# Speed Up Boot-up Time for Guest in Alibaba Cloud

Li Weinan william.lwn@alibaba-inc.com
Guo Cheng hanyu.gc@alibaba-inc.com
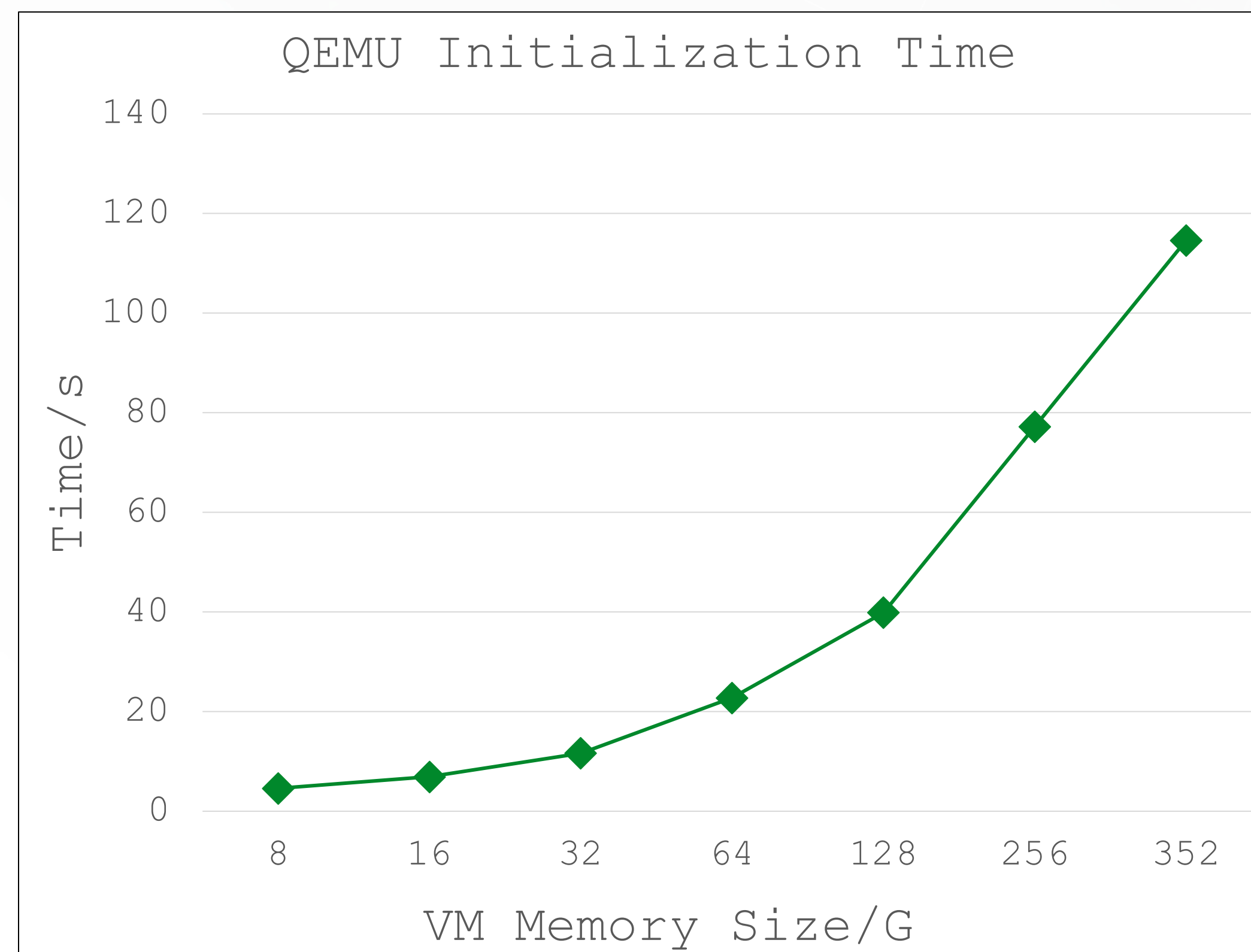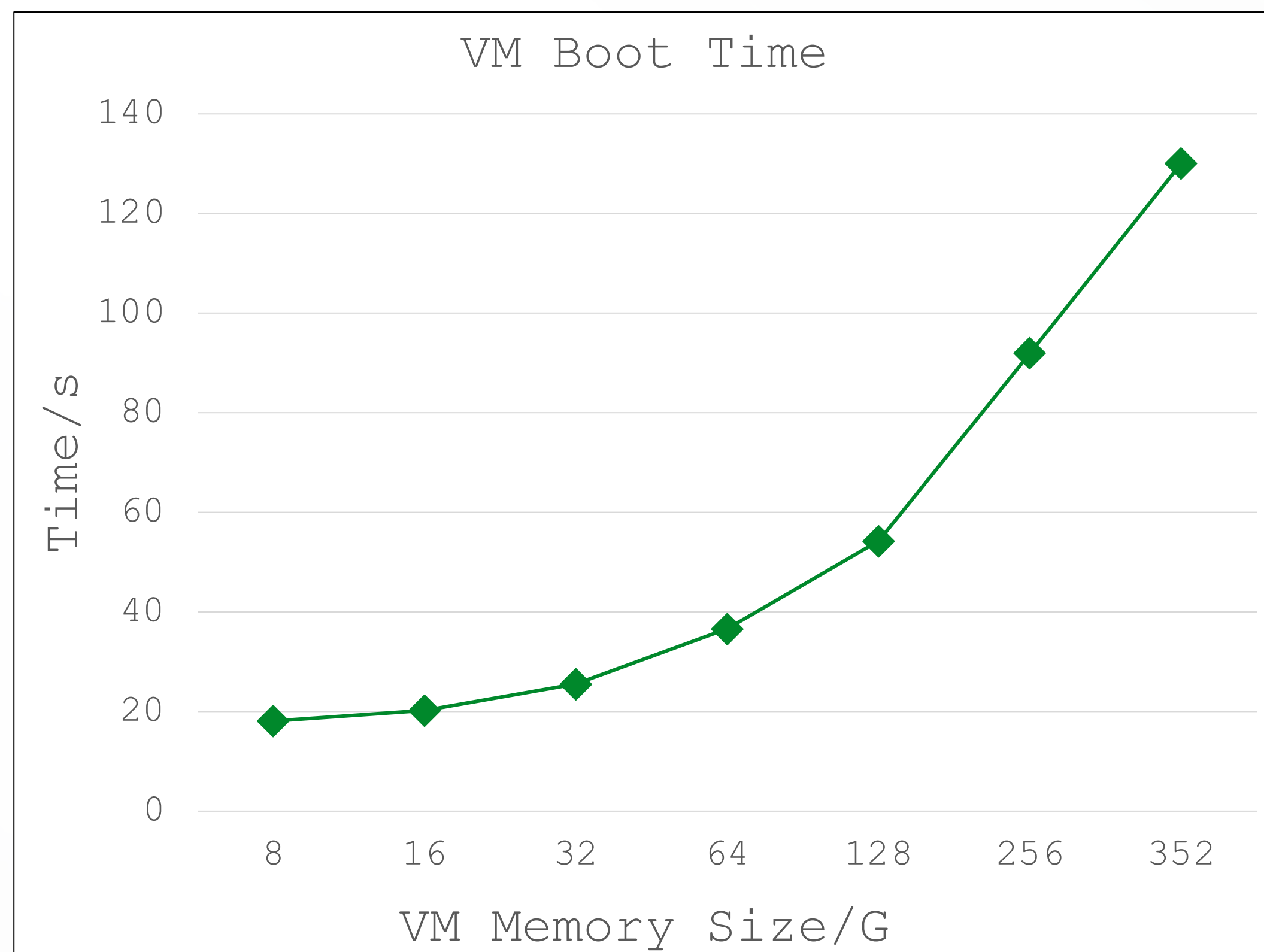
Alibaba Cloud
Worldwide Cloud Services Partner

# Agenda

➢ Background

➢ Async dma map

➢ Guest boot process with async dma map

➢ Optimization design

➢ Achievements

**What is the problem?**

➢Dma_map all the guest memory when there is passthrough device

➢8G->384G
Dma_map time is one big problem!

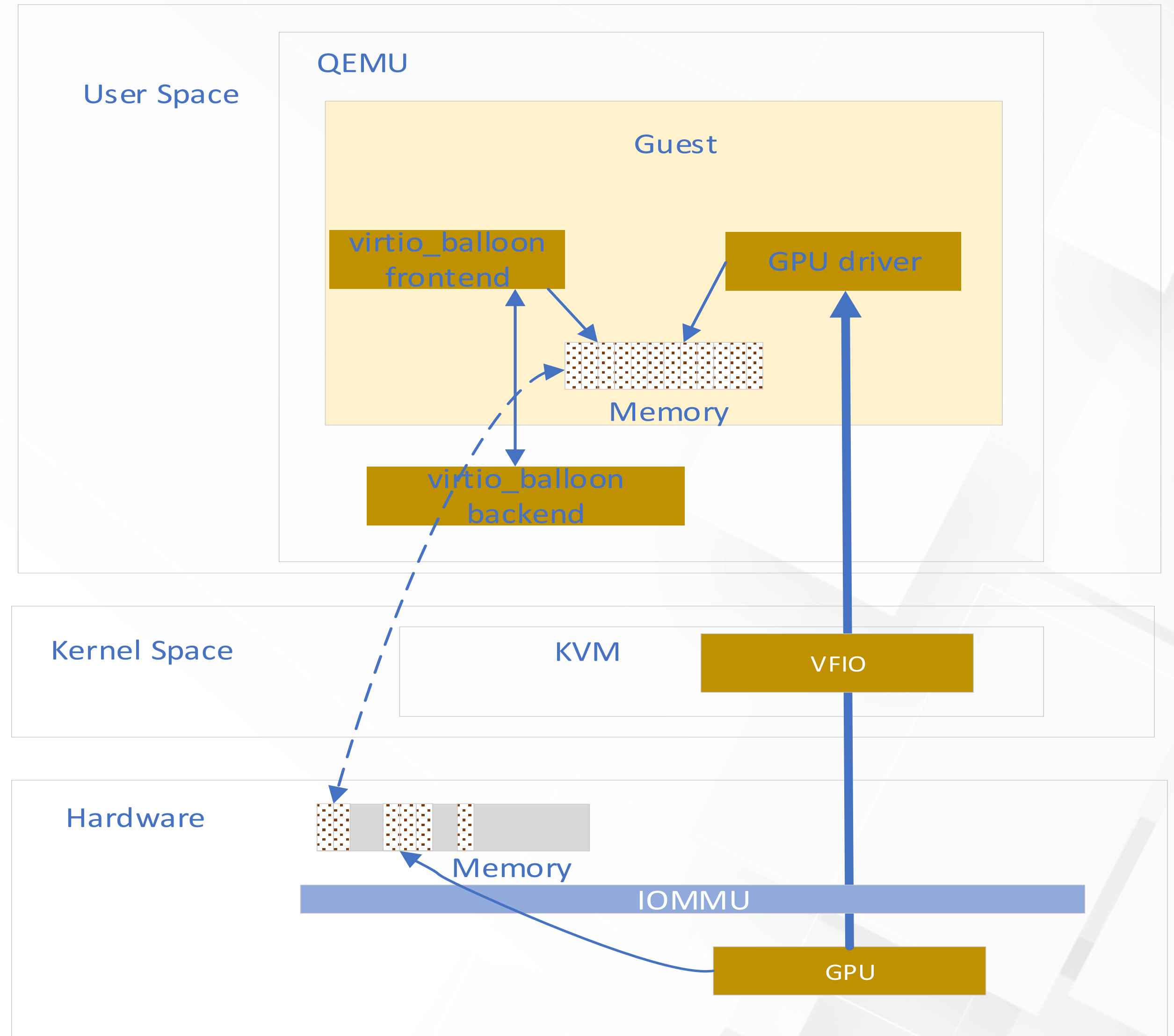# Guest boot & QEMU initialization time vs memory size

# Conditions

➤ More time costs along with more memory

➤ No DMA no dma_map

➤ DMA specific range memory

# Options

➢ virtual IOMMU

➢ Async dma map
- – Only map necessary memory first

- Map asynchronously in the background

# Async dma_map
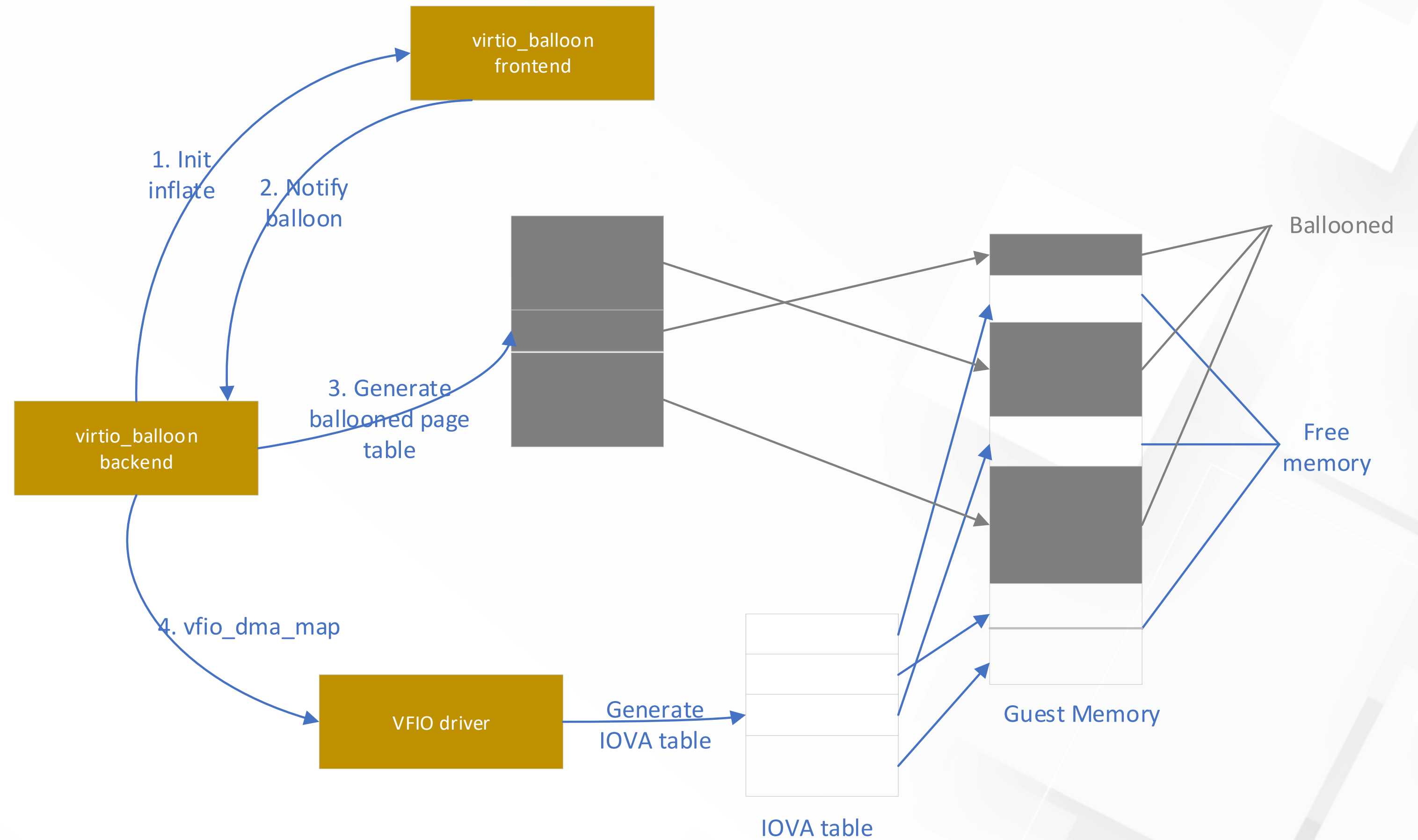
⭐ Balloon memory before allocated for DMA



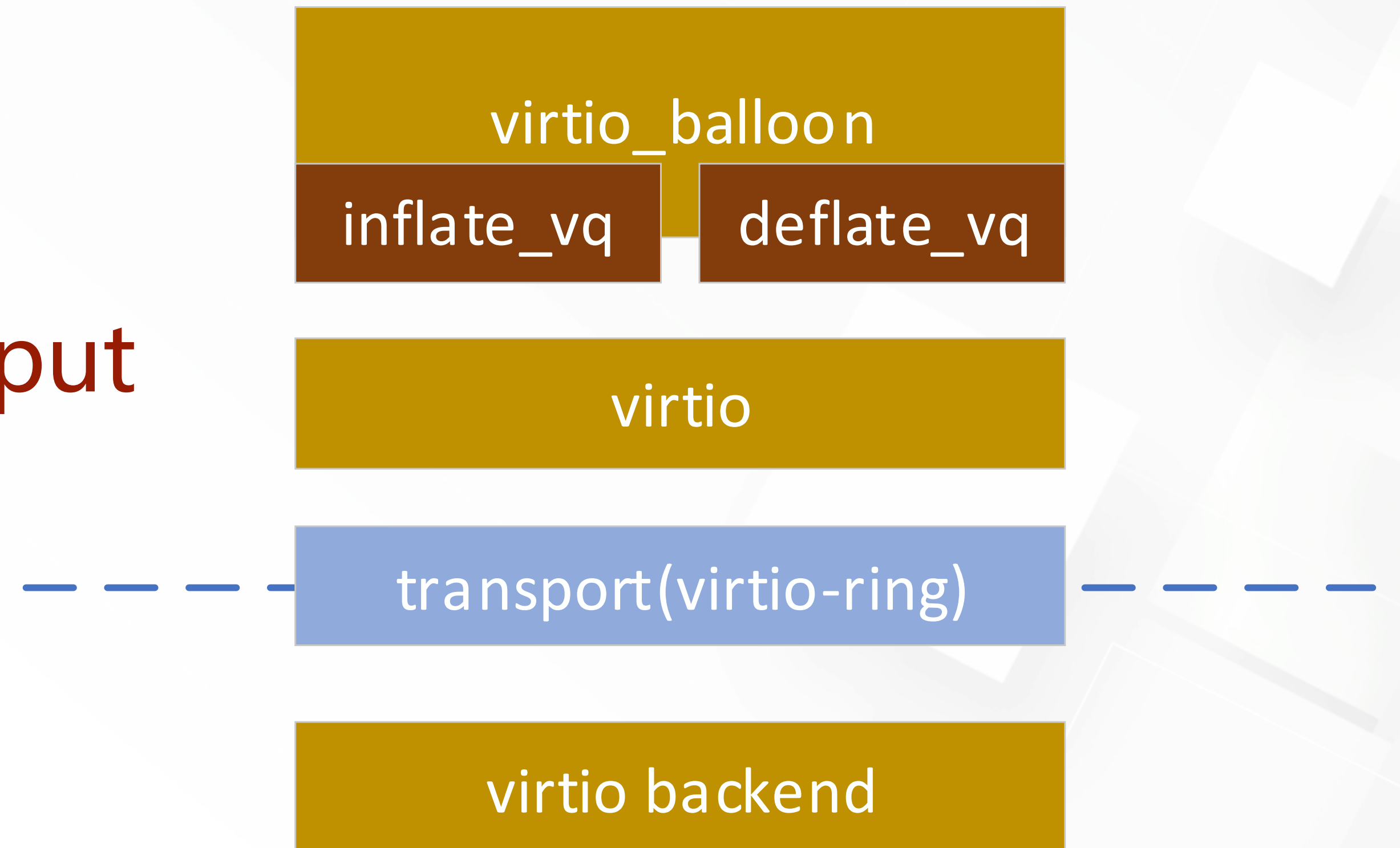Overview of memory access with a passthrough device

# Architecture Overview

➢ QEMU
  Trigger vfio_dma_map
  Trigger balloon change
  Track ballooned pages

➢ virtio_balloon driver
  Balloon pages
  Tell to host (QEMU)

➢ VFIO driver
  do vfio_pin_map_dma

virtio_balloon
frontend

1. Init
inflate

2. Notify
balloon

virtio_balloon
backend

3. Generate
ballooned page
table

4. vfio_dma_map

VFIO driver

Generate
IOVA table

IOVA table
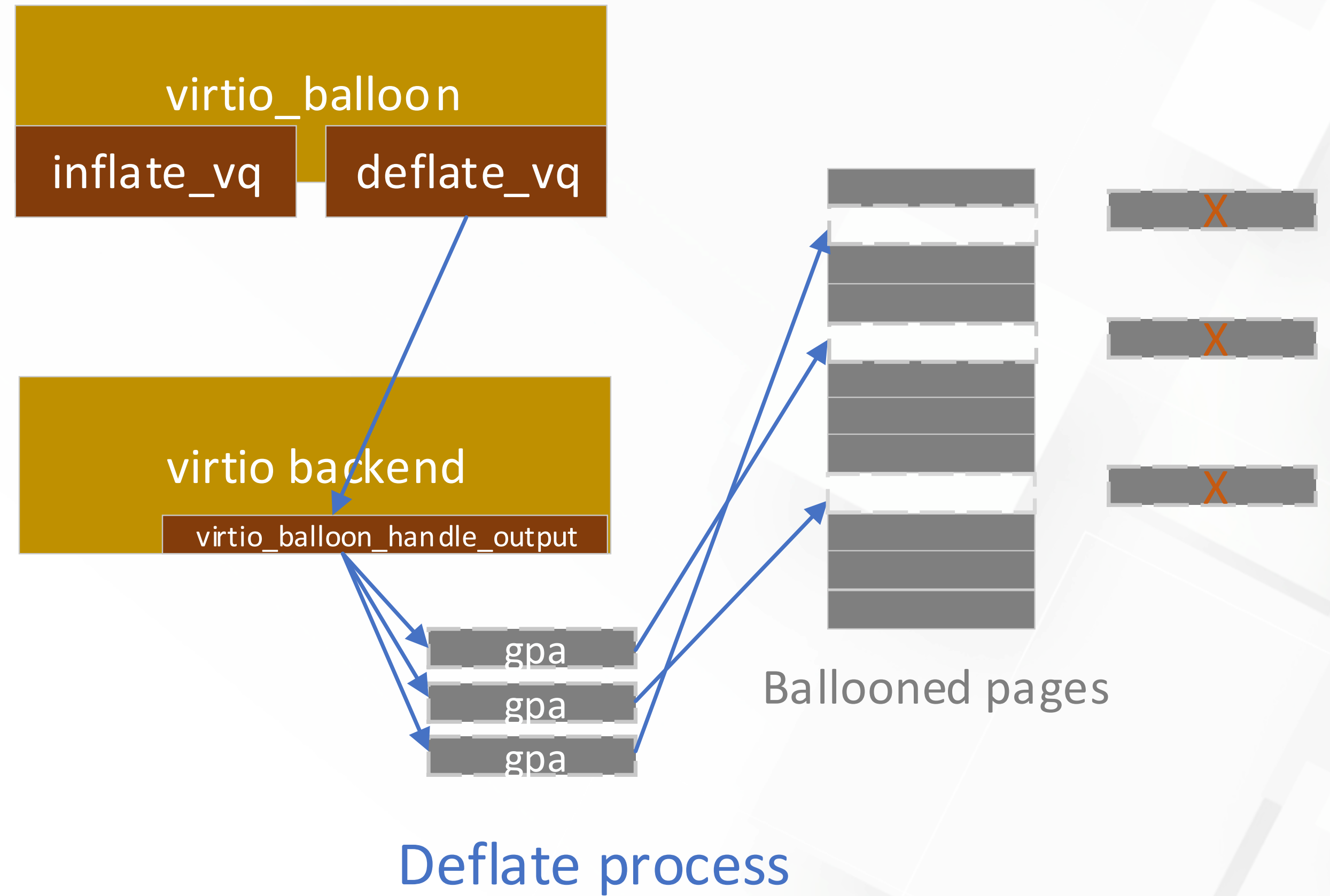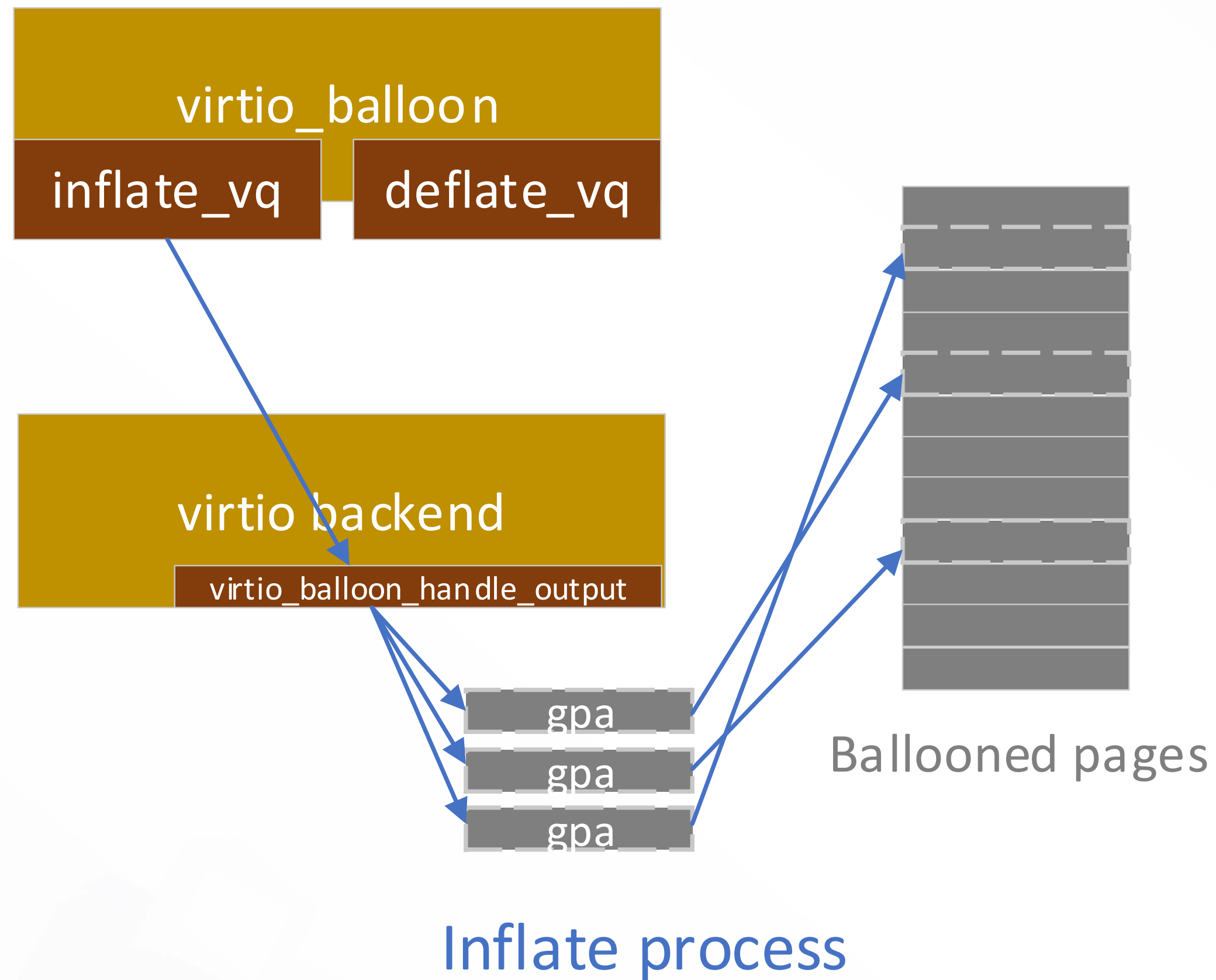
Ballooned

Free
memory

Guest Memory

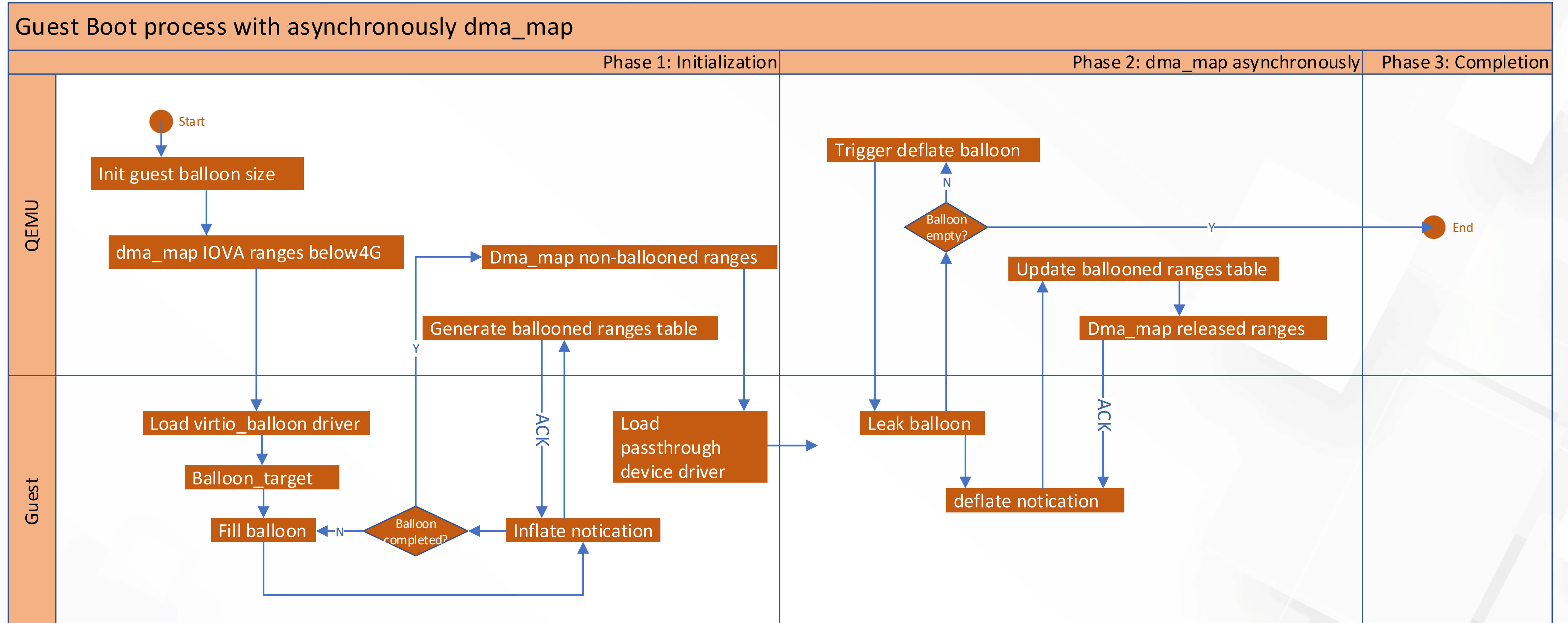# virtio_balloon communication

## Related functions and struct
- inflate_vq
- deflate_vq
- virtio_balloon_handle_output
- VirtQueueElement:

guest PFN

page_num

# Balloon range tracking workflow



Inflate process

Deflate process

# Guest boot process with async dma map

# Optimization design
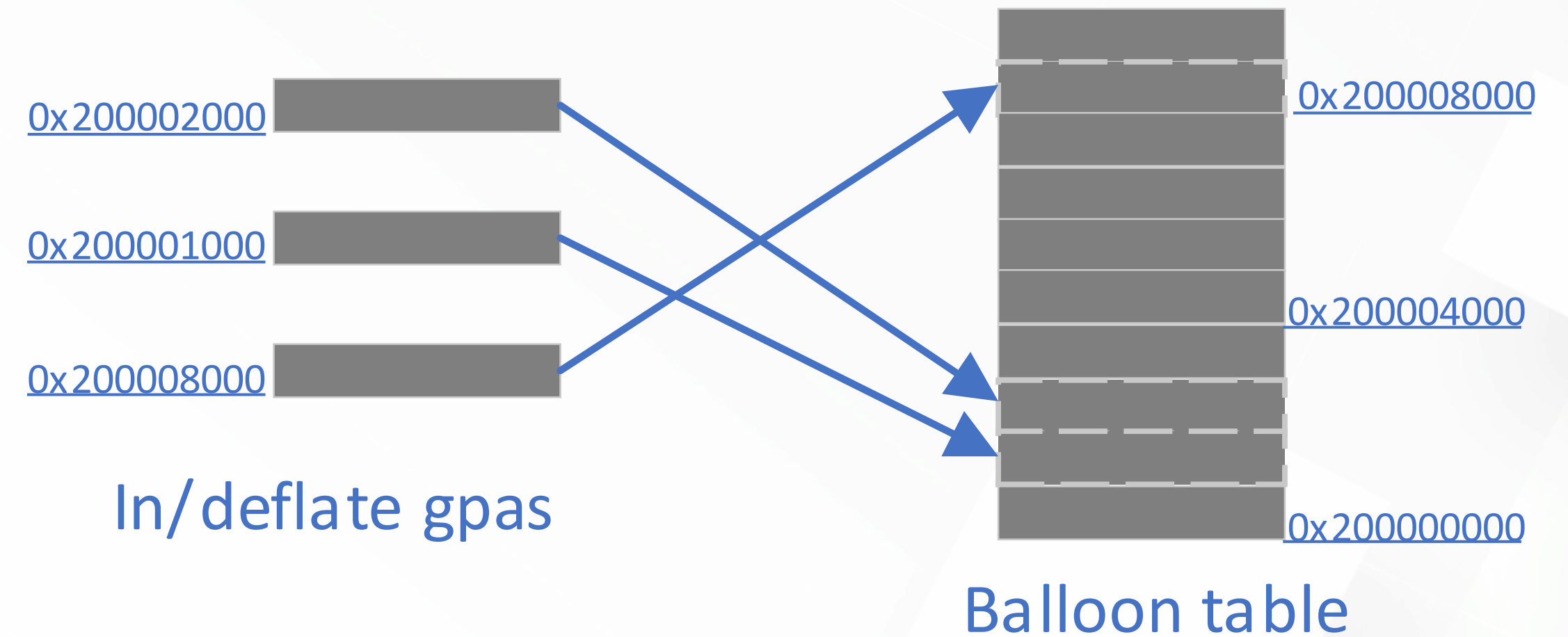
## Auto-combination
**Problem:**
☐ one page at a time

☐ 256 pages  per cycle

**Practice:**
✓ **Combine adjacent pages**

- Most of the memory ranges are adjacent

- dma_map after inflate balloon process finished

0x200002000

0x200001000

0x200008000

In/deflate gpas

0x200008000

0x200004000

0x200000000

Balloon table

# Optimization design

Increase balloon page size

```
#define VIRTIO_BALLOON_ARRAY_PFNS_MAX 256

struct page *balloon_page_alloc(void)
{
    struct page *page = alloc_page(balloon_mapping_gfp_mask() |
                        __GFP_NOMEMALLOC | __GFP_NORETRY |
                        __GFP_NOWARN);
    return page;
}
```
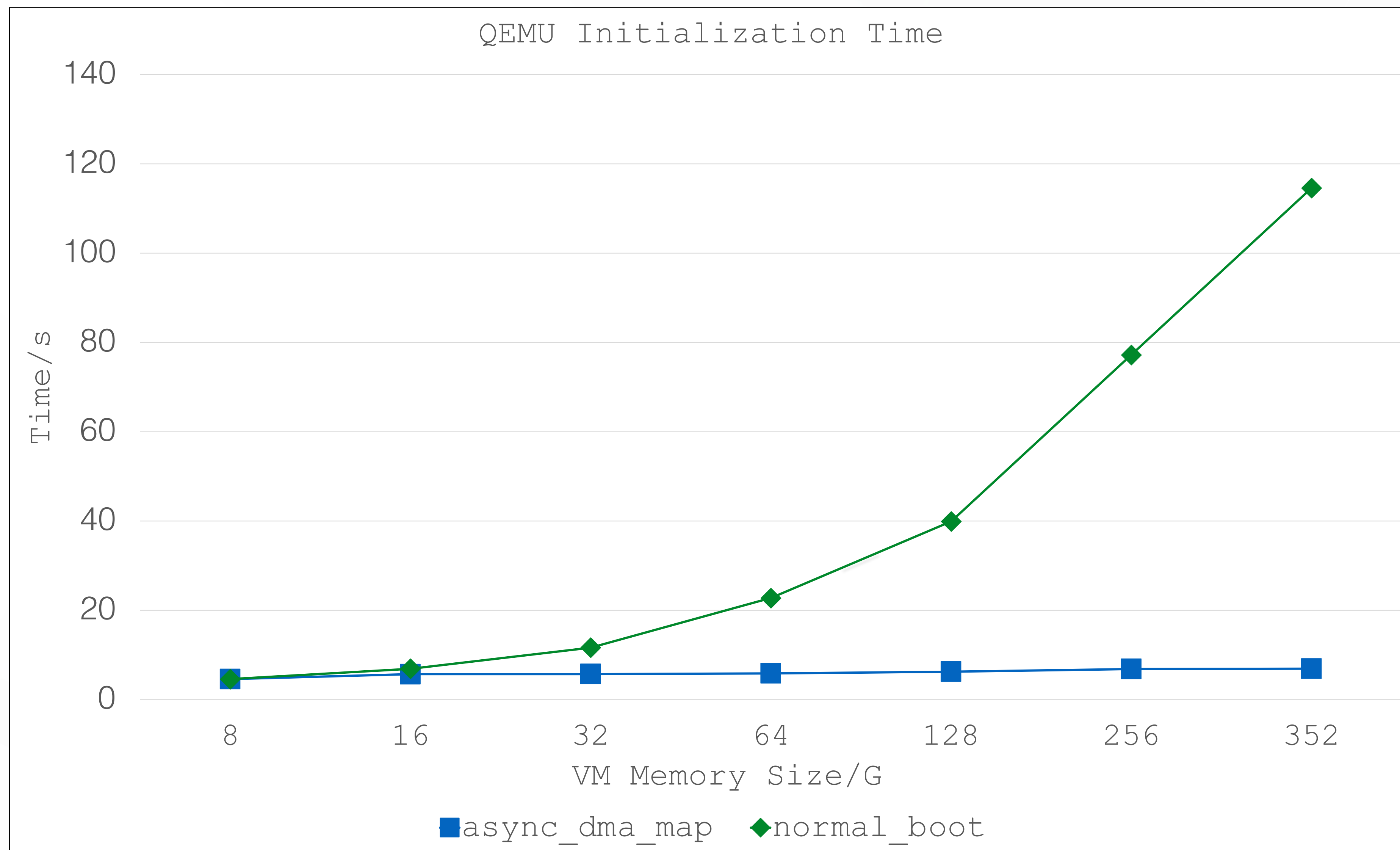
☐ 4K page is too small which will import heavy but unnecessary communication between guest and host

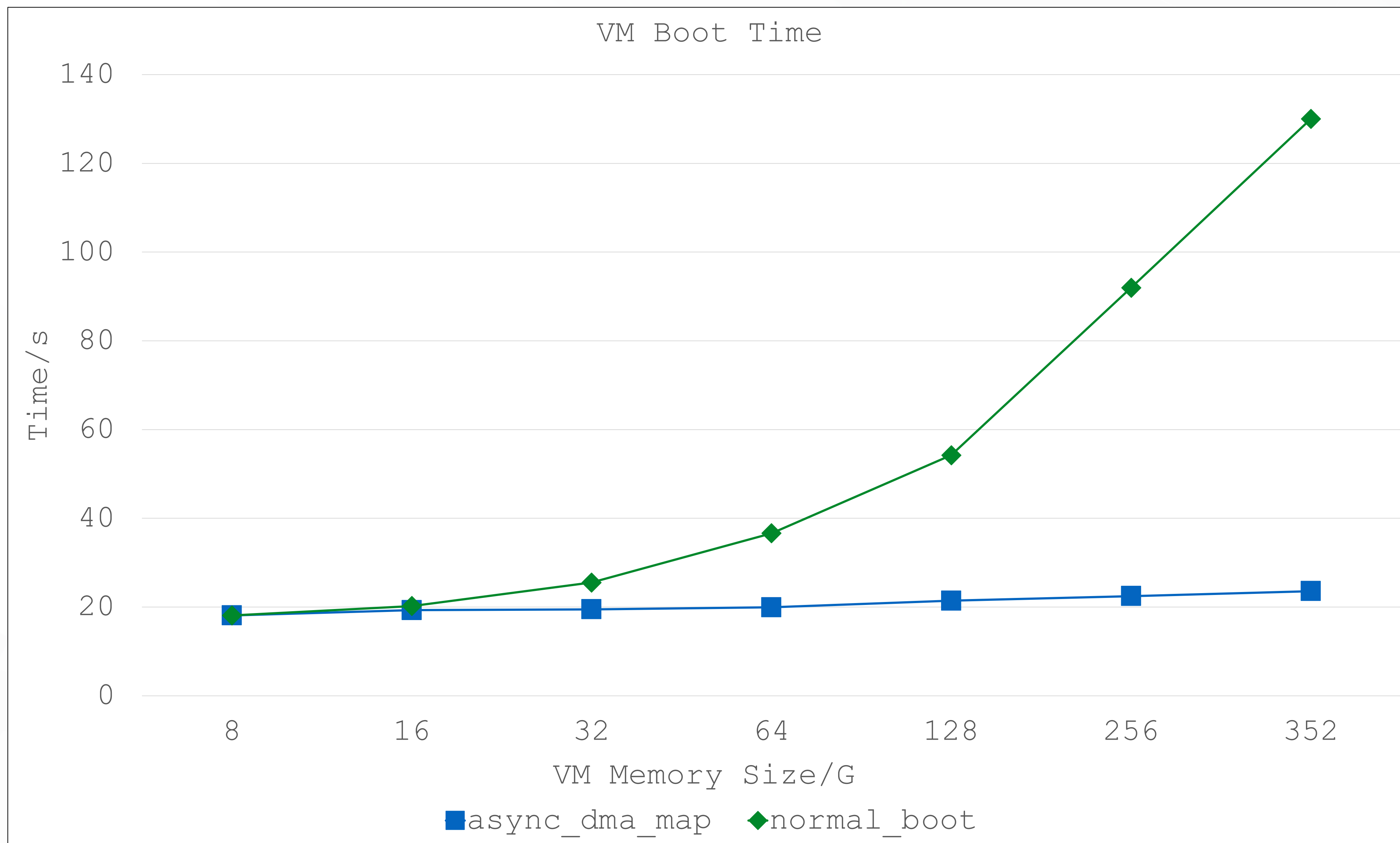✓ 4KB -> 2MB, one virtio talk can in/deflate 2 x 256 = 512MB memory

**Optimization design**

✓ Pre-map to perform dma_map

- Asynchronously dma_map can start early independent of deflating notification

- Insert new dma_map range if the released pages beyond mapped ranges

# Achievements

QEMU Initialization Time

- async_dma_map
- normal_boot

# Achievements



VM Boot Time

Time/s vs VM Memory Size/G

async_dma_map — normal_boot

# Q&A

Li Weinan
william.lwn@alibaba-inc.com

Guo Cheng
hanyu.gc@alibaba-inc.com