



Scalevisor: A distributed hypervisor for rack-scale computing

(日) (四) (문) (문) (문)

Supervisors: Gaël Thomas, Mathieu Bacou

Pipereau Yohan

Télécom SudParis - Institut Polytechnique de Paris

July 7, 2021

How does VM hosting works ?



Figure: Use case





2/28 July 7, 2021

Cloud provider problem: fragmentation



Figure: Resource manager - VM Scheduler

イロン イボン イヨン イヨン 三日

TELECOM SudParis POLYTECHNIQUE



Wasted resources

Eolas reported 40% unused resources in its datacenter !

Improving consolidation

Saving money

Reducing Energy Consumption

Problem - Complexity

Finding minimum number of servers for a set of VM is NP-hard !





Solution: memory mutualization



Figure: Resource manager - VM Scheduler

イロン イボン イヨン イヨン 三日



Solution: memory disaggregation



Figure: Resource manager - VM Scheduler



(日)



Virtual machines consolidation

Definition - Consolidation

- **scenario 1**: run more VMs on same number of PMs
- **scenario 2**: run same number of VMs on less PMs

2 configurations for VM placement

- Static VM Consolidation: initial placement of VMs
- Dynamic VM Consolidation, reallocation of existing VMs

Some techniques

VM migration: move 1 VM from PM1 to PM2

memory overcommitment: ΣVM_{memory} > PM_{memory} (memory ballooning, transparent page sharing, memory compression, swapping)



Dynamic resource management - VM migration

Similar Problem

Similar to static resource management for VM allocation

Useful

Upgrade a server

Consolidate virtual machines







イロン イロン イヨン イヨン 三日



How can we mutualize virtual machines memory in a datacenter ?





(ロ)(同)(E)(E)(E)(E)



Design & Implementation



July 7, 2021 9/28

Pipereau Yohan

Scalevisor: A distributed hypervisor for rack-scale computing

Guest memory view



(a) Classical VM memory usage

















July 7, 2021

Design & Implementation

rmem layer





Pipereau Yohan

Scalevisor: A distributed hypervisor for rack-scale computing

イロン イボン イヨン イヨン 三日



```
struct remote_memory_region {
    u64 remote_addr;
    u32 remote_rkey;
    u32 remote_len;
};
```

```
alloc : 1-sided
```

```
free : 1-sided
```

```
write : 2-sided
```

```
read : 2-sided
```

イロン イロン イヨン イヨン 三日



Design & Implementation

paravirtualization layer





15/28 July 7, 2021

Pipereau Yohan

Scalevisor: A distributed hypervisor for rack-scale computing

イロン イロン イヨン イヨン 三日

July 7, 2021

What is paravirtualization

Definition - Paravirtualization

A technique of communication between a VM and a hypervisor



Scalevisor: A distributed hypervisor for rack-scale computing

イロン イボン イヨン イヨン 三日





July 7, 2021

Why do we need paravirtualization ?

Semantic gap

Hypervisor has no visibility on memory management in the guest



Figure: false page anonymity



イロン イロン イヨン イヨン 三日

17/28

July 7, 2021

Design & Implementation

swap layer





Pipereau Yohan

Scalevisor: A distributed hypervisor for rack-scale computing

イロン イロン イヨン イヨン 三日

swap layer





Design & Implementation

How to use ?



18/28 July 7, 2021

Pipereau Yohan

Scalevisor: A distributed hypervisor for rack-scale computing

イロン イロン イヨン イヨン 三日



VM allocation steps

- $1. \ \mbox{In the host, connect with a memory server}$
- 2. In the host, launch VM vith paravirtualized driver
- 3. In the guest, swapon a block device or swapfile



イロン イロン イヨン イヨン 三日



POLYTECHNIQUE

Design & Implementation

Live VM migration



19/28 July 7, 2021

Pipereau Yohan

Scalevisor: A distributed hypervisor for rack-scale computing

イロン イロン イヨン イヨン 三日

Why do we need live VM migration ?

Currently

✓ Memory mutualization for static VM consolidation

Remember - Live VM migration

For upgrade and dynamic VM consolidation



イロン イロン イヨン イヨン 三日



INTON

Design & Implementation Live VM migration

Support for Live VM migration



Figure: RX tree in guest for live migration



21/28 July 7, 2021

Evaluation

Evaluation





POLYTECHNIQUE

Pipereau Yohan

Scalevisor: A distributed hypervisor for rack-scale computing

Number of lines of code

project	SLOC
qemu	713
linux	6
virtio-rmem	993
vhost-rmem	1713
rmem_server	1103

Table: Number of lines of code





< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □



Evaluation



kmeans 1thread: wall clock runtime as a function of local memory ratio

Figure: Wall clock runtime as a function of how much memory is local



Scalevisor: A distributed hypervisor for rack-scale computing

イロン イロン イヨン イヨン 三日







Different degradation profile under different memory disaggregation ratio !





24/28 July 7, 2021

Pipereau Yohan

Scalevisor: A distributed hypervisor for rack-scale computing

Conclusion

Conclusion



TELECOM SudParis

24/28 July 7, 2021

Pipereau Yohan

Scalevisor: A distributed hypervisor for rack-scale computing

・ロト ・四ト ・ヨト ・ヨト ・ヨ



Remember one thing

We swap on remote memory to use remaining free memory on other hosts.

Expectations

- 1. Remote memory latency: 5-10 µs
- 2. Improve consolidation ratio
- 3. Reduce VM migration time

(ロ)(同)(E)(E)(E)(E)



July 7, 2021

26/28



- Support Live VM migration of zerocopy devices
- What happens when allocations on remote node fails for a size because of memory fragmentation





イロン イボン イヨン イヨン 三日



Conclusion

Questions ?



TELECOM SudParis

18110

27/28 July 7, 2021

Pipereau Yohan

Scalevisor: A distributed hypervisor for rack-scale computing

・ロト ・四ト ・ヨト ・ヨト ・ヨ

Appendix

Appendix



Pipereau Yohan

Scalevisor: A distributed hypervisor for rack-scale computing

・ロン ・回入 ・ヨン ・ヨン 三日



POLYTECHNIQUE DE PARIS

Distributed Shared Memory 1/2

How ? Deviate TDP page fault

Advantages

July 7, 2021

- Enable to build a real address space
- More flexibility to implement cache coherency operations





イロン イボン イヨン イヨン 三日





Distributed Shared Memory 2/2

Drawbacks

- Lot of engineering
- Vendor-specific
- what if shadow page table ?
- In-kernel modifications
- No distinction between hot memory and cold memory
- Cache coherence protocols are expensive

(ロ)(同)(E)(E)(E)(E)



Appendix

Influence of merging requests



merging enabled (2020 05 01 1), merging disabled (2020 05 14 2)

2020 05 01 1 read 2020 05 14 2 read 2020 05 01 1 write 2020 05 14 2 write

Figure: Enabling/Disabling requests merging







Wall time as a function of disaggregation ratio for kmeans

Figure: Enabling/Disabling requests merging

(日)

POLYTECHNIQUE

TELECO SudPari

A qemu-kvm virtual machine



Figure: Enabling/Disabling requests merging

イロン イロン イヨン イヨン 三日

POLYTECHNIQUE

paravirtualization layer



Figure: paravirtualization



(日)

