# Honey, I Shrunk the Guests

## Page Access Tracking using a Minimal Virtualisation Layer

30.03.2025

Dustin Tien Nguyen[1], Sebastian Rußer[1], Maximilian Ott[1],
Rüdiger Kapitza[1], Wolfgang Schröder-Preikschat[1], Jörg Nolte[2]

[1] Friedrich-Alexander-Universität Erlangen-Nürnberg
[2] Brandenburgische Technische Universität Cottbus-Senftenberg

## Plenty of memory, hardly any knowledge

### Many memory technologies

- HBM
- DRAM
- Persistent Memory
- CXL-attached memory
- NUMA

### Different characteristics

- Capacity
- Latency
- Bandwidth
- Persistence
- Cost

## Plenty of memory, hardly any knowledge

### Many memory technologies

- HBM
- DRAM
- Persistent Memory
- CXL-attached memory
- NUMA

### Different characteristics

- Capacity
- Latency
- Bandwidth
- Persistence
- Cost

### How can we efficiently utilise memory?

- ⚡ "Expert interfaces" are difficult (PMDK, libnuma)
- ⚡ Language support does not expand to legacy programs (NV-Heaps)

# Plenty of memory, hardly any knowledge

## Many memory technologies

- HBM
- DRAM
- Persistent Memory
- CXL-attached memory
- NUMA

## Different characteristics
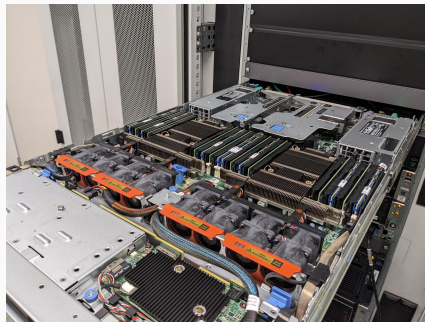
- Capacity
- Latency
- Bandwidth
- Persistence
- Cost

## How can we efficiently utilise memory?

⚡ "Expert interfaces" are difficult (PMDK, libnuma)

⚡ Language support does not expand to legacy programs (NV-Heaps)

> → Operating system support for transparent memory placement

- Fast-tier memory is rare
  $\rightarrow$ it must be distributed efficiently
- Process workloads shift over time
  - in intensity
  - in locality
- Memory placement decisions should be adaptable



$\rightarrow$ Detailed runtime information on memory utilisation is necessary

**Conventional approaches**

- Instrumentalisation
  $\rightarrow$ expensive, very accurate

**Conventional approaches**

- Instrumentalisation
  → expensive, very accurate
- Page table scanning/manipulation
  → expensive, coarse granularity

# Identifying which memory works best for a given process?

**Conventional approaches**

- Instrumentalisation
  $\rightarrow$ expensive, very accurate
- Page table scanning/manipulation
  $\rightarrow$ expensive, coarse granularity
- Sampling (PEBS)
  $\rightarrow$ low overhead, acceptable accuracy

# Identifying which memory works best for a given process?

## Conventional approaches

- Instrumentalisation
  $\rightarrow$ expensive, very accurate
- Page table scanning/manipulation
  $\rightarrow$ expensive, coarse granularity
- Sampling (PEBS)
  $\rightarrow$ low overhead, acceptable accuracy

## Page Modification Logging (PML)

- Virtualisation extension for VM checkpointing, VM migration
- Hardware-based logging of write accesses
- Only works within virtualisation

# Identifying which memory works best for a given process?

## Conventional approaches

- Instrumentalisation
  $\rightarrow$ expensive, very accurate
- Page table scanning/manipulation
  $\rightarrow$ expensive, coarse granularity
- Sampling (PEBS)
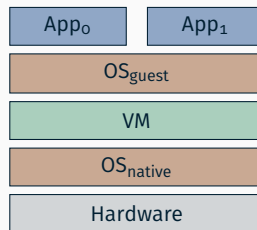  $\rightarrow$ low overhead, acceptable accuracy

## Page Modification Logging (PML)

- Virtualisation extension for VM checkpointing, VM migration
- Hardware-based logging of write accesses
- Only works within virtualisation

> $\rightarrow$ Is virtualisation viable for gathering memory access statistics?
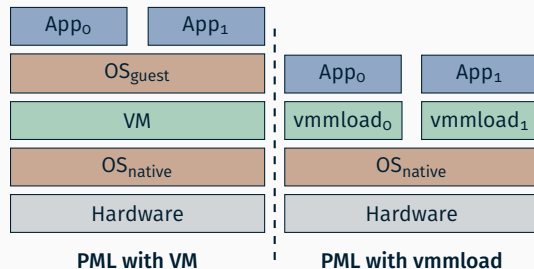
## vmmload: A minimal virtualisation layer

- Guest OS introduces overhead
- Processes should communicate with the host OS
- Processes should communicate with other processes

| $App_0$ | $App_1$ |
|---------|---------|
| $OS_{guest}$ | |
| VM | |
| $OS_{native}$ | |
| Hardware | |

**PML with VM**

- Guest OS introduces overhead
- Processes should communicate with the host OS
- Processes should communicate with other processes
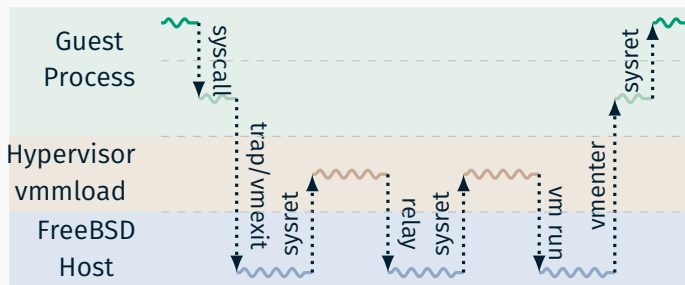- → **vmmload** as minimal hypervisor

| $App_0$ | $App_1$ |
|---------|---------|
| $OS_{guest}$ | |
| VM | |
| $OS_{native}$ | |
| Hardware | |

**PML with VM**

| $App_0$ | $App_1$ |
|---------|---------|
| $vmmload_0$ | $vmmload_1$ |
| $OS_{native}$ | |
| Hardware | |

**PML with vmmload**

# Relaying of System Calls
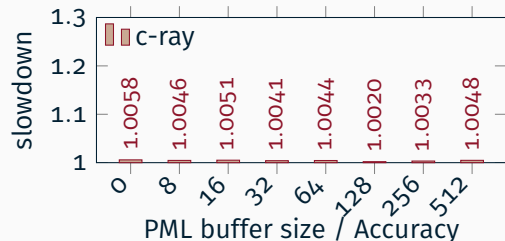
## Requirements
- Isolate hypervisor from guest
- Interface with host OS

## Implications
- Emulate system calls that manipulate the issuers process' state
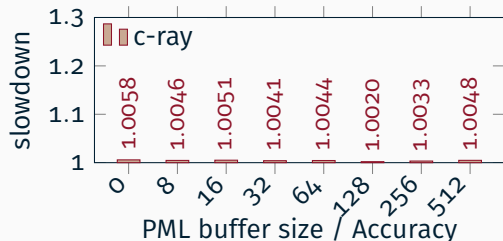- Translation of memory addresses

# Slowdown

## CPU-bound processes



- CPU-bound processes are hardly affected
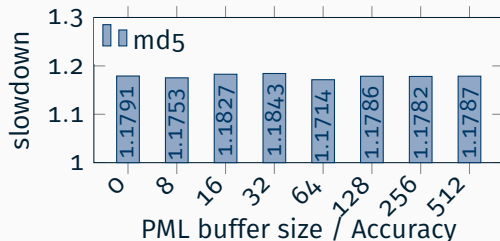- Delay incurred by greater PML buffer size is low
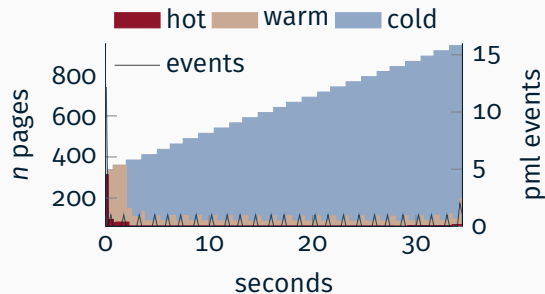
# Slowdown

**CPU-bound processes**



slowdown vs PML buffer size / Accuracy

c-ray

| PML buffer size | value |
|---|---|
| 0 | 1.0058 |
| 8 | 1.0046 |
| 16 | 1.0051 |
| 32 | 1.0041 |
| 64 | 1.0044 |
| 128 | 1.0020 |
| 256 | 1.0033 |
| 512 | 1.0048 |

**I/O bound processes**



slowdown vs PML buffer size / Accuracy

md5

| PML buffer size | value |
|---|---|
| 0 | 1.1791 |
| 8 | 1.1753 |
| 16 | 1.1827 |
| 32 | 1.1843 |
| 64 | 1.1714 |
| 128 | 1.1786 |
| 256 | 1.1782 |
| 512 | 1.1787 |

- CPU-bound processes are hardly affected
- Delay incurred by greater PML buffer size is low

- Highly interactive processes
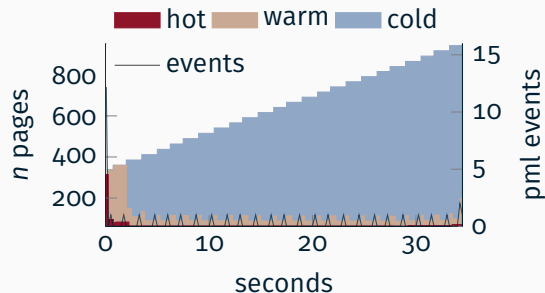- System call overhead introduces great slowdown

# Access Patterns for *c-ray*
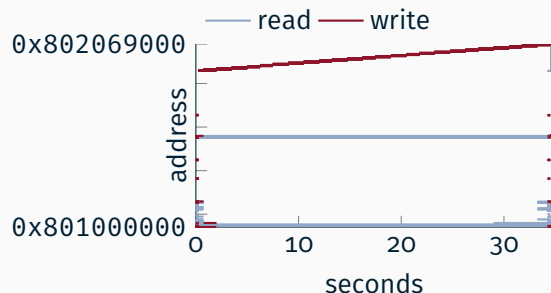


## Temperature

- Identify total number of pages
- Identify frequently accessed pages
- Frequency of PML events

# Access Patterns for *c-ray*



**Temperature**

- Identify total number of pages
- Identify frequently accessed pages
- Frequency of PML events

**Distribution**

- Distinguish between *read*/*write*
- Alteration of working set
- Sparsity/density of accesses

# Summary

### We have

- ✔ shown that **vmmload** can collect memory access statistics
- ✔ achieved statistics over read/write accesses

### Next, we plan to

- 🕐 Reduce system call overhead to $\frac{1}{4}$ with kernel integration
- 🕐 Derive memory placement/migration decisions

> → Source code is freely available

## Summary

**We have**

- ✔ shown that **vmmload** can collect memory access statistics
- ✔ achieved statistics over read/write accesses

**Next, we plan to**

- 🕐 Reduce system call overhead to $\frac{1}{4}$ with kernel integration
- 🕐 Derive memory placement/migration decisions

➜ Source code is freely available

➜ Thank you for your attention