# KVM GUEST FREE PAGE HINTING
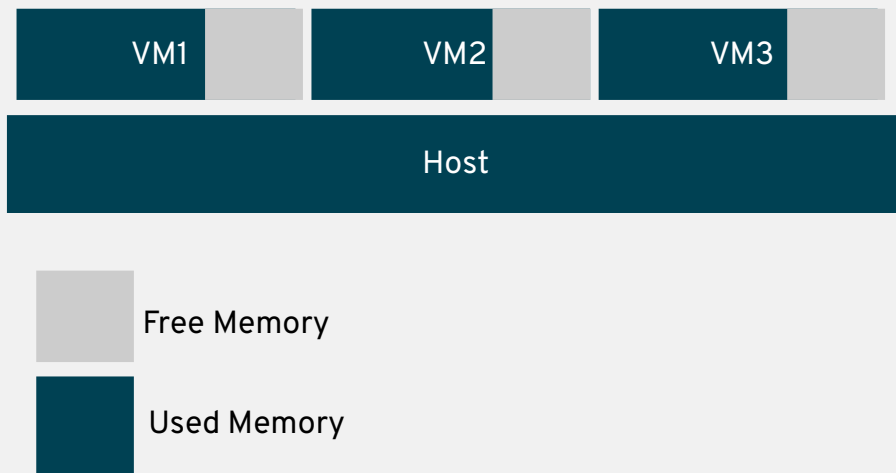
Nitesh Narayan Lal
Software Engineer
October 26th 2018

# AGENDA

- Issue

- Objective

- Implementation

# ISSUE

Inability of the guest to report the free memory back to the host even if the host is running out of memory.
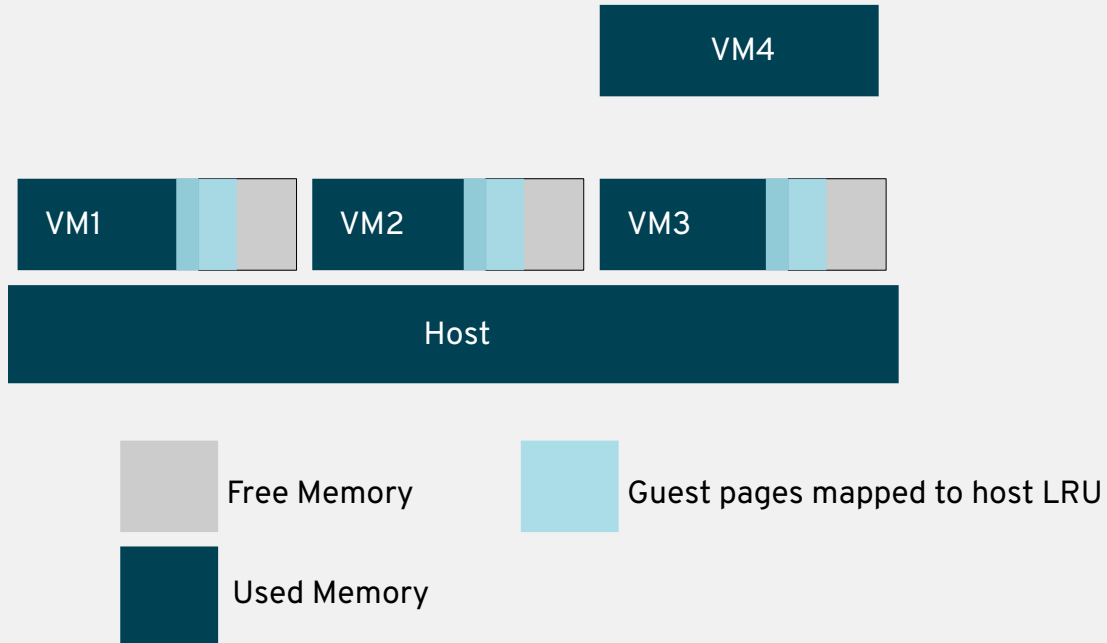
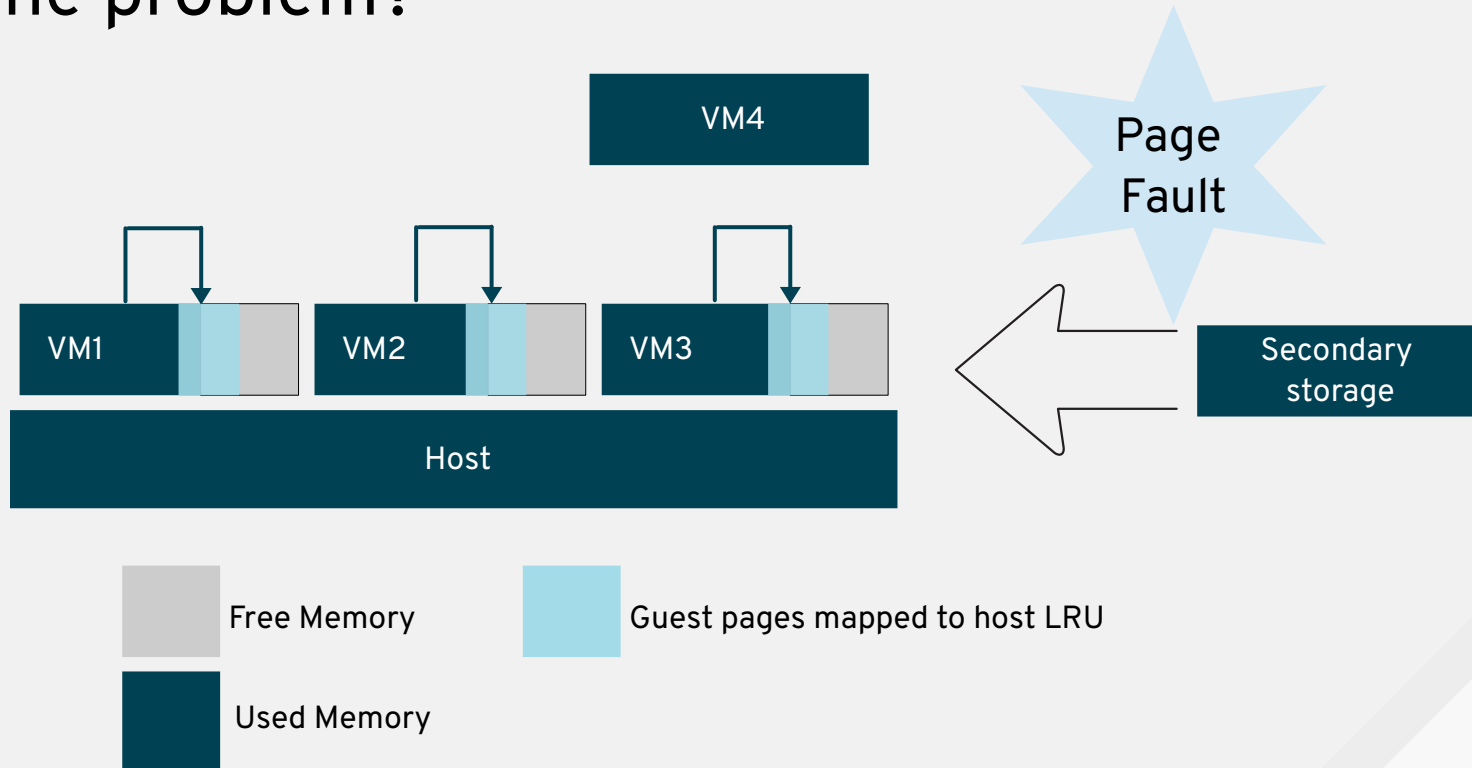# Host wants to launch another guest...

- Is it possible?
  - YES

Photo by pixabay.com from Pexels

# How?

# What's the problem?



VM4

VM1  VM2  VM3

Page Fault

Secondary storage

Host

Free Memory

Guest pages mapped to host LRU

Used Memory

redhat.

# Issues?

vCPU
Stall


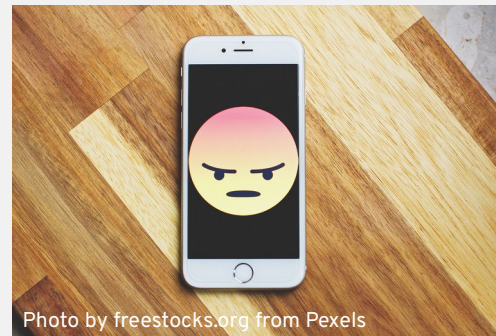Photo by Artem Bali from Pexels

IO is expensive


Photo by freestocks.org from Pexels

IO on free pages???

Why don't you discard its data?
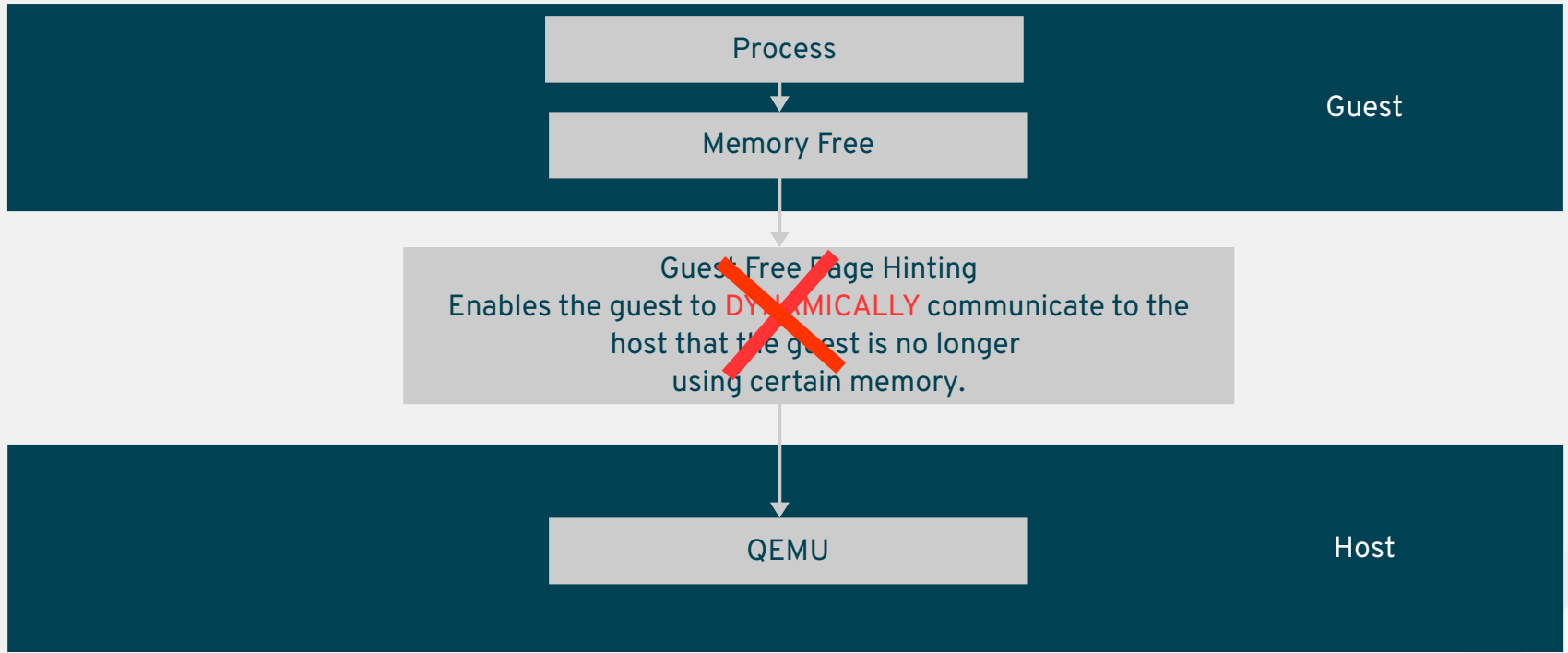
# OBJECTIVE

Enhancement in the kernel to fit more virtual machines/process on each system, by removing unused memory from virtual machines.



Free Memory

Used Memory

# CURRENT SITUATION

Process

Memory Free

Guest

Guest Free Page Hinting
Enables the guest to DYNAMICALLY communicate to the
host that the guest is no longer
using certain memory.

QEMU

Host

redhat.

# INITIAL IMPLEMENTATION

1. Per CPU array is full
2. Acquire seqlock

Process → Page Free → guest_free_page → Check reallocation    Slowpath

Global array

Release lock

number of entries > threshold

Notify host

arch_alloc_page
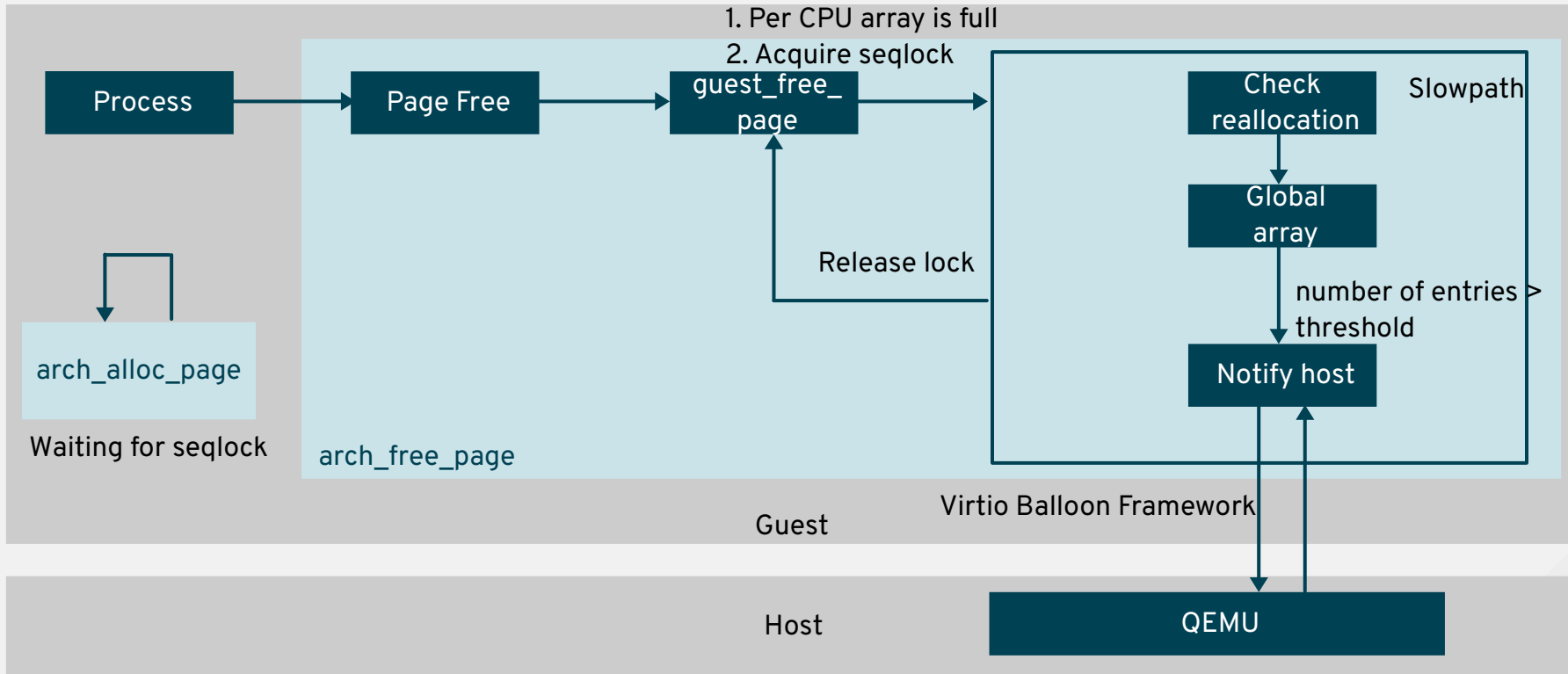
Waiting for seqlock

arch_free_page

Virtio Balloon Framework

Guest

Host

QEMU

redhat.
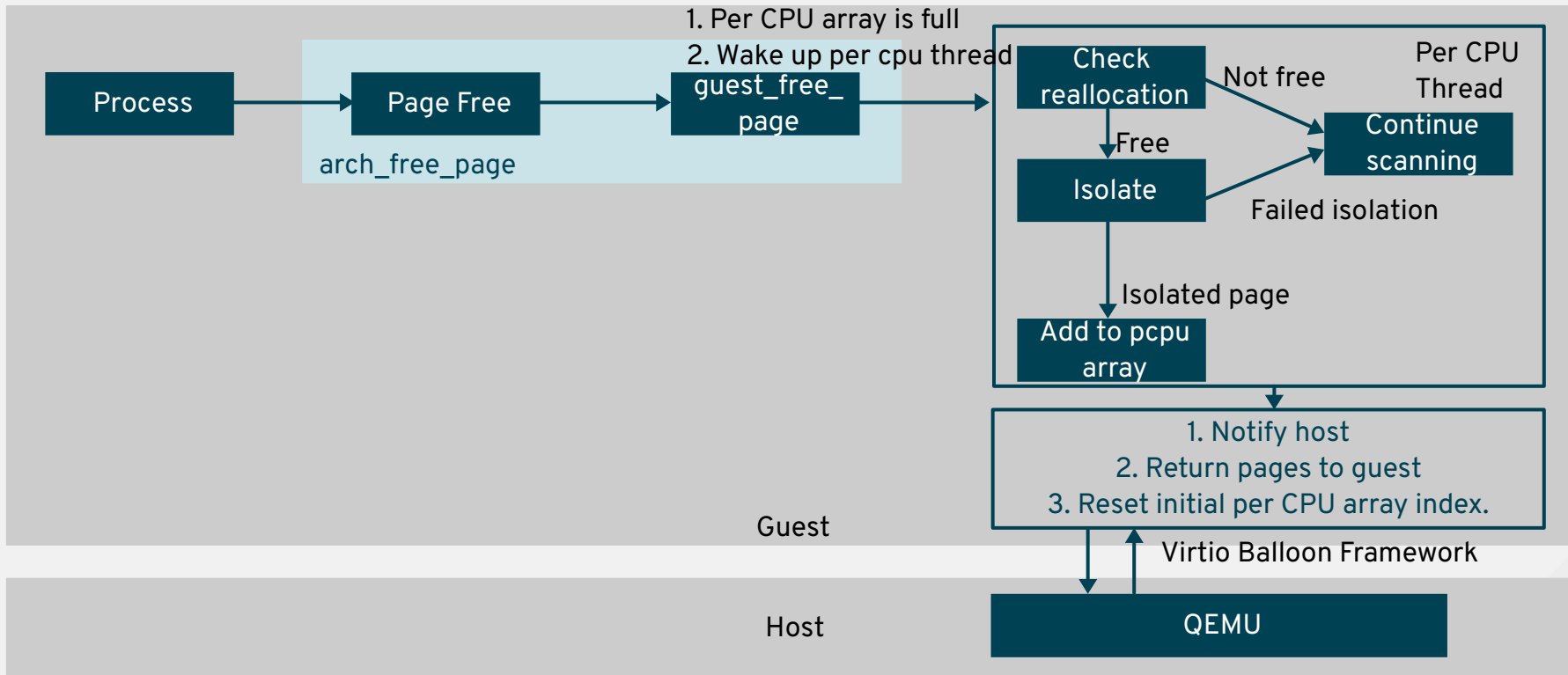
# ISSUES

- Blocks all allocations until the host frees the captured free pages.
    - To avoid reallocation of page which we are sending to the host.
- Blocks arch_free_page() & arch_alloc_page flow.

redhat.

# CURRENT IMPLEMENTATION

Process → Page Free → guest_free_page

arch_free_page

1. Per CPU array is full
2. Wake up per cpu thread

Per CPU Thread

Check reallocation → Not free → Continue scanning

Free

Isolate → Failed isolation → Continue scanning

Isolated page

Add to pcpu array

1. Notify host
2. Return pages to guest
3. Reset initial per CPU array index.

Guest

Virtio Balloon Framework

Host

QEMU

redhat.

# BENEFITS

- Still reports the free memory dynamically without human intervention.

- Zone lock is acquired only for the page to be scanned.

  - Other zones are still free to do allocations.

- Doesn't block the arch_free_page() or arch_alloc_page() flow.

# CHALLENGES

- Failing to capture good number of freed pages due to limited array size.
- Failing to isolate a good number of pages.
    - Probably because they are still in Per-CPU Page Frame Cache list and not in buddy.
- Find a good test-case that's real-world and easy to observe in development.
- Analyze and reduce the per CPU thread overhead.

redhat.

# LIMITATIONS

- We depend on buddy free list for freeing a page.
  - If a page is not present in buddy free list we can not isolate.
- Isolation failures under high memory pressure on guest.

redhat.

# FUTURE OPTIMIZATIONS

- Kicking host only with a larger set of isolated pages.

- Sorting per CPU entries based on zone number to avoid repetitive locks on the same zone.

- Making the threshold condition configurable based on the user-case/requirement.

redhat.

# REFERENCES

- Last posted upstream patch-set: https://www.spinics.net/lists/kvm/msg170113.html
- Development Linux repository: https://github.com/niteshnarayanlal/linux-hinting
- Development QEMU repository: https://github.com/niteshnarayanlal/qemu

redhat.

Backup Slide

# How is this different from "Virtio-balloon: support free page reporting"?

- Similarity?
- Differences?
  - Implementation
  - Use-case
- Will "Guest free page hinting" approach work for migration use case?
  - Can we compare?

redhat.

# INITIAL IMPLEMENTATION

1. Per CPU array is full
2. Acquire seqlock

Process → Page Free → guest_free_page → Check reallocation

1. Global list is full
2. Compress & pack

Global array → Pack array

1. If number of entries < threshold
2. Copy them to per CPU array
3. Release lock

number of entries > threshold

Notify host

1. Clear global array
2. Release lock

arch_alloc_page

Waiting for seqlock

arch_free_page

Slowpath

Virtio Balloon Framework

Guest

Host

QEMU

redhat.