# Memory virtualization overhead mitigation using contiguous memory virtual machines

Peterson Yuhala

Supervised by : Prof. Alain Tchana

University of Neuchâtel - June 2018

# Introduction

- Address translation represents a significant part of virtualization overhead [ISCA'16,ISCA'17].
- Our main goal is to nullify that overhead in virtualized systems.

Institut de Recherche
en Informatique
de Toulouse

32

# Virtualized System Architecture

▶ Memory virtualization depends on VM type. Three main mechanisms exist.

# Memory virtualization

- ▶ Memory virtualization depends on VM type. Three main mechanisms exist.

- ▶ In paravirtualised systems, VMs use a technique called direct paging [Xen'15].

- ▶ Here the guest uses a Physical to Machine mapping maintained within the hypervisor to map virtual addresses directly to physical addresses in its page table [Xen'15].

# Memory virtualization

- Memory virtualization depends on VM type. Three main mechanisms exist.
- In paravirtualised systems, VMs use a technique called direct paging [Xen'15].
- Here the guest uses a Physical to Machine mapping maintained within the hypervisor to map virtual addresses directly to physical addresses in its page table [Xen'15].
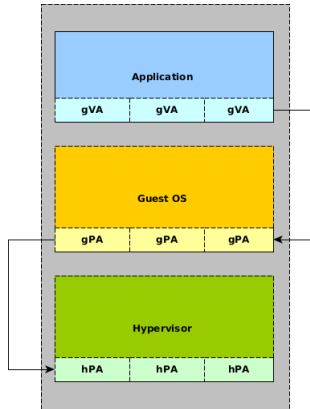- For hardware assisted virtualised systems, we have two main techniques used: shadow paging and nested paging.

# Memory virtualization

▶ Memory virtualization depends on VM type. Three main mechanisms exist.

▶ In paravirtualised systems, VMs use a technique called direct paging [Xen'15].

▶ Here the guest uses a Physical to Machine mapping maintained within the hypervisor to map virtual addresses directly to physical addresses in its page table [Xen'15].

▶ For hardware assisted virtualised systems, we have two main techniques used: shadow paging and nested paging.

▶ Our main focus here is hardware assisted virtualised systems as HVMs are becoming the norm [Amazon AWS].

▶ The main issue here is Memory Address Translation in vitualized systems. ie Obtaining the corresponding machine address in host address space given a virtual address in guest address space.

# Paging
## Basic Definitions

▶ Main memory addresses are called host physical addresses (hPA) and process addresses are called guest virtual addresses (gVA). Addresses in the guest address space are referred to as guest physical addresses (gPA).
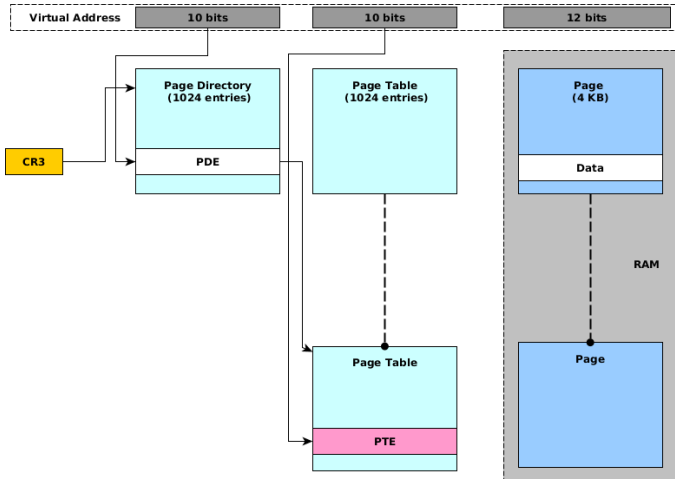
- Native systems (No Virtualization) require one-level address translations. ie virtual address → physical address.

► The following figure illustrates the page walking mechanism in a two-level page table.

▶ Virtualized systems on the other hand require two-level
address translations. ie gVA → gPA AND gPA → hPA.

gVA → [ Guest Page Table (per process) ] → gPA → [ Host Page Table (per VM) ] → hPA

- ► Two approaches are used : **Shadow Paging** and **Nested Paging**.

# Shadow Paging

- In shadow paging the hypervisor maintains a shadow page table for gVA $\rightarrow$ hPA mappings and the **CR3** register points to the shadow page table.
- The hypervisor captures any attempts by the VM to update its page tables and uses gVA $\rightarrow$ gPA and gPA $\rightarrow$ hPA to construct the shadow page table.

# Shadow Paging

- In shadow paging the hypervisor maintains a shadow page table for gVA → hPA mappings and the **CR3** register points to the shadow page table.

- The hypervisor captures any attempts by the VM to update its page tables and uses gVA → gPA and gPA → hPA to construct the shadow page table.

- On a TLB miss, the hardware performs a 1D page walk on the shadow page table to obtain the machine address corresponding to a guest virtual address.

11/32

► Only 4 memory references (same as base native) required for a complete page walk on a TLB miss [ISCA'16].

- Only 4 memory references (same as base native) required for a complete page walk on a TLB miss [ISCA'16].
- No need for any extra hardware support for page walks.

► Every page table update requires a costly trap in the
hypervisor to update shadow page table entries [ISCA'16].

Institut de Recherche
en Informatique
de Toulouse

12/32

► Every page table update requires a costly trap in the hypervisor to update shadow page table entries [ISCA'16].

► Costly traps in the hypervisor on context switches [ISCA'16].

13/32

► Nested paging is a widely used hardware technique to virtualize memory and here the MMU a two dimensional page table consisting of the **guest page table** (gPT) and **host page table** (hPT).

13/32

▶ Nested paging is a widely used hardware technique to virtualize memory and here the MMU a two dimensional page table consisting of the **guest page table** (gPT) and **host page table** (hPT).

▶ On a TLB miss the processor performs a 2D page walk for any gVA → hPA translation.

# Nested Paging
## 2D page walk

Figure

15/32

► Allows fast direct updates to both page tables (gPT and hPT) without any traps in the hypervisor (ie no VMM intervention).

16/32

► On a TLB miss the MMU performs a 2D page walk that
  requires up to 24 memory references. Figure 1 shows how
  up to 24 memory references can be obtained for a
  complete 2D page walk.

| | Base Native | Nested Paging | Shadow Paging |
|---|---|---|---|
| **TLB Hit** | fast | fast | fast |
| **Max. memory access on TLB miss** | 4 | 24 | 4 |
| **Page table updates** | fast direct | fast direct | slow mediated by VMM |
| **Hardware support** | 1D page walk | 2D page walk | 1D page walk |

# Our Solution

► From the table, Nested and Shadow paging both present significant overhead due to memory accesses and page table updates respectively.
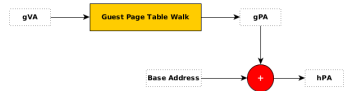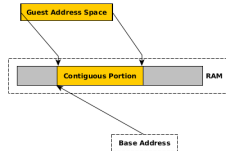
18/32

- ▶ From the table, Nested and Shadow paging both present significant overhead due to memory accesses and page table updates respectively.

- ▶ We propose a new mechanism of address translation which maintains the advantages of both Nested and Shadow paging, and at the same time nullifies the overhead of both mechanisms.

19/32

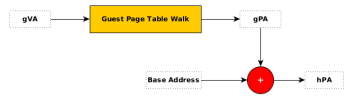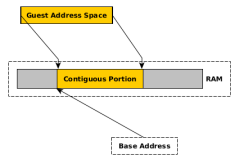▶ Our basic idea is to give contiguous machine memory to VMs and introduce a hardware feature to store the base address of the **contiguous chunk**.

19/32

▶ Our basic idea is to give contiguous machine memory to VMs and introduce a hardware feature to store the base address of the **contiguous chunk**.



▶ Using our solution, the **gPT** and the **base address** are sufficient for a complete gVA → hPA translation.

|  | Base Native | Nested Paging | Shadow Paging | Our Solution |
|---|---|---|---|---|
| TLB Hit | fast | fast | fast | fast |
| Max. memory access on TLB miss | 4 | 24 | 4 | 4 |
| Page table updates | fast direct | fast direct | slow mediated by VMM | fast direct |
| Hardware support | 1D page walk | 2D page walk | 1D page walk | 1D page walk |

- ▶ We made a **simulator** to show that it is possible to give contiguous memory to a high percentage of VMs in a datacenter.

- ▶ We simulated VM traces from Microsoft Azure [SOSP'17] and Bitbrains.

- ▶ The VM traces contain information about VM creation and destruction timestamps, CPU utilization and VM memory.

21/32

22/32

- ▶ The simulator simulates 2 VM placement algorithms :
    - ▶ **Traditional placement algorithm** which takes into account resource availability (spread and stack).
    - ▶ **Contiguity-Aware placement algorithm** which prioritizes servers which can provide contiguous memory for VMs.

22/32

- ▶ The simulator simulates 2 VM placement algorithms :
  - ▶ **Traditional placement algorithm** which takes into account resource availability (spread and stack).
  - ▶ **Contiguity-Aware placement algorithm** which prioritizes servers which can provide contiguous memory for VMs.
- ▶ It takes as input VM traces and then calculates the percentage of VMs with contiguous memory at the end of the simulation.

23/32

- ▶ We use different generations of MS Azure servers in our simulator. ie. generation 3 to 6 and HPC servers so as to reflect the reality in real datacenters.
- ▶ A mix of servers with the following composition was used in our simulator :

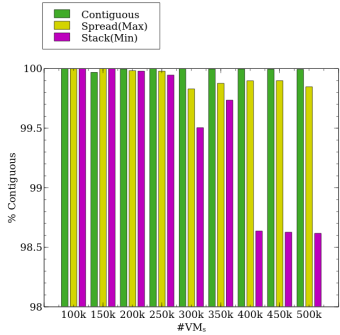| Gen. | % |
|------|-----|
| Gen3 | 20 |
| HPC | 20 |
| Gen4 | 20 |
| Gen5 | 20 |
| Gen6 | 10 |
| Godzilla | 10 |

Figure : % contiguous-memory VMs using different placement
algorithms

Figure : Server consolidation rate for both contiguous-aware
and traditional placement algorithms.

26/32

▶ We have a high % of VMs with contiguous memory
  because VMs sizes are relatively small compared to
  server sizes [SOSP'17].

▶ The relatively small and discrete size of VMs reduces the
  impact of memory fragmentation and prevents resource
  sprawl in the datacenter.

► We implement our solution in Xen hypervisor. Here we modify the default Xen memory allocator so it allocates contiguous chunks of machine memory to VMs.

27/32

- ► We implement our solution in Xen hypervisor. Here we modify the default Xen memory allocator so it allocates contiguous chunks of machine memory to VMs.

- ► To provide contiguous memory to a VM in Xen, we calculate the **order** corresponding to the VM's memory demand. ie **VM Total Memory =** $2^{order}$ **\* 4 KB** .

- ► This order is then passed to the **Xen heap allocator** which provides a contiguous chunk of memory containing $2^{order}$ contiguous 4kb pages.

28/32

▶ We succeed already in giving contiguous memory to PV
  VMs. We would then eliminate unneccessary hypercalls
  knowing that VMs have contiguous memory.

28/32

- ▶ We succeed already in giving contiguous memory to PV VMs. We would then eliminate unneccesary hypercalls knowing that VMs have contiguous memory.
- ▶ We succeed in mapping 99% of HVM memory to contiguous machine memory. We would modify the Xen heap allocator to give 100% contiguity. The memory translation algorithm would then be modified as per our basic idea.

28/32

- ▶ We succeed already in giving contiguous memory to PV VMs. We would then eliminate unneccessary hypercalls knowing that VMs have contiguous memory.

- ▶ We succeed in mapping 99% of HVM memory to contiguous machine memory. We would modify the Xen heap allocator to give 100% contiguity. The memory translation algorithm would then be modified as per our basic idea.

- ▶ The last step would be the implementation of the whole solution in a simulator and processing of final results.

▶ **Memory isolation** : Each VM must be confined within its memory region. How do we ensure isolation knowing that the hypervisor has minimal interference.

29/32

- **Memory isolation** : Each VM must be confined within its memory region. How do we ensure isolation knowing that the hypervisor has minimal interference.
- How to **enforce contiguity**.

29/32

- **Memory isolation** : Each VM must be confined within its memory region. How do we ensure isolation knowing that the hypervisor has minimal interference.
- How to **enforce contiguity**.
- **Cohabitation** of our solution with Nested Paging or Shadow Paging. ie Contiguous memory VMs use our solution and other VMs use Nested Paging/Shadow Paging.

30/32

**Memory Isolation** :

▶ Introduction of a base address (**gBA**) for each VM
(included within the VM Control Structure (VMCS)).

**Memory Isolation** :

► Introduction of a base address (**gBA**) for each VM
(included within the VM Control Structure (VMCS)).

► Once scheduled, the VM's **gBA** and VM's memory size
(**gMS**) are stored within dedicated processor registers.

**Memory Isolation** :

► Introduction of a base address (**gBA**) for each VM
(included within the VM Control Structure (VMCS)).

► Once scheduled, the VM's **gBA** and VM's memory size
(**gMS**) are stored within dedicated processor registers.

► To ensure isolation, the TLB is enhanced with a **security
module** such that during insertion into the TLB, the
security module checks if $hPA \leq gBA + gMS$

**Memory Contiguity** :

▶ At boot time, we allocate contiguous memory to VMs if possible.

Institut de Recherche
en Informatique
de Toulouse

**Memory Contiguity** :

▶ At boot time, we allocate contiguous memory to VMs if possible.

▶ For VMs with discontinuous memory, the hypervisor reorganizes the VM's memory to be contiguous. This can be done either by migrating the VM's pages during idle periods or by migrating the VM to another host.

Institut de Recherche
en Informatique
de Toulouse

▶ Memory virtualization overhead can be very high when
  using traditional techniques (Nested and Shadow Paging).

- ▶ Memory virtualization overhead can be very high when using traditional techniques (Nested and Shadow Paging).
- ▶ We propose to eliminate the overhead of these traditional techniques using contiguous-memory VMs.

32/32

- ▶ Memory virtualization overhead can be very high when using traditional techniques (Nested and Shadow Paging).
- ▶ We propose to eliminate the overhead of these traditional techniques using contiguous-memory VMs.
- ▶ Through simulation we showed that it is possible to give contiguous memory to a high % of VMs in datacenters.

32/32

- ▶ Memory virtualization overhead can be very high when using traditional techniques (Nested and Shadow Paging).
- ▶ We propose to eliminate the overhead of these traditional techniques using contiguous-memory VMs.
- ▶ Through simulation we showed that it is possible to give contiguous memory to a high % of VMs in datacenters.
- ▶ Our solution is being prototyped in the Xen hypervisor; we would then evaluate its performance to show that it is indeed better than Nested and Shadow paging.

Thank you for listening !