# Enabling GPU Memory Oversubscription via Transparent Paging to an NVMe SSD
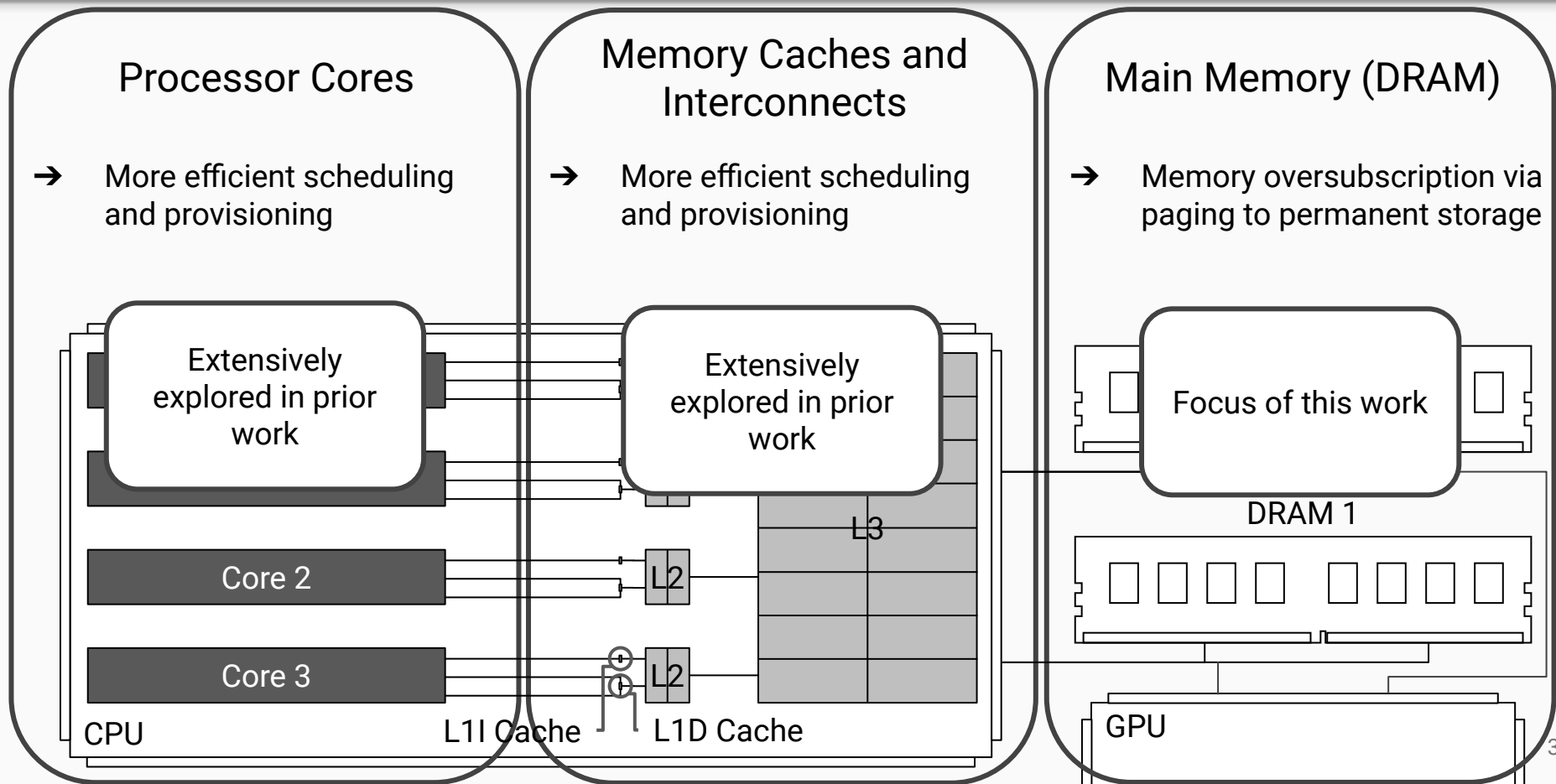
Joshua Bakita and James H. Anderson

Department of Computer Science
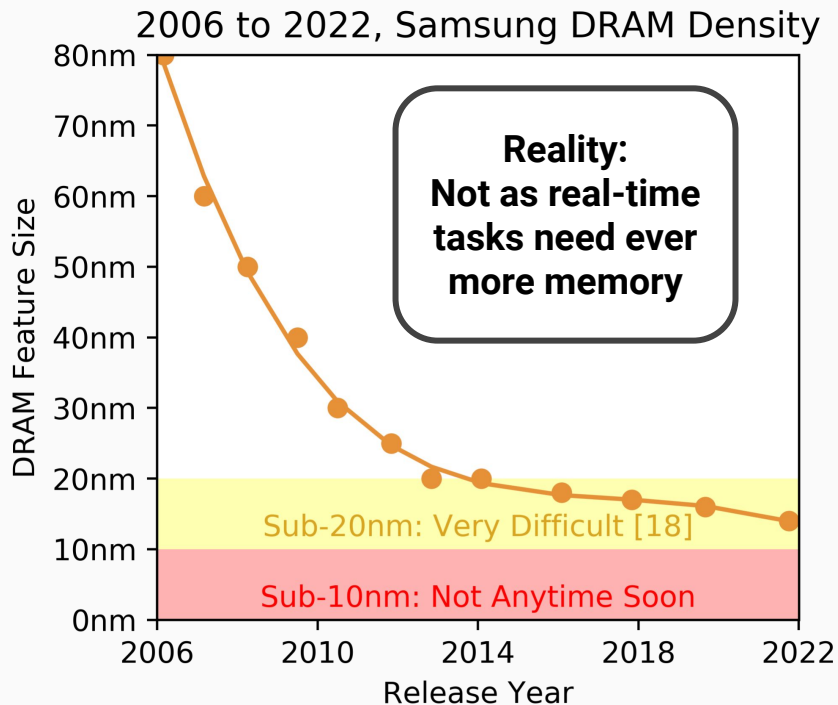University of North Carolina, Chapel Hill
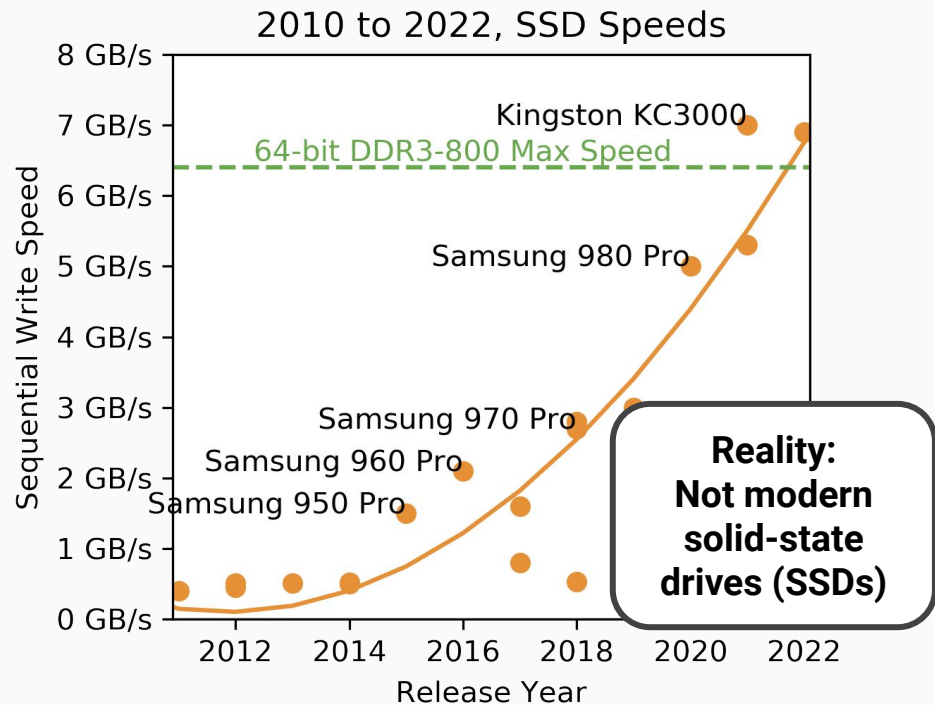
1

# How can we do more, with less?

# How can we do more, with less?

## Processor Cores

→ More efficient scheduling and provisioning

## Memory Caches and Interconnects

→ More efficient scheduling and provisioning

## Main Memory (DRAM)

→ Memory oversubscription via paging to permanent storage

Extensively explored in prior work

Extensively explored in prior work

Focus of this work

Core 2

Core 3

CPU

L2

L2

L3

L1I Cache

L1D Cache

DRAM 1

GPU

## Assumption: DRAM is plentiful



2006 to 2022, Samsung DRAM Density

**Reality: Not as real-time tasks need ever more memory**

Sub-20nm: Very Difficult [18]

Sub-10nm: Not Anytime Soon

## Assumption: Storage is too slow



2010 to 2022, SSD Speeds

Kingston KC3000

64-bit DDR3-800 Max Speed

Samsung 980 Pro

Samsung 970 Pro

Samsung 960 Pro

Samsung 950 Pro

**Reality: Not modern solid-state drives (SSDs)**

# Key Goals

Prior work [15-17, 44-45] limited to this scope

Memory oversubscription for a real-time system that is:

**Predictable**          **Fast**          **Easily Applicable**

With key insights drawn from **technology trends**, **real-time scheduling**, and **GPU architecture**, we achieve all three for real-time CPU+GPU+SSD systems.

# Enabling <u>Predictable</u> Memory Oversubscription

Goal 1 of 3

# Dangers to avoid

Demand paging combined with
least-recently-used (LRU) eviction

Four-Page
DRAM:

| | | | |
|---|---|---|---|
| 1 | 2 | 3 | 4 |

Many-Page
Storage:

| | | | |
|---|---|---|---|
| | | | … |

Page of
Task 1

Page of
Task 2

# Dangers to avoid

Demand paging combined with least-recently-used (LRU) eviction

Four-Page DRAM:

| | | | |
|---|---|---|---|
| 1 | 2 | 3 | 4 |

Many-Page Storage:

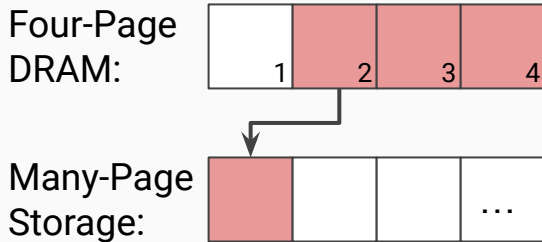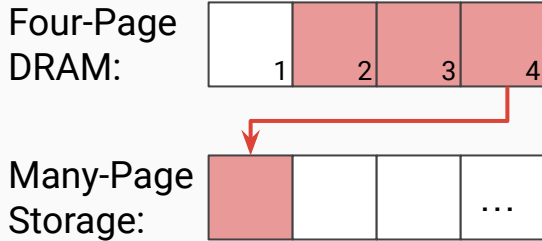| | | | |
|---|---|---|---|
| | | | … |

■ Page of Task 1    ■ Page of Task 2

1.  Task 1 runs a job, accessing pages 3 -> 4 -> 2

2.  Task 2 runs a job, needing two pages of DRAM. The OS selects the LRU page of Task 1 and moves it to storage.

8

# Dangers to avoid

Demand paging combined with least-recently-used (LRU) eviction

Four-Page DRAM:



Many-Page Storage:

Page of Task 1

Page of Task 2

1. Task 1 runs a job, accessing pages 3 -> 4 -> 2

2. Task 2 runs a job, needing two pages of DRAM. The OS selects the LRU page of Task 1 and moves it to storage.

Case 1:
**Page 2**

# Dangers to avoid

Demand paging combined with least-recently-used (LRU) eviction

Four-Page DRAM:

Many-Page Storage:

■ Page of Task 1    ■ Page of Task 2

1. Task 1 runs a job, accessing pages 3 -> 4 -> 2

2. Task 2 runs a job, needing two pages of DRAM. The OS selects the LRU page of Task 1 and moves it to storage.
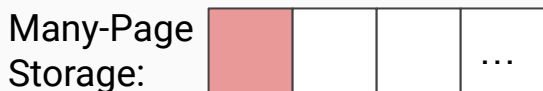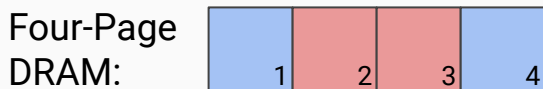
Case 1:
**Page 2**

Case 2:
**Page 4**

# Dangers to avoid

Demand paging combined with least-recently-used (LRU) eviction

Four-Page DRAM:

| | | | |
|---|---|---|---|
| 1 | 2 | 3 | 4 |

Many-Page Storage:

| | | | |
|---|---|---|---|
| | | | … |

■ Page of Task 1    ■ Page of Task 2

1. Task 1 runs a job, accessing pages 3 -> 4 -> 2

2. Task 2 runs a job, needing two pages of DRAM. The OS selects the LRU page of Task 1 and moves it to storage.

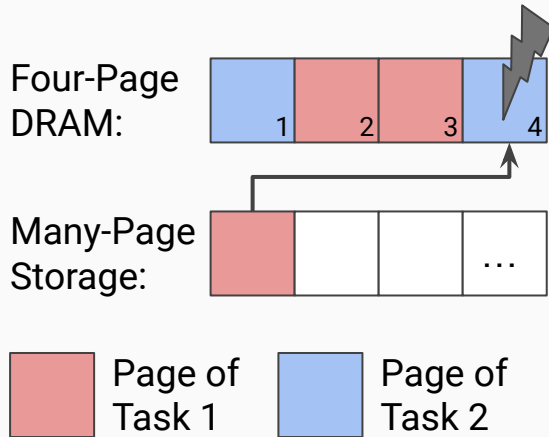   Case 1:          Case 2:
   **Page 2**       **Page 4**

3. Task 2 runs its job to completion.

## Predictable Oversubscription
# Dangers to avoid

Demand paging combined with least-recently-used (LRU) eviction

Four-Page DRAM:

Many-Page Storage:

Page of Task 1

Page of Task 2

1. Task 1 runs a job, accessing pages 3 -> 4 -> 2

2. Task 2 runs a job, needing two pages of DRAM. The OS selects the LRU page of Task 1 and moves it to storage.

   Case 1:          Case 2:
   **Page 2**       **Page 4**

3. Task 2 runs its job to completion.

4. Task 1 runs its next job, accessing pages 3 and 4. It's execution time will vary greatly depending on what was moved to storage.
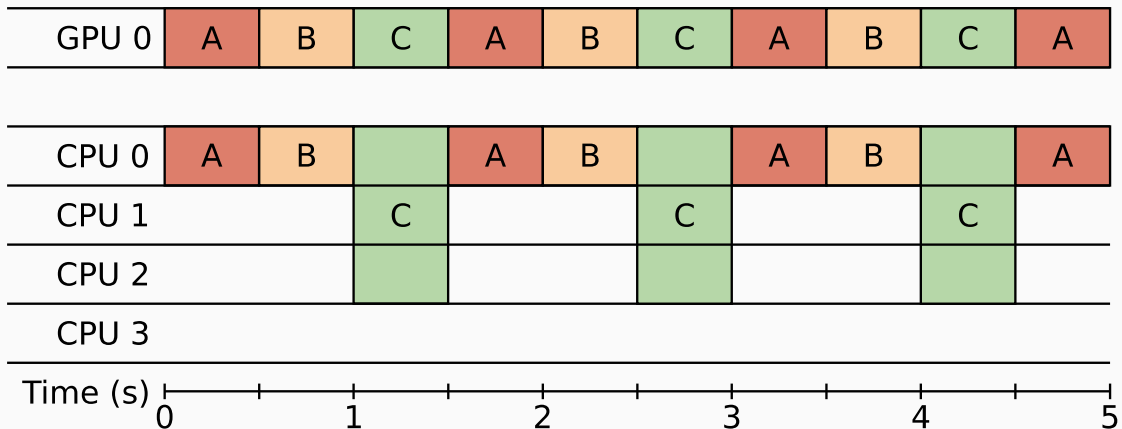
# **Predictable** Oversubscript.
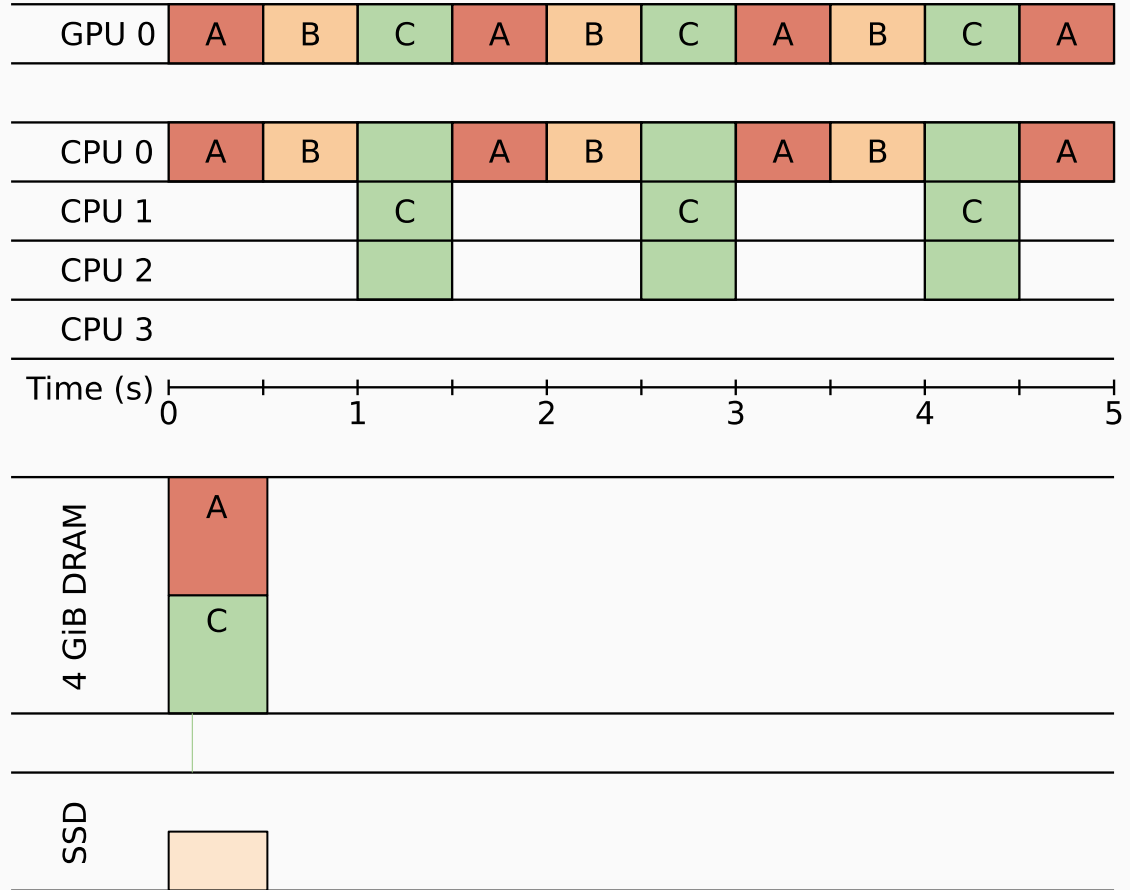
## Using schedule foreknowledge

After a job completes in a real-time system, we know the minimum amount of time before the next job arrives.

With a table-driven scheduler, we know exactly.

Consider an example…

Three Components: A, B, C
500ms budget, 1500ms period, 2GiB memory each

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| GPU 0 | A | B | C | A | B | C | A | B | C | A |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| CPU 0 | A | B | | A | B | | A | B | | A |
| CPU 1 | | | C | | | C | | | C | |
| CPU 2 | | | | | | | | | | |
| CPU 3 | | | | | | | | | | |

Time (s)  0    1    2    3    4    5

Works fine, given 6 GiB of DRAM.

What if we only have 4 GiB?

## **Predictable** Oversubscript.
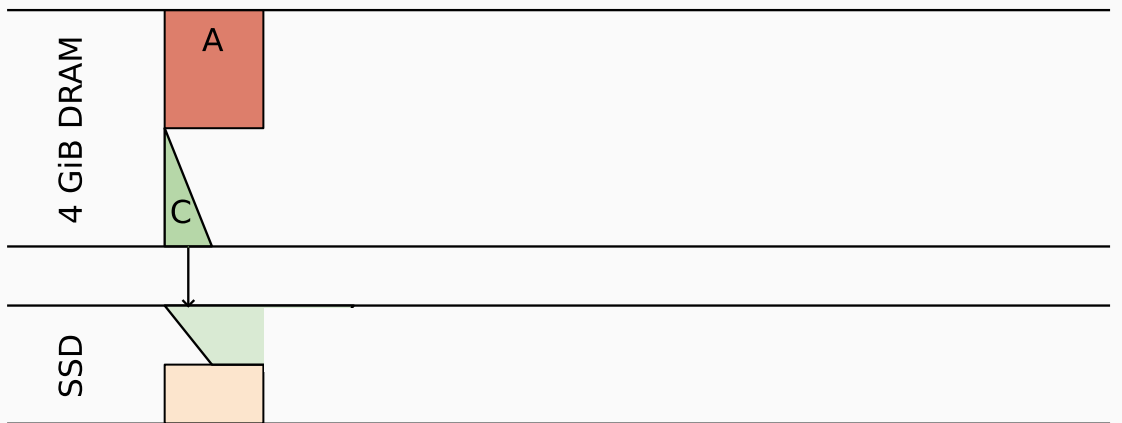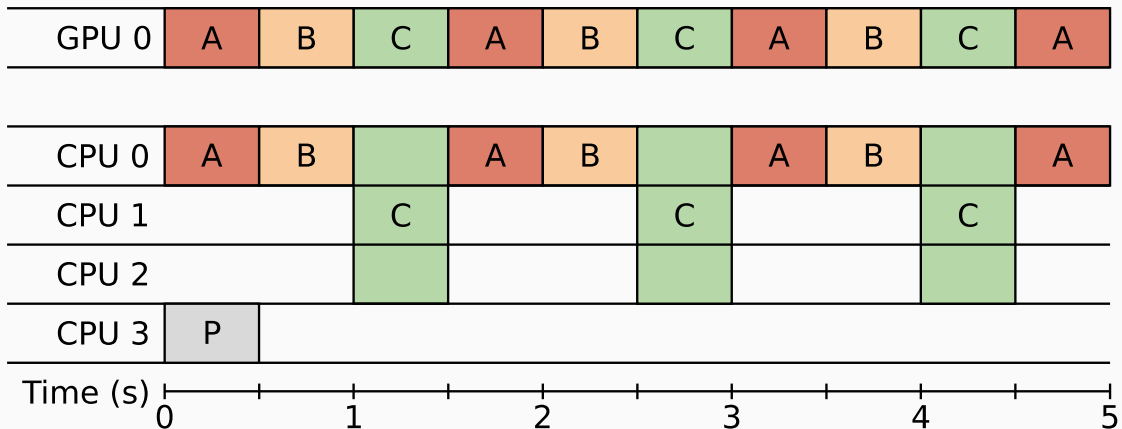
# Using schedule foreknowledge

After a job completes in a real-time system, we know the minimum amount of time before the next job arrives.

With a table-driven scheduler, we know exactly.

Consider an example...

Creating a memory schedule.

Three Components: A, B, C
500ms budget, 1500ms period, 2GiB memory each

## **Predictable** Oversubscript.
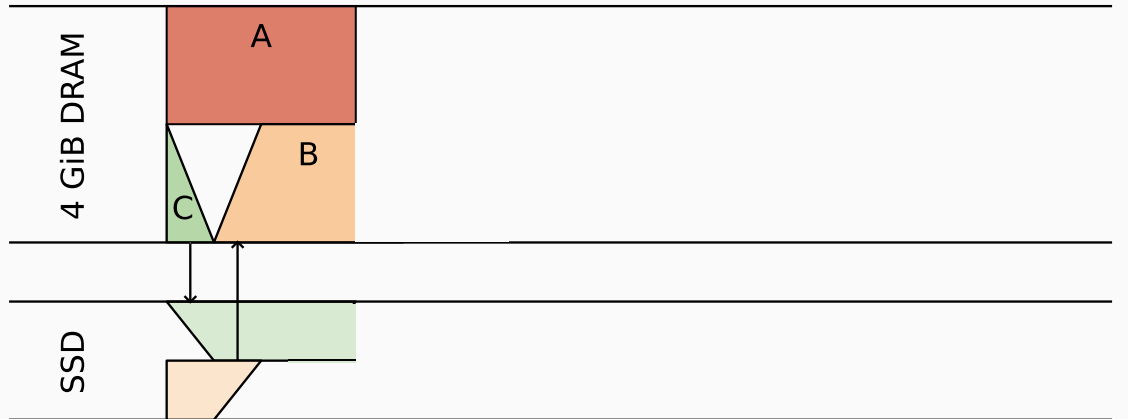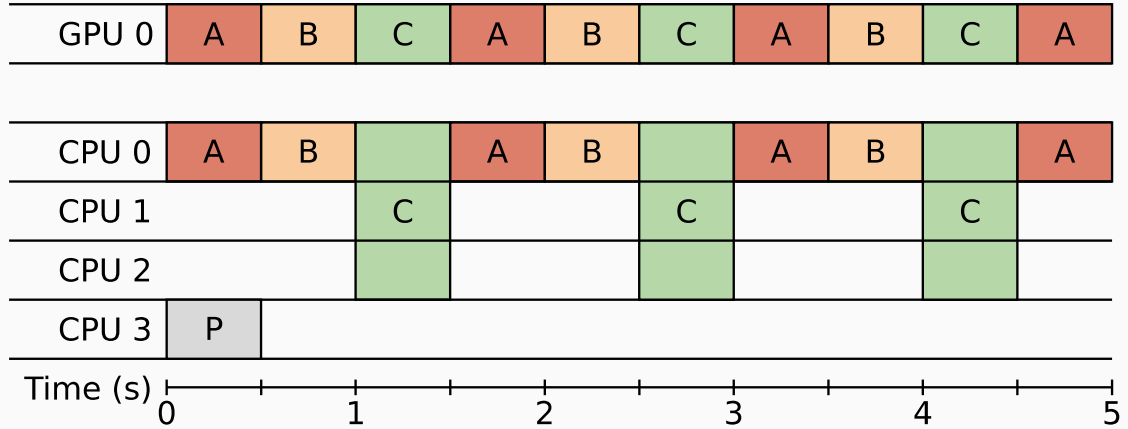# Using schedule foreknowledge

After a job completes in a real-time system, we know the minimum amount of time before the next job arrives.

With a table-driven scheduler, we know exactly.

Consider an example…

Creating a memory schedule.

### Three Components: A, B, C
### 500ms budget, 1500ms period, 2GiB memory each

| GPU 0 | A | B | C | A | B | C | A | B | C | A |
|-------|---|---|---|---|---|---|---|---|---|---|

| CPU 0 | A | B |   | A | B |   | A | B |   | A |
|-------|---|---|---|---|---|---|---|---|---|---|
| CPU 1 |   |   | C |   |   | C |   |   | C |   |
| CPU 2 |   |   | C |   |   | C |   |   | C |   |
| CPU 3 | P |   |   |   |   |   |   |   |   |   |

Time (s): 0   1   2   3   4   5

4 GiB DRAM: A, C

SSD

# **Predictable** Oversubscript.
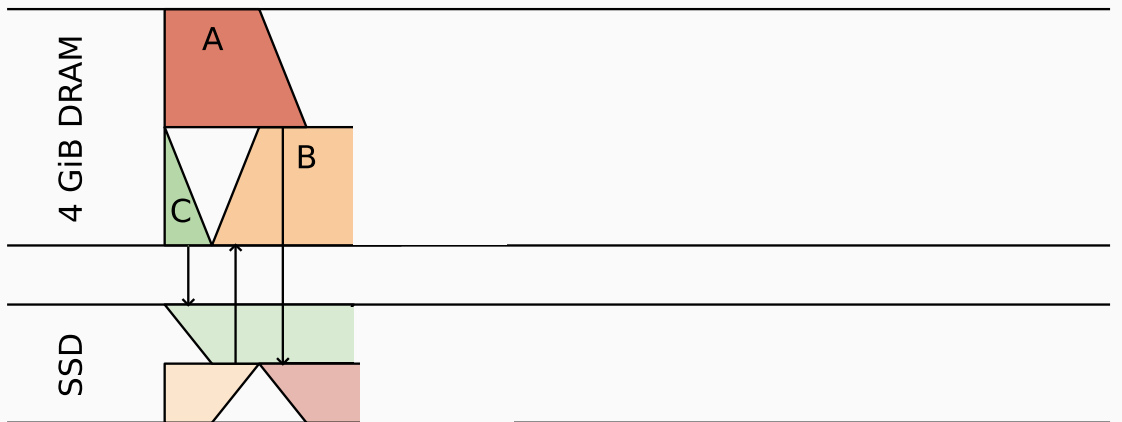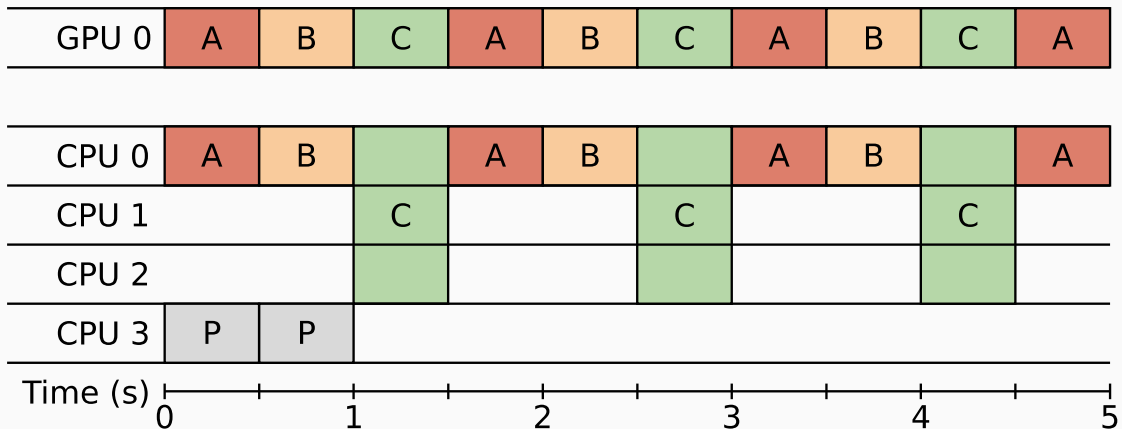
## Using schedule foreknowledge

After a job completes in a real-time system, we know the minimum amount of time before the next job arrives.

With a table-driven scheduler, we know exactly.

Consider an example...

Creating a memory schedule.

Three Components: A, B, C
500ms budget, 1500ms period, 2GiB memory each

# **Predictable** Oversubscript.

# Using schedule foreknowledge

After a job completes in a real-time system, we know the minimum amount of time before the next job arrives.

With a table-driven scheduler, we know exactly.

Consider an example...

Creating a memory schedule.



Three Components: A, B, C
500ms budget, 1500ms period, 2GiB memory each

**Predictable Oversubscript.**
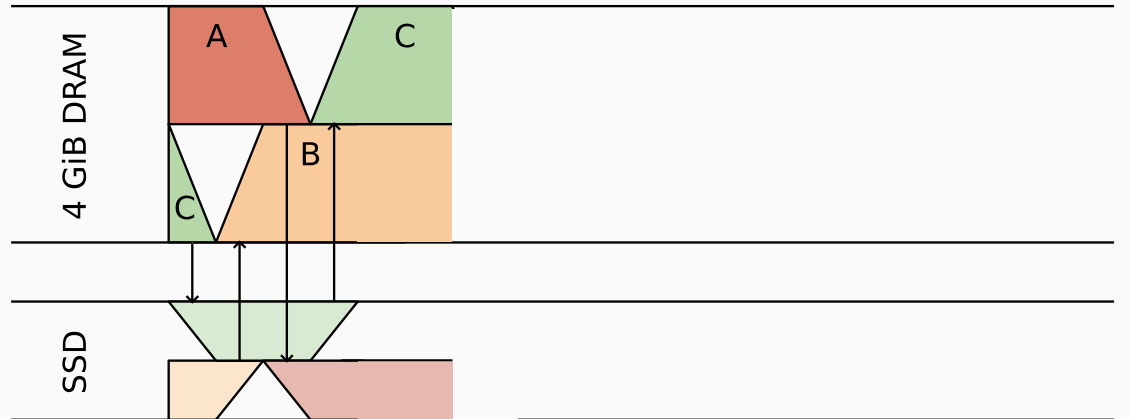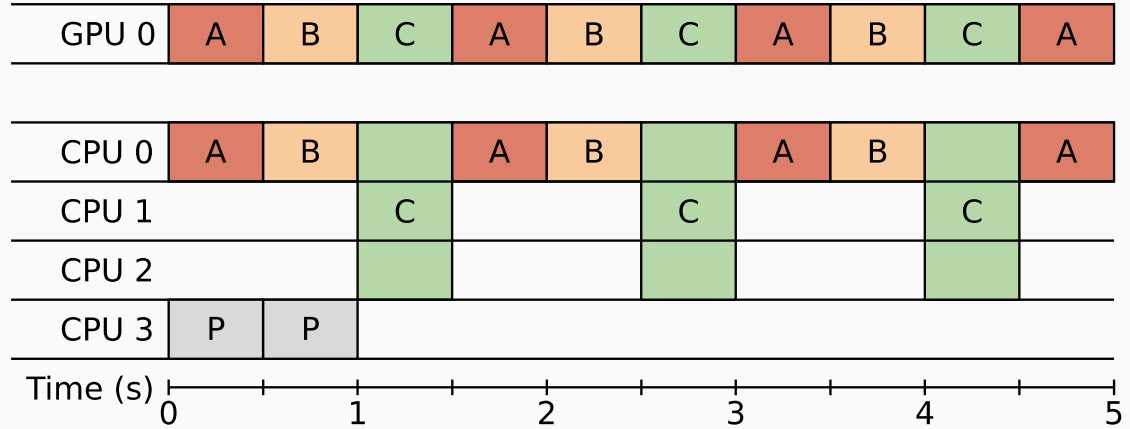
# Using schedule foreknowledge

After a job completes in a real-time system, we know the minimum amount of time before the next job arrives.

With a table-driven scheduler, we know exactly.

Consider an example...

Creating a memory schedule.

Three Components: A, B, C
500ms budget, 1500ms period, 2GiB memory each

# **Predictable** Oversubscript.
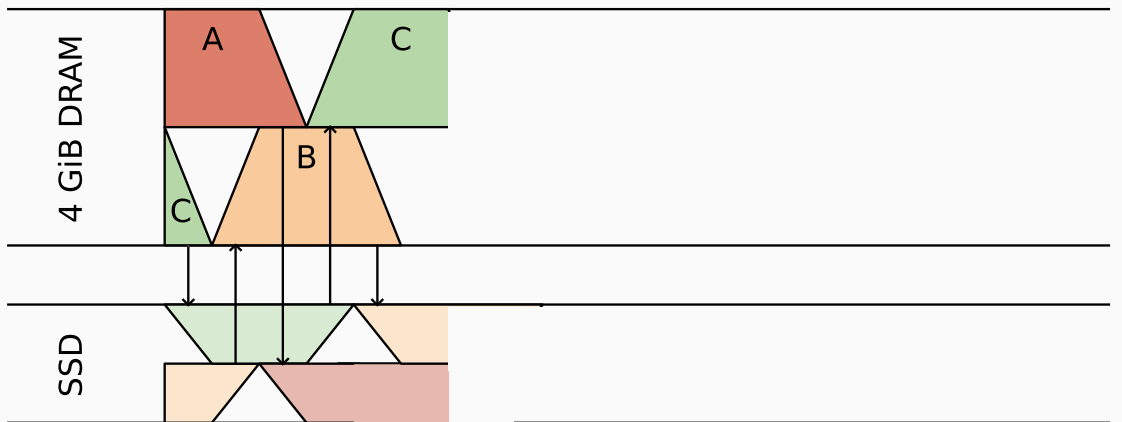
## Using schedule foreknowledge

After a job completes in a real-time system, we know the minimum amount of time before the next job arrives.

With a table-driven scheduler, we know exactly.

Consider an example...

Creating a memory schedule.

Three Components: A, B, C
500ms budget, 1500ms period, 2GiB memory each

**Predictable Oversubscript.**
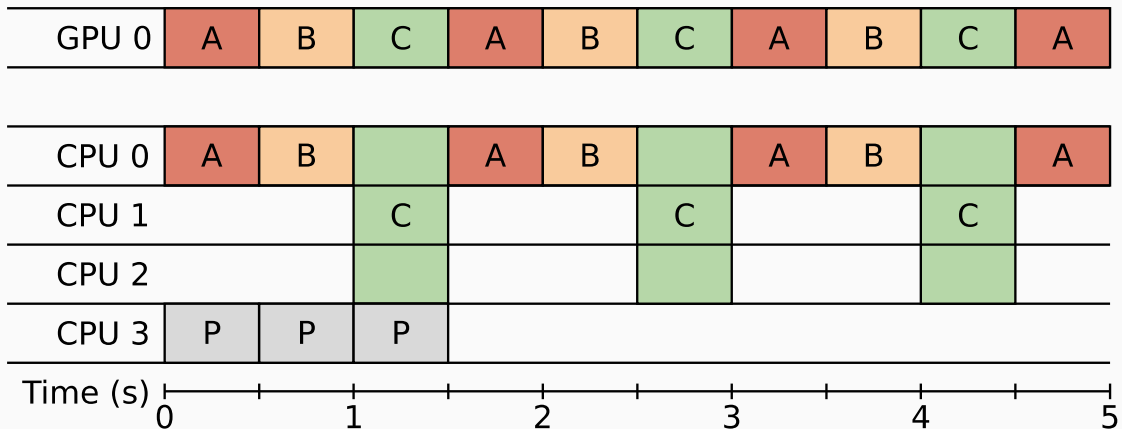
# Using schedule foreknowledge

After a job completes in a real-time system, we know the minimum amount of time before the next job arrives.

With a table-driven scheduler, we know exactly.

Consider an example…

Creating a memory schedule.

Three Components: A, B, C
500ms budget, 1500ms period, 2GiB memory each

# **Predictable** Oversubscript.

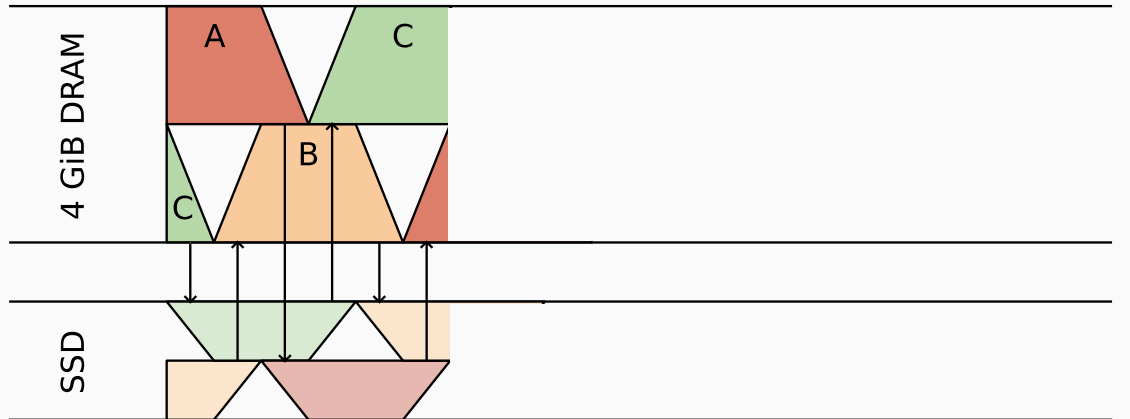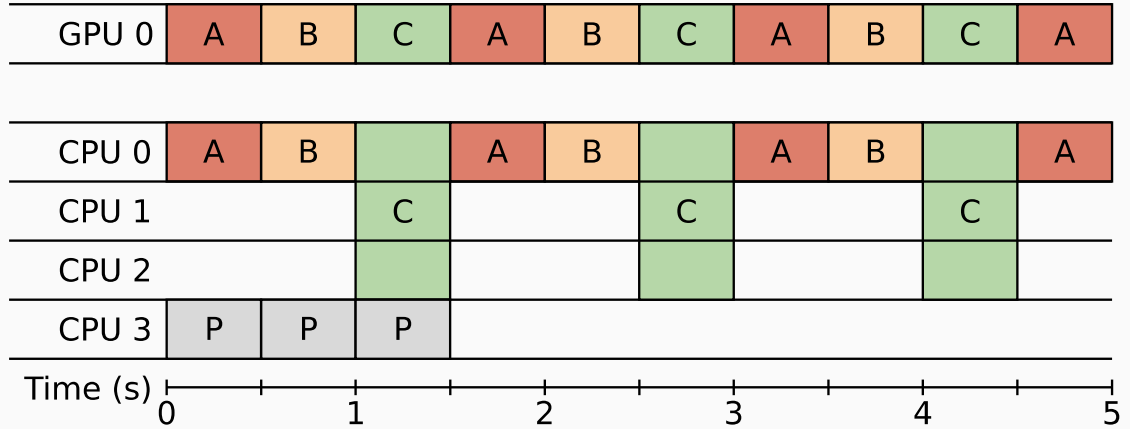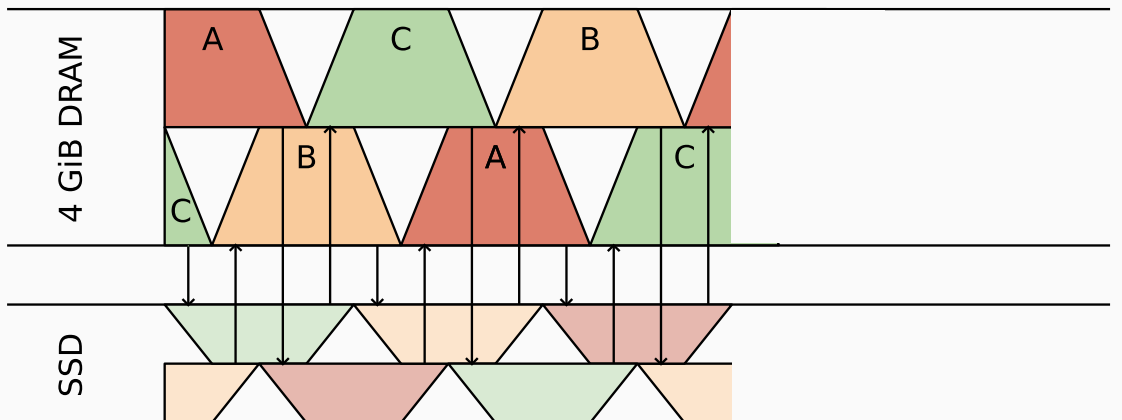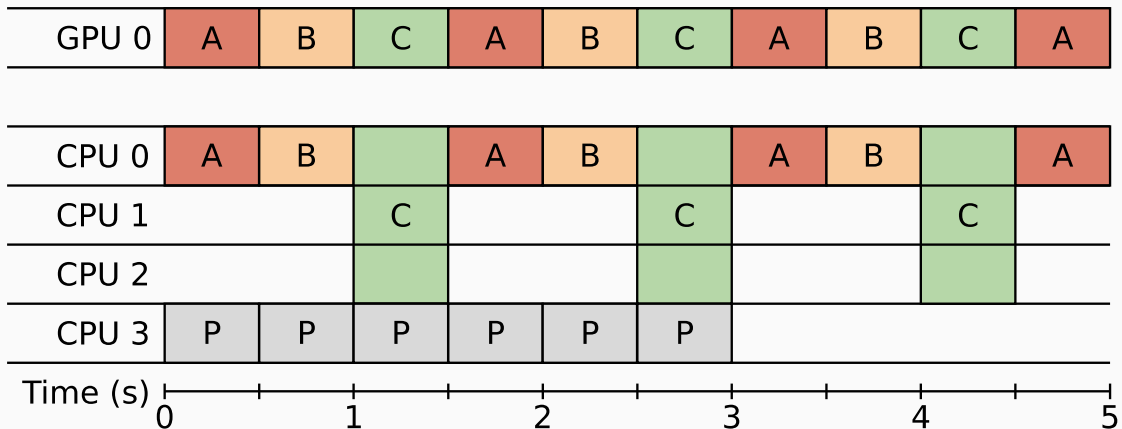## Using schedule foreknowledge

After a job completes in a real-time system, we know the minimum amount of time before the next job arrives.

With a table-driven scheduler, we know exactly.

Consider an example…

Creating a memory schedule.

Three Components: A, B, C
500ms budget, 1500ms period, 2GiB memory each

# Predictable Oversubscript.

## Using schedule foreknowledge

After a job completes in a real-time system, we know the minimum amount of time before the...

W...
kn...

Co...

Creating a memory schedule.

**Must be very fast (8 GiB/s in this example)**

**Key Insight:**
By using period information to schedule paging operations, we can make them predictable.

Three Components: A, B, C
500ms budget, 1500ms period, 2GiB memory each

22

# Enabling <u>Fast</u> Memory Oversubscription

Goal 2 of 3

# **Fast** Oversubscription

## What's already available?

Demand paging for loading data from storage.

Platform: NVIDIA Jetson AGX Xavier

Three Components: A, B, C
500ms budget, 1500ms period, 2GiB memory each



→ Need 8 GiB/s to meet needs of example

Unacceptable

# What slows demand paging?

**Key Insight:**
Synchronization costs make paging multicore CPU memory unacceptable in a real-time Linux system

Overheads?

Profiling results

## 14%
Page allocation

## 19%
Page mapping

## 27%
Bookkeeping and actual I/O

## 40%
Locking and retry commit

On Jetson Xavier AGX running Linux 4.9 with our Sabrenet Rocket 4 Plus SSD.

# Paging GPU memory

**Key Insight:**
GPU APIs do not allow for cross-application shared pages

**Key Insight:**
Pages can be moved directly from GPU virtual memory to and from an SSD without a mapping on the CPU

**Key Insight:**
GPU memory is most commonly used to store *read-only* weights

Our GPU paging method is 3 times faster than demand paging

# Fast Oversubscription

## Utilizing the SSD controller

SSDs support command offloading. Can we use this instead of the CPU?
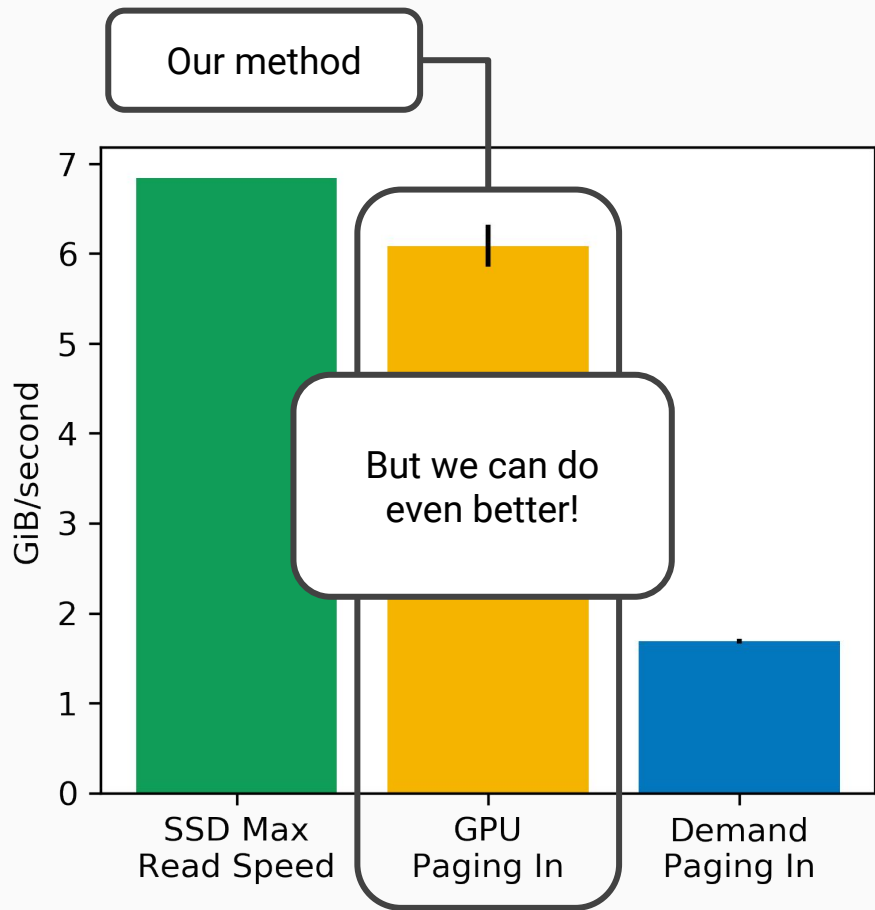
Allows for:
➔ 30% asynchronous page-in
➔ 80% asynchronous page-out

**Physical DRAM Pages**

| | 1 | 5 | 9 | 13 | 17 | 21 |
| | 2 | | | | | |
| | 3 | | | | | |
| | 4 | | | | | |

other cmd

From pages: 7, 8, 5
To sector: 55

Write Command Contents

Key Insight:
SSD controllers can process much of a paging operation asynchronously

MMU

CPUs

Misc PCIe Devices

PCIe Bus

NVMe SSD

① SSD Queue Setup
② Write Command Queued
③ SSD Executes Write Asynchronously

Scatter-gather write operation on an NVMe SSD

1. Allocate Pages
2. Map to I/O MMU
3. Initiate Transfer

5. Map Pages to GPU
6. Flush GPU TLB+L2
7. Complete Transfer

Asynchronous

Time (ms)    0    50    100    150

4. SSD Controller Reads Data

Timeline for GPU page-in operation

# Enabling <u>Easily Applicable</u> Memory Oversubscription

Goal 3 of 3

**Easily Applicable** Oversubscription

# A simple API

On Linux:

0.  Use `strace` to identify address space (*AS_ID*) and allocation ID (*BUF_ID*)
1.  `ioctl(`*AS_ID*`, NVGPU_AS_IOCTL_WRITE_SWAP_BUFFER, {`*BUF_ID*`});`
2.  `ioctl(`*AS_ID*`, NVGPU_AS_IOCTL_READ_SWAP_BUFFER, {`*BUF_ID*`});`

Asynchronous variants available.

Code is open source and documented. See
https://www.cs.unc.edu/~jbakita/rtss22-ae.html to get started.

# Conclusions

**We can make memory oversubscription in a real-time system:**

**<u>Predictable</u>**                                 **<u>Fast</u>**                               **<u>Easily Applicable</u>**

How can we know what          How can synchronization          Can we make GPU paging
and when to page?                 overheads be avoided?                 easy to use?

Use foreknowledge                 Page GPU mem                    Yes, via our 2-line
present in schedule                                                        Linux API

                                         How fast can we go?

                                         >6GB/s read,
                                         >5GB/s write

# What you have to read the paper for…

Evaluation:
- Comparison to direct I/O, and how we manage to be faster
- Benchmarks which demonstrate our minimal impact on memory bandwidth
- Exact distributions for all benchmark results
- Full details on our supported API calls

Regarding SSDs:
- Details on how we offload paging operations onto the SSD controller
- How we ensure SSD caches don't bottleneck
- How we utilize real-time GPU scheduling invariants to speed up page-out operations

Regarding GPUs:
- Details on how memory allocations work on NVIDIA's embedded boards
- Details on NVIDIA GPU virtual memory capabilities
- History of NVIDIA page table formats
- How to determine GPU address space and buffer IDs
- A version of strace supporting detailed tracing of all NVIDIA driver syscalls on Jetson boards

+ More details and background on everything covered in this presentation

# Thanks!
# Questions?

Future work:
➔ OS scheduler integration
➔ Application to mode changes, DNN layers, etc.
➔ Increase performance, portability, and SSD space allocation algorithm

Contact:
Email: jbakita@cs.unc.edu
Twitter: @JJBakita
Web: https://cs.unc.edu/~jbakita

Old Well, University of North Carolina at Chapel Hill, Winter 2017