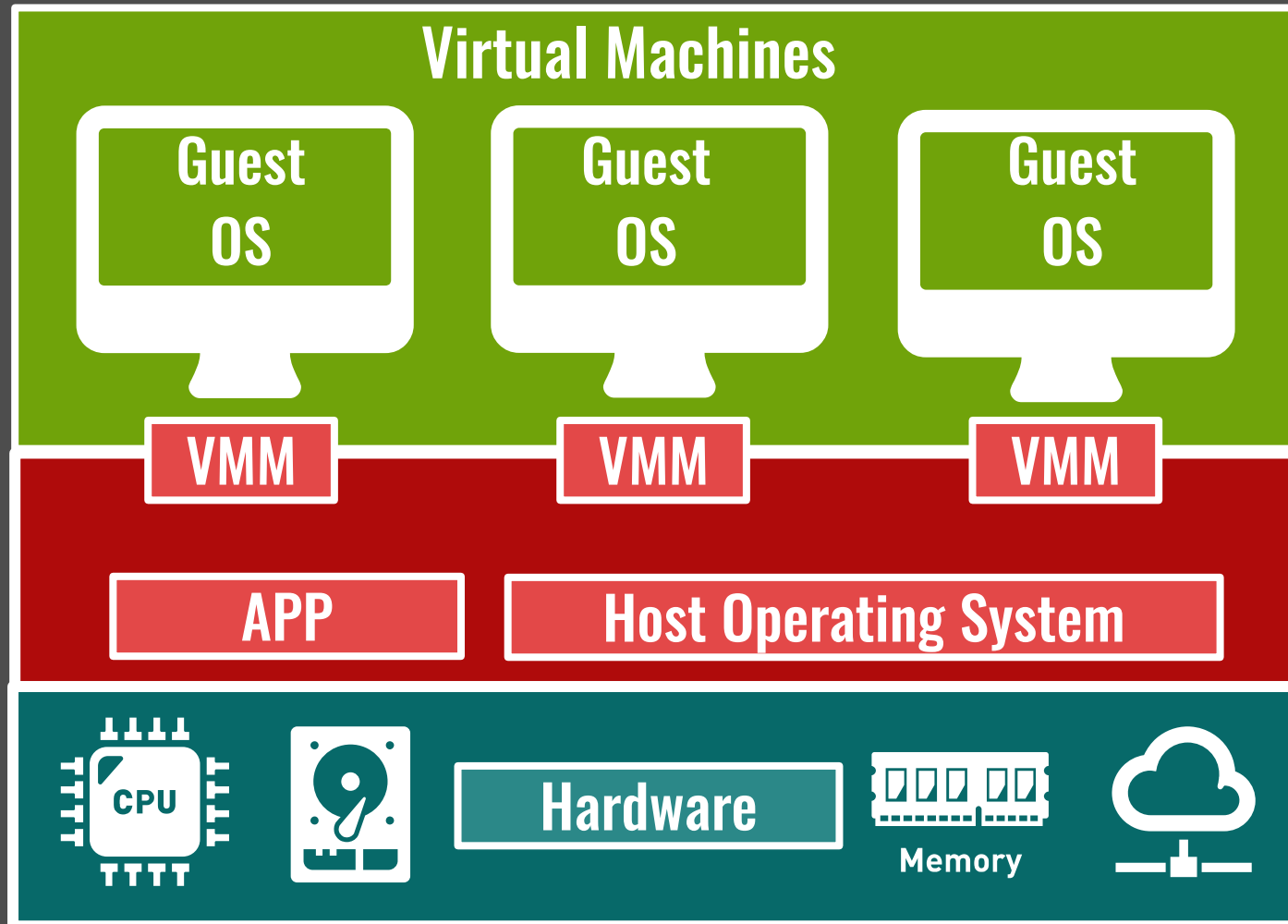### VM Clustering

- K-MEANS
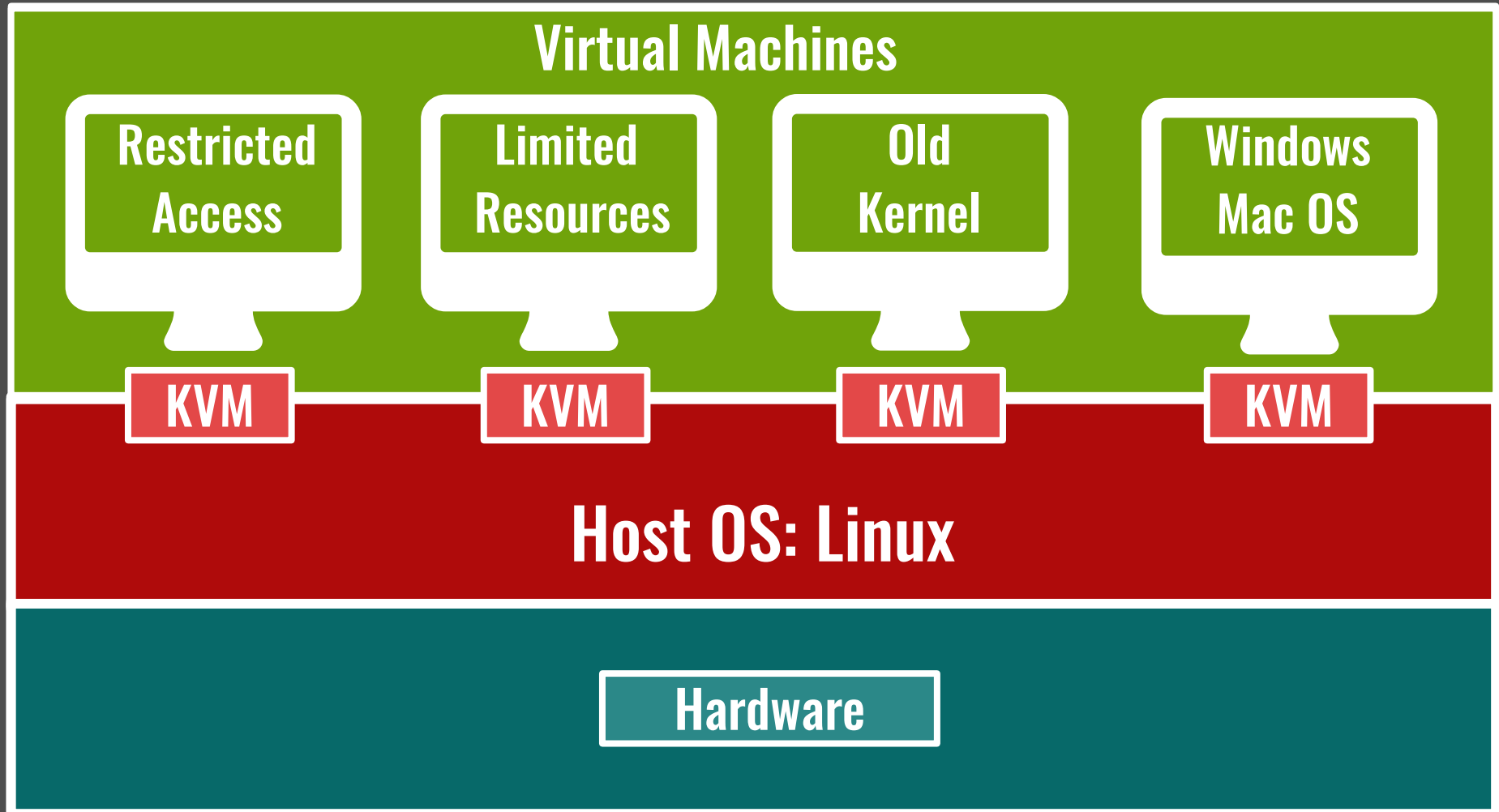- CPU Intensive
- Disk Intensive
- Network Intensive

## 5

### Conclusion & Future Work

- Trace mining
- VM Clustering
- vCPU features
- Process Clustering

# VM Architecture

## Virtual Machines

| Guest OS | Guest OS | Guest OS |
|----------|----------|----------|

VMM    VMM    VMM

| APP | Host Operating System |
|-----|------------------------|

CPU    Hardware    Memory

# Motivation

## Virtual Machines

| Restricted Access | Limited Resources | Old Kernel | Windows Mac OS |

**KVM**    **KVM**    **KVM**    **KVM**

## Host OS: Linux

### Hardware

# Previously on VM Analysis

Step **01-Monitoring**

10 Dec 2015   →   A

- **Fine-grain resource Monitoring**
  Retrieve guest state from VMCS
  Analyzing VM interference
  Disk View, vCPU view , Netrowk view

# Previously on VM Analysis

- ## Resource Overcommitment Detection

  CPU Running State Analysis
  CPU overcommitment Analysis
  Memory overcommitment Analysis

B    5 May 2016

Step 02-Overcommitment

## Investigations

### Use Case – Preemptive Virtual Machine

30 runs of Sysbench to find primes < 10000

virtual Machine CPU view - Without Preemption
CPU 0

Average = 324ms    STD = 5ms

Virtual Machine CPU view - With Preemption
CPU 0

Average = 443ms    STD = 116ms



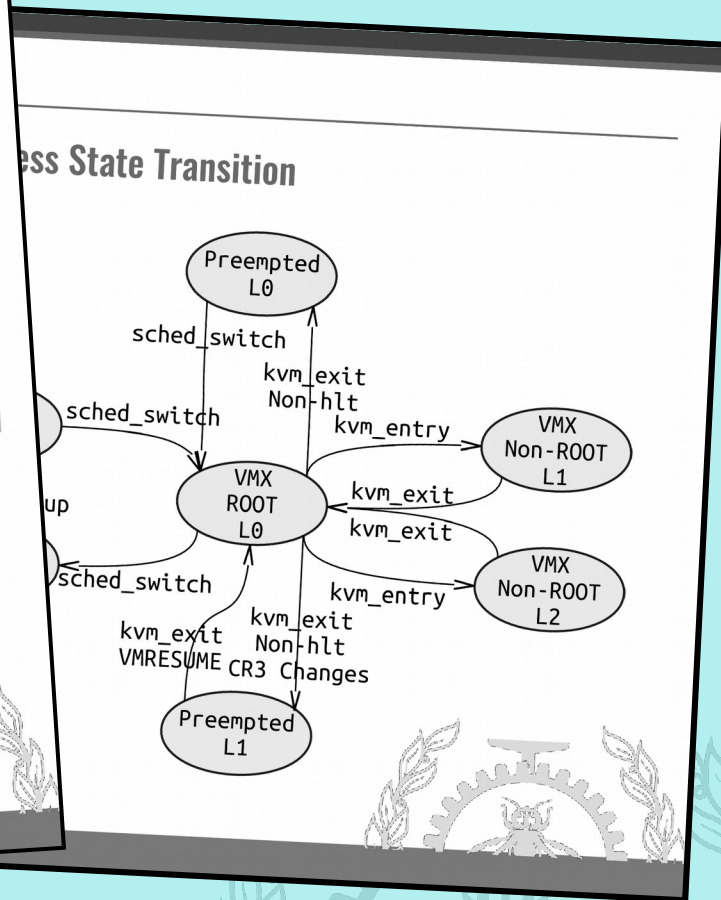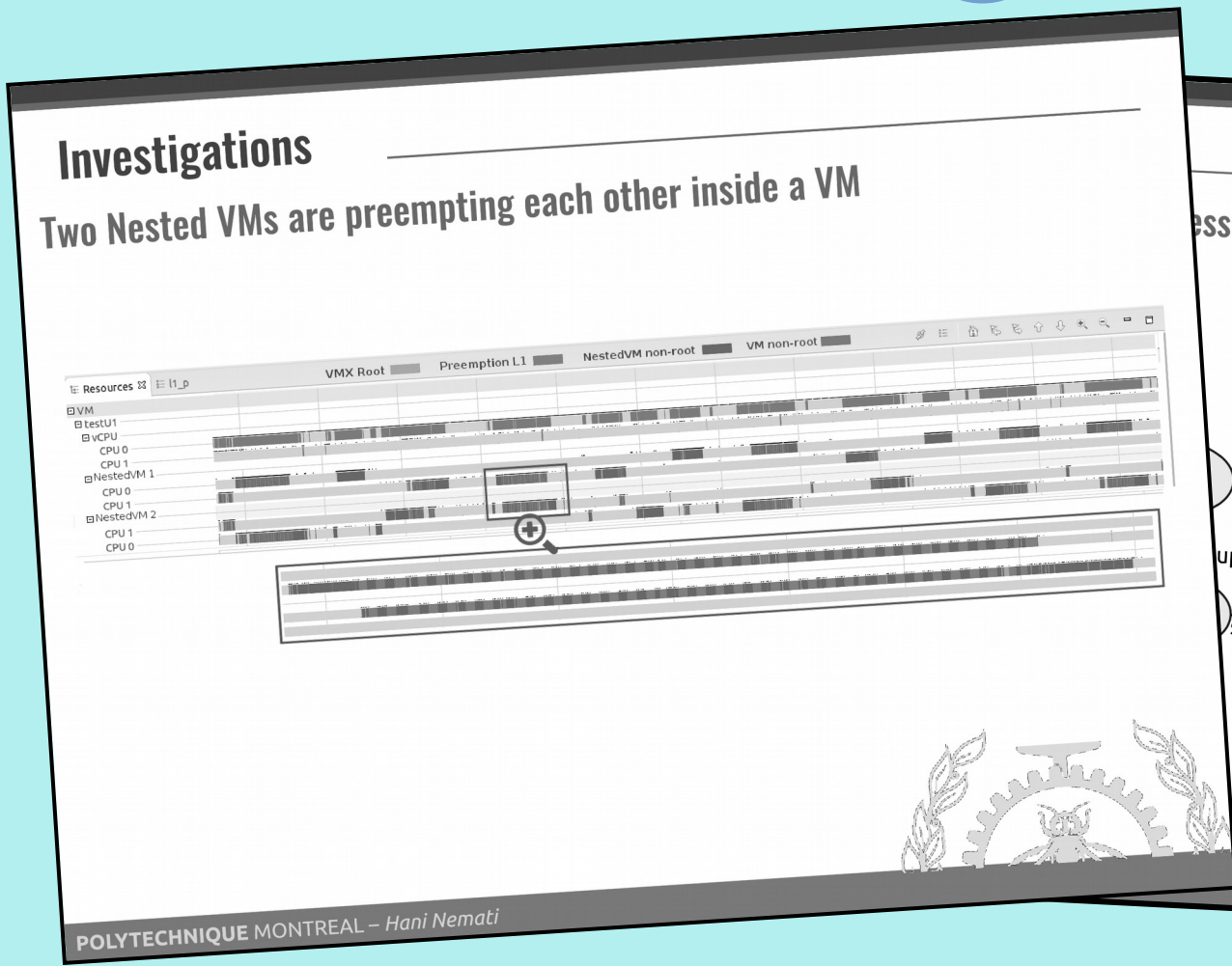POLYTECHNIQUE MONTREAL – Hani Nemati

# Previously on VM Analysis

12 Dec 2016    C

Step **03**-Nested VM

- **Any level Nested VM State Analysis**
  Nested VM vCPU state analysis
  Nested VM overcommitment detection
  Virtualization overhead analysis

## Investigations
Two Nested VMs are preempting each other inside a VM



Resources ≋ ≣ l1_p
- VM
  - testU1
    - vCPU
      - CPU 0
      - CPU 1
  - NestedVM 1
    - CPU 0
    - CPU 1
  - NestedVM 2
    - CPU 1
    - CPU 0

VMX Root    Preemption L1    NestedVM non-root    VM non-root

## ...ess State Transition

# Previously on VM Analysis

- ## vCPU idle State Analysis
  Virtual interrupt injection Analysis
  Analyzing idle State of VM
  eBPF based VM analysis

D  5 May 2017

Step **04**-Wait Analysis



**Investigations**
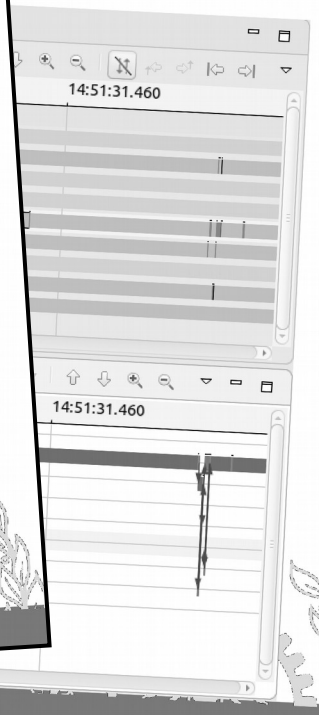
Vec from kvm_inj_virq
CR3 from vcpu_enter_guest

Vec from kvm_inj_virq
CR3 from vcpu_enter_guest

CR3 from vcpu_enter_guest

| | | | | |
| Block | Running | Block | Running | Block |

Main Thread

vCPU Thread — Running | Block | Running | Block | Running

Block I/O Thread — Block | Running | Block | Running | Block

Network Thread — Block

Process #CR3 — Running | Block I/O | Running | Network | Running

```
If (Vec == (Block I/O irq)) {
    Block State = Block I/O State
} else if (Vec == (network irq)) {
    Block State = Network State

}
```

14:51:31.460

14:51:31.460

POLYTECHNIQUE MONTREAL – Hani Nemati

# Previously on VM Analysis

7 Dec 2017    E

Step **05**-Critical Path

- **Execution Path Analysis of VM**
  Contention of Disk
  Cap on virtual resources
  VM boot-up comparison

# Previously on VM Analysis

- ### Hierarchical Distributed Critical Path

  Critical Path for Nested VM
  Critical Path for distributed VMs

F

10 May 2018

Step **06**-Distributed VM

# Previously on VM Analysis

- Hierarchical Distributed Critical Path

**F** — 10 May 2018

Step **06**-Distributed VM

- Execution Path Analysis of VM

**E** — 7 Dec 2017

Step **05**-Critical Path

- vCPU idle State Analysis

**D** — 5 May 2017

Step **04**-Wait Analysis

- Any level Nested VM State Analysis

**C** — 12 Dec 2016

Step **03**-Nested VM

- Resource Overcommitment Detection

**B** — 5 May 2016

Step **02**-Overcommitment

**A** — 10 Dec 2015

Step **01**-Monitoring

- Fine-grain resource Monitoring

**Automation**

How about

# Automation?

# VM Clustering Architecture



**Private** Clouds

**Hybrid** Clouds

**Public** Clouds

*Report Engine*

*Agent-less tracing*
*LTTng*

*VM clustering*
*KMEANS*

*Data Collection*
*Feature Extraction*
*Tracecompass*

*VM Data*

# Disk Metrics

**Control Flow view**

idle     syscall     userspace

main thread: idle | futex | userspace | futex | idle | futex | userspace | futex | idle

vCPU thread: userspace | idle | userspace

IO thread: idle | futex | User space | preadv | User space | futex | idle

① vm_exit =30

② 

③ 

Disk latency

inject_disk_interrupt

④ Wait for handling disk request

⑤

① VM exits from guest mode to host mode to handle disk request

② Qemu main thread submits the disk request

③ Qemu IO thread handles the disk request

④ Two metrics: wait for handling disk request and disk latency

⑤ Virtual disk interrupt is injected to VM

# CPU Metrics

## Control Flow view



1. Virtual timer interrupt is injected to VM

2. A VM process wakes up another VM process

3. Virtual task interrupt is injected to VM

# Network Metrics

Control Flow view  (2)

idle    syscall    userspace

Wait for packet to receive

| vCPU thread | userspace | idle | userspace |
| vhost-net | idle | userspace | idle |

(1)   (3)

inject_network_interrupt

(1)  A packet is received by vhost-net

(2)  Wait for a packet to receive

(3)  Virtual network interrupt is injected to VM

# Experimental Evaluation

## k-means
## Clustering

# Experimental Evaluation



**CPU**
**Intensive**

# Experimental Evaluation



## CPU
## Intensive
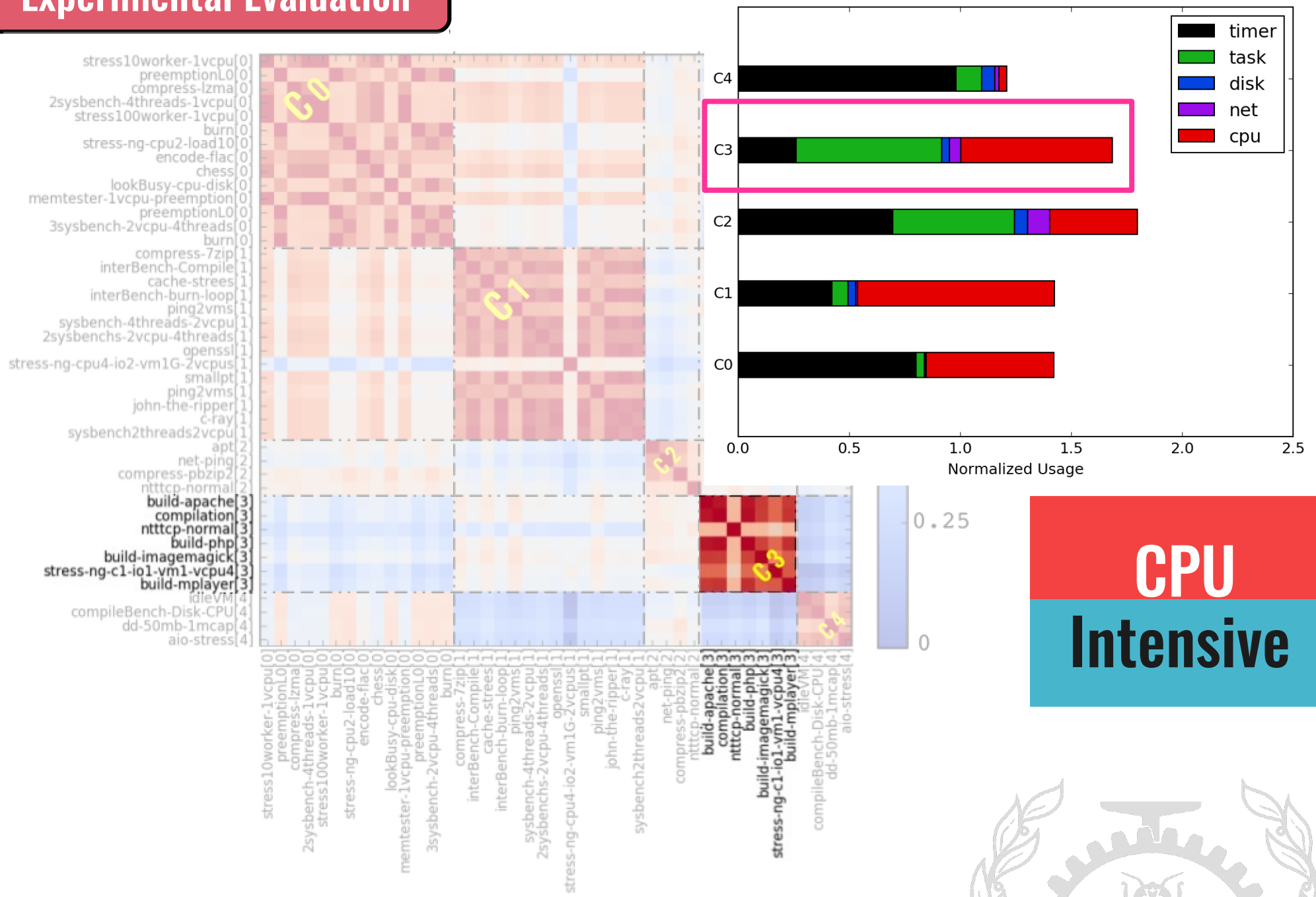
# Experimental Evaluation

## CPU
## Intensive

# Experimental Evaluation

## CPU
## Intensive

# Experimental Evaluation



**CPU Intensive**

# Experimental Evaluation
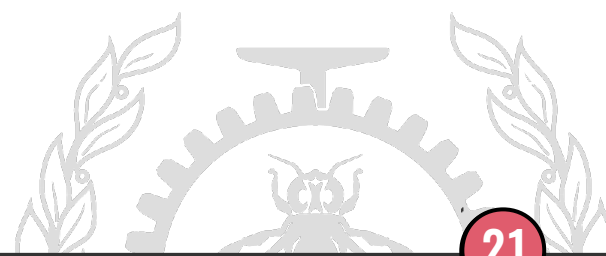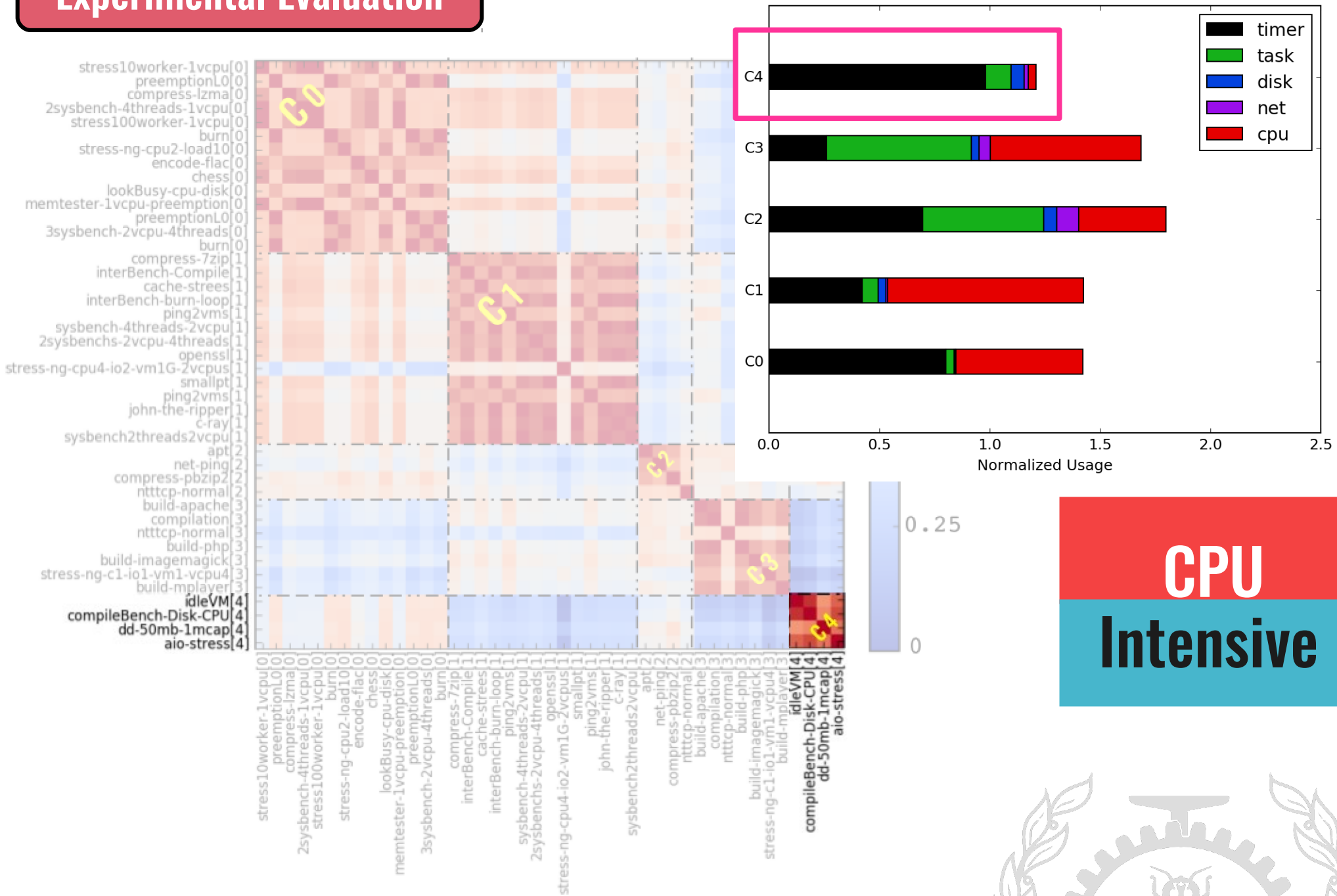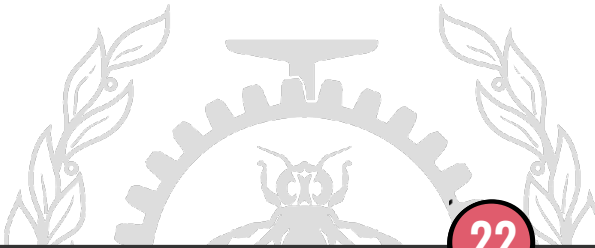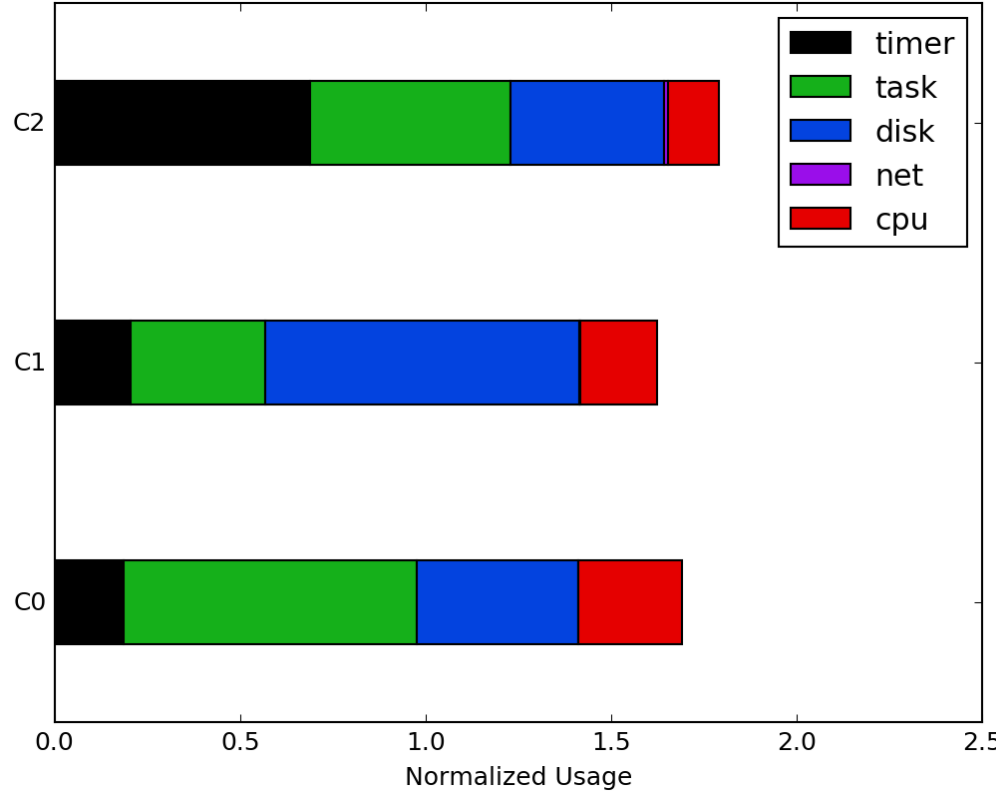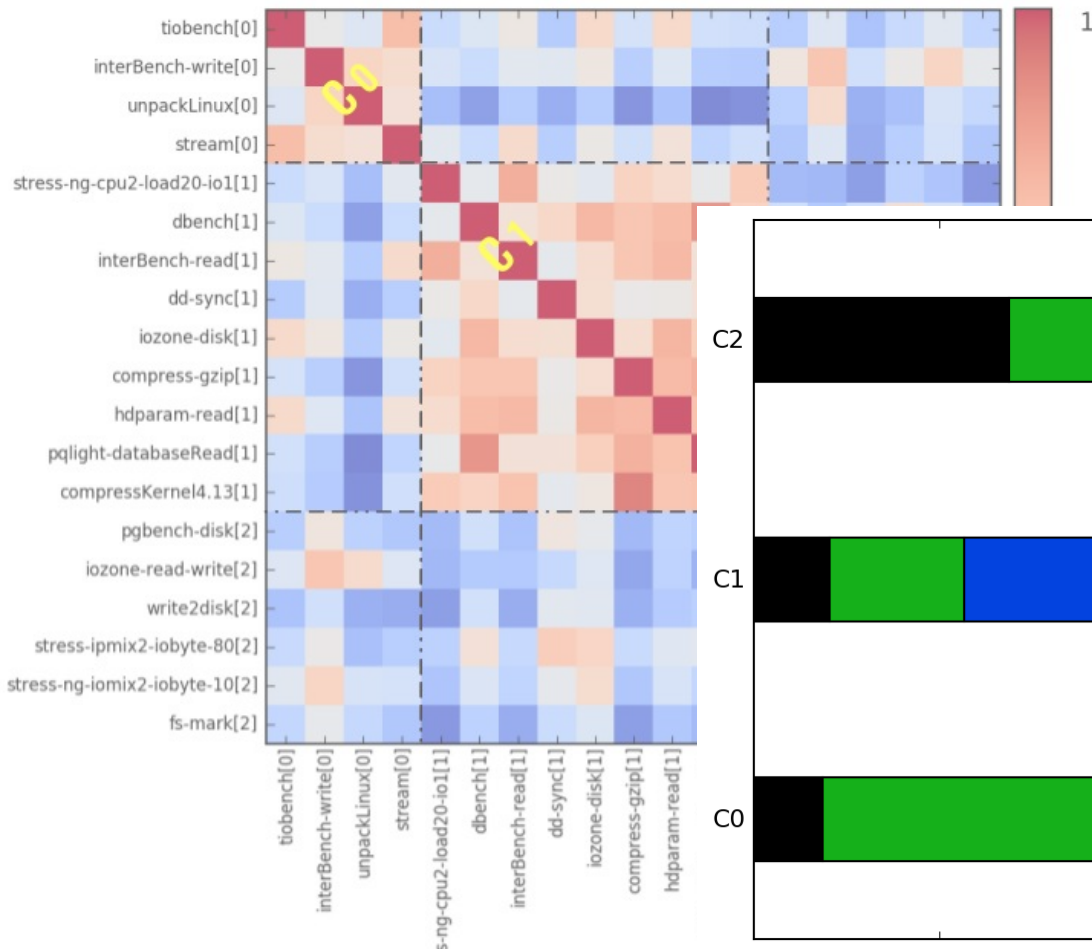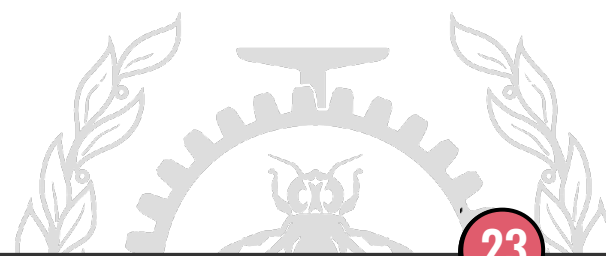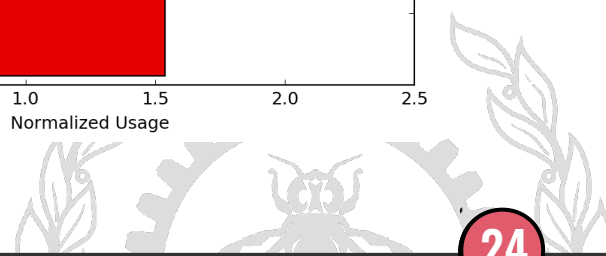


## CPU
## Intensive

# Experimental Evaluation



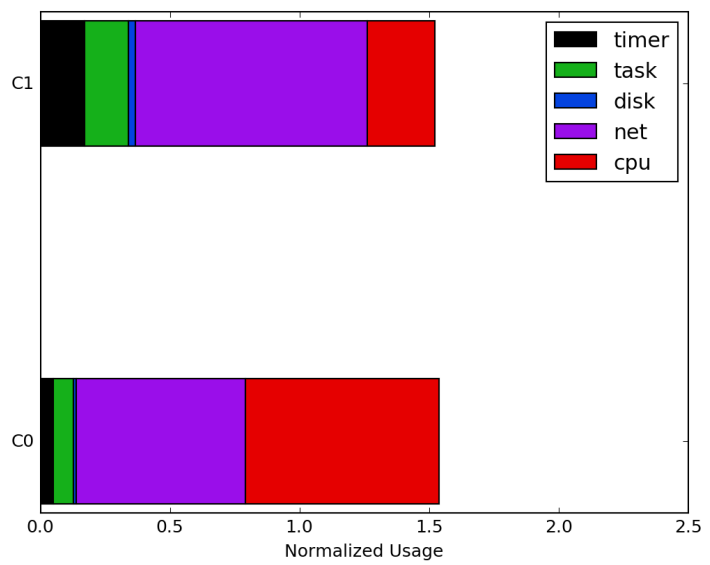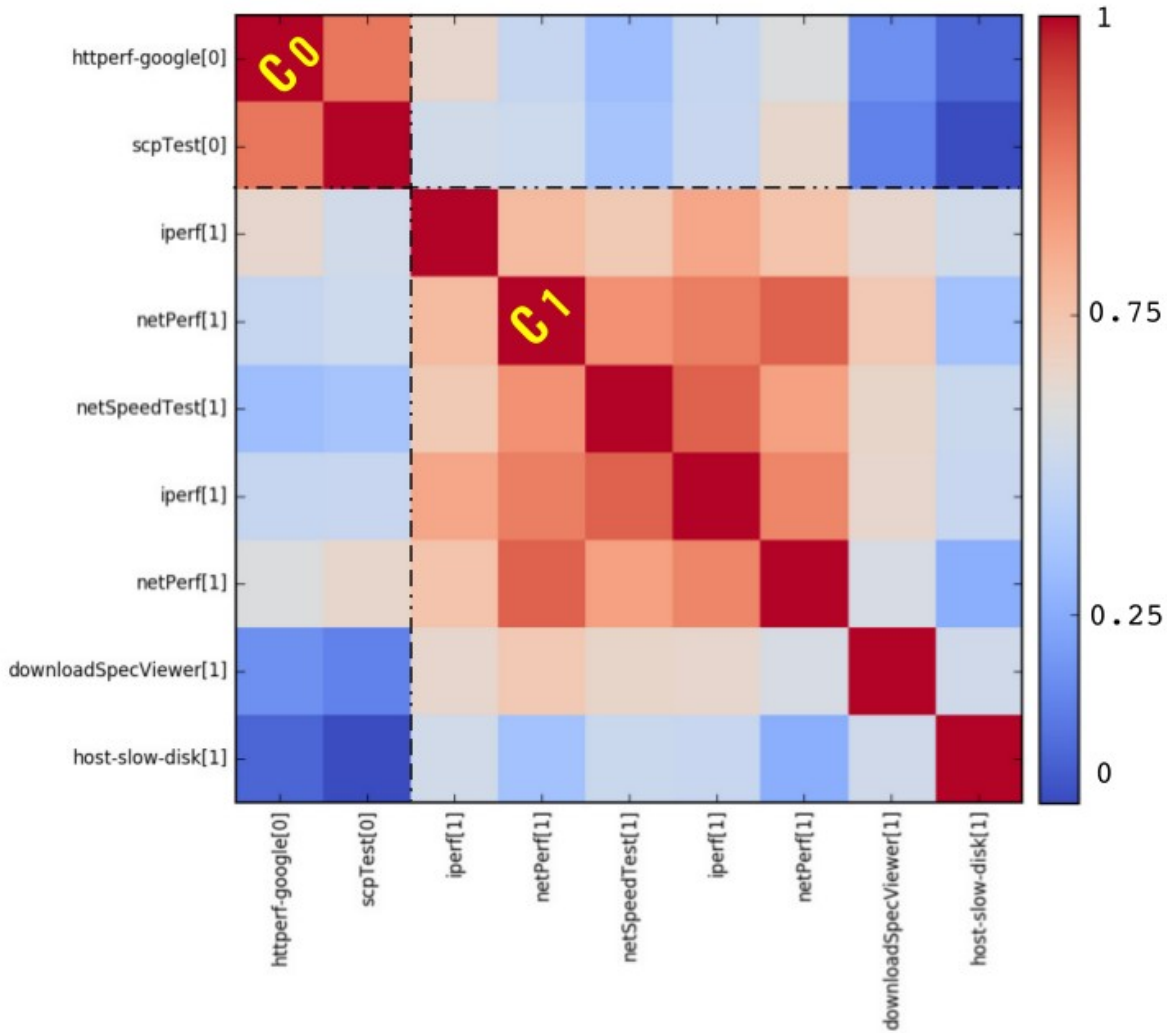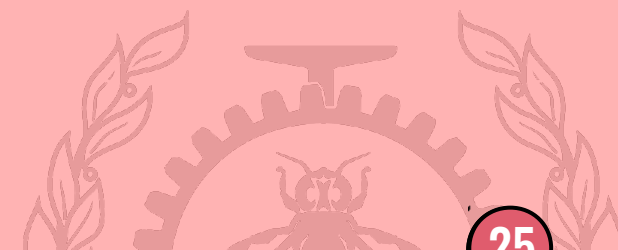**Disk Intensive**

# Experimental Evaluation

## Network Intensive

# Experimental Evaluation
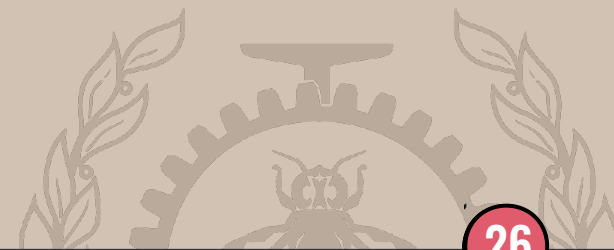
## How to Try

**1**   Access to **Host** only

**2**   Run **LTTng** on Host with my new added tracepoint (vcpu_enter_guest)

**3**   Clone **TraceCompass** from github
(incubator:https://github.com/Nemati/org.eclipse.tracecompass.incubator.git )

**4**   Git checkout vm2

**5**   Clone **lamiminer** from github (https://github.com/azharivs/lamiminer.git)
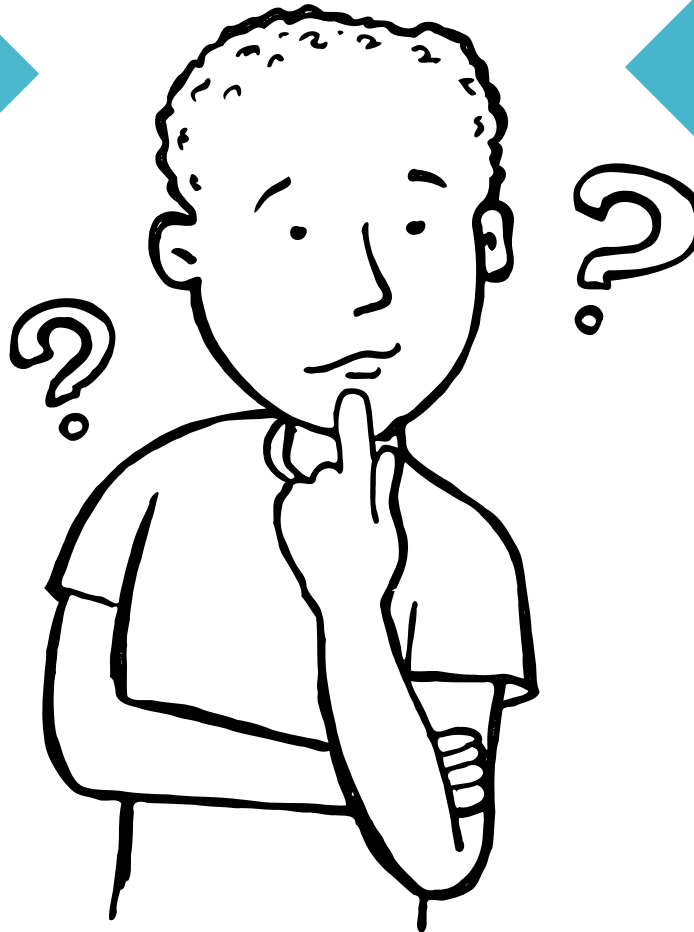
## Conclusions

**1**   Clustering VM based on injected interrupts to VM

**2**   CPU intensive group has high timer interrupt  and high CPU usage

**3**   Disk intensive group has high disk interrupt and high task interrupt

**4**   Network intensive group has high network interrupt and high CPU usage

POLYTECHNIQUE MONTREAL – Hani Nemati & Vahid Azhari

# Questions

azharivs@gmail.com
http://github.com/azharivs

Hani.nemati@polymtl.ca
http://github.com/Nemati