

# Nvshare

A Kubernetes Use Case

George Alexopoulos <giorgosalexo0@gmail.com>



# Interactive environment



# Interactive environment



user



# Interactive environment





# Interactive environment



user

write

```
In [ ]: # Initialize denominator
k = 1

# Initialize sum
s = 0

for i in range(1000000):
    # even index elements are positive
    if i % 2 == 0:
        s += 4/k
    else:
        # odd index elements are negative
        s -= 4/k

    # denominator is odd
    k += 2

print("pi =", s)
```

Browse



```
In [ ]: print("Do something else.")
```



# Interactive environment



user

write

run

cell

Browse



```
In [ ]: # Initialize denominator
k = 1

# Initialize sum
s = 0

for i in range(1000000):

    # even index elements are positive
    if i % 2 == 0:
        s += 4/k
    else:
        # odd index elements are negative
        s -= 4/k

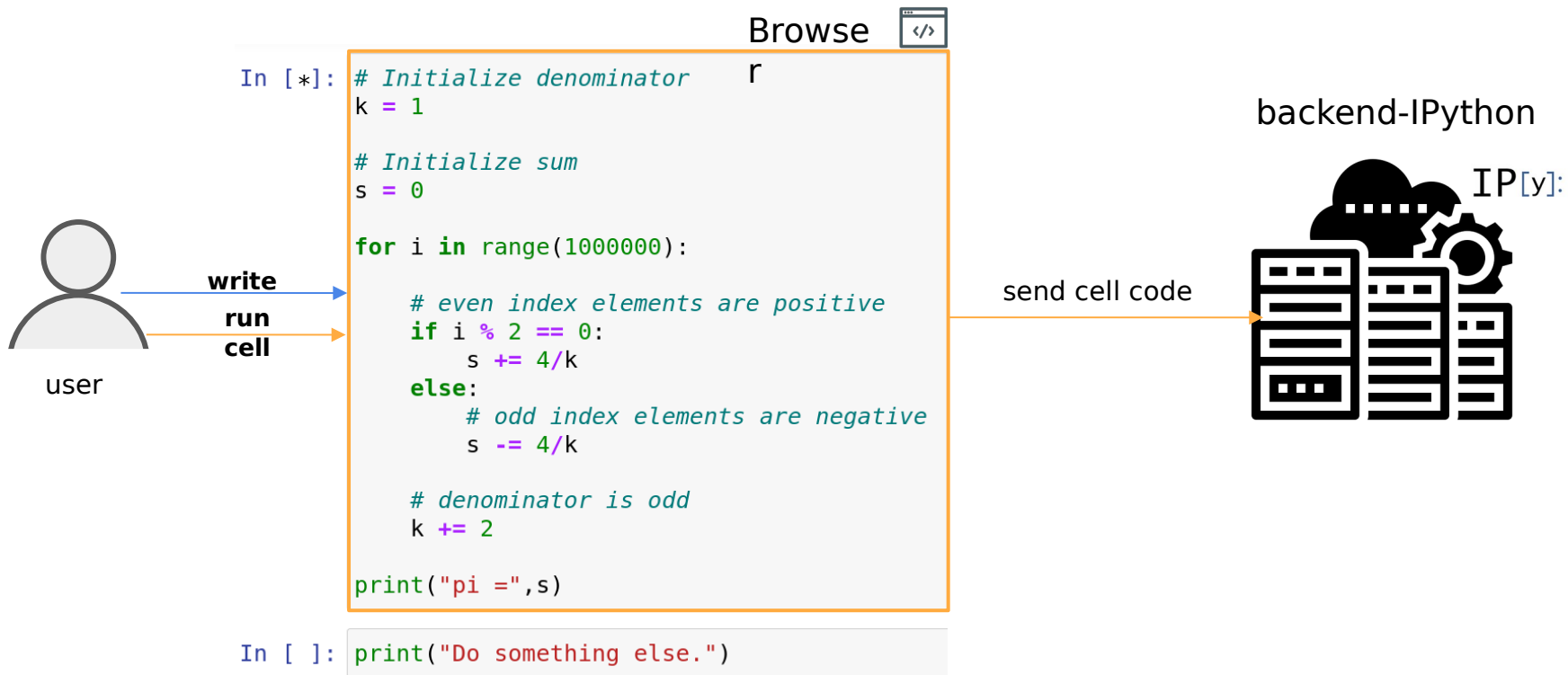
    # denominator is odd
    k += 2

print("pi =", s)
```

```
In [ ]: print("Do something else.")
```

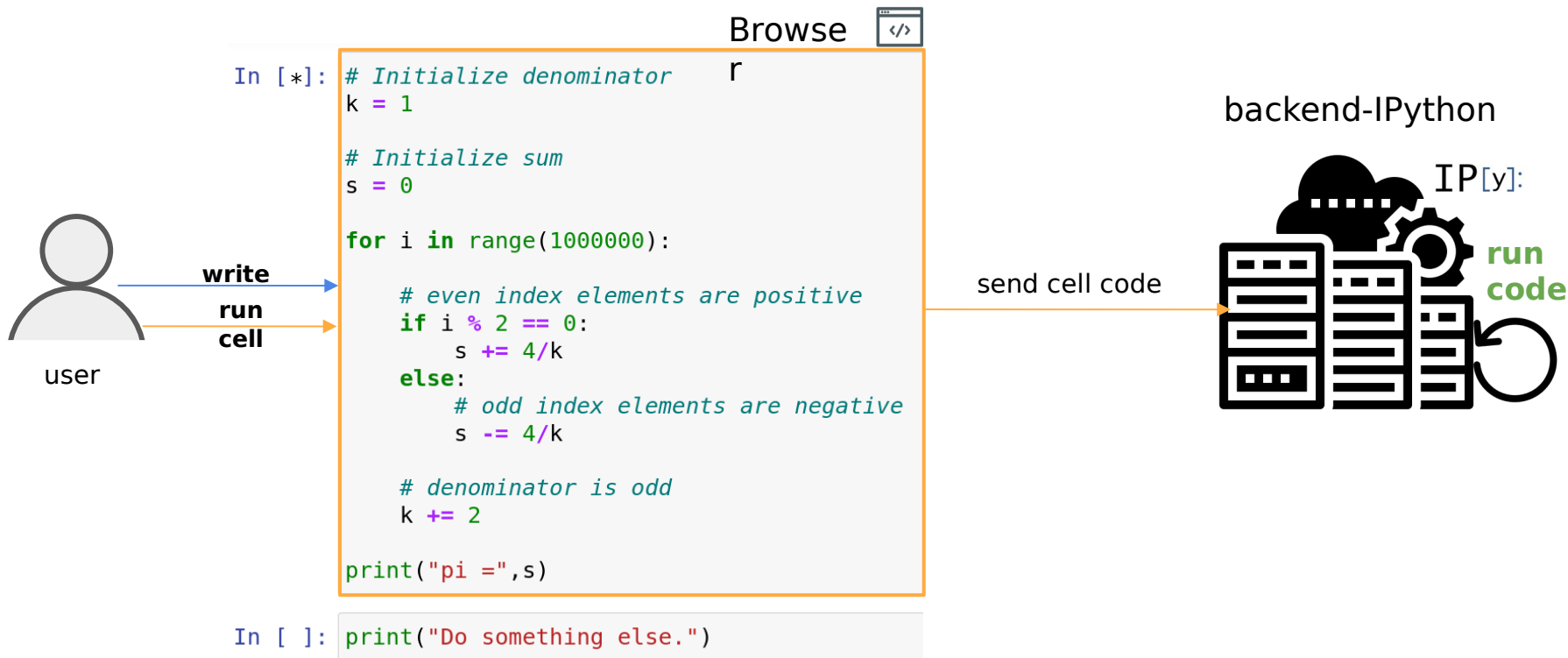


# Interactive environment





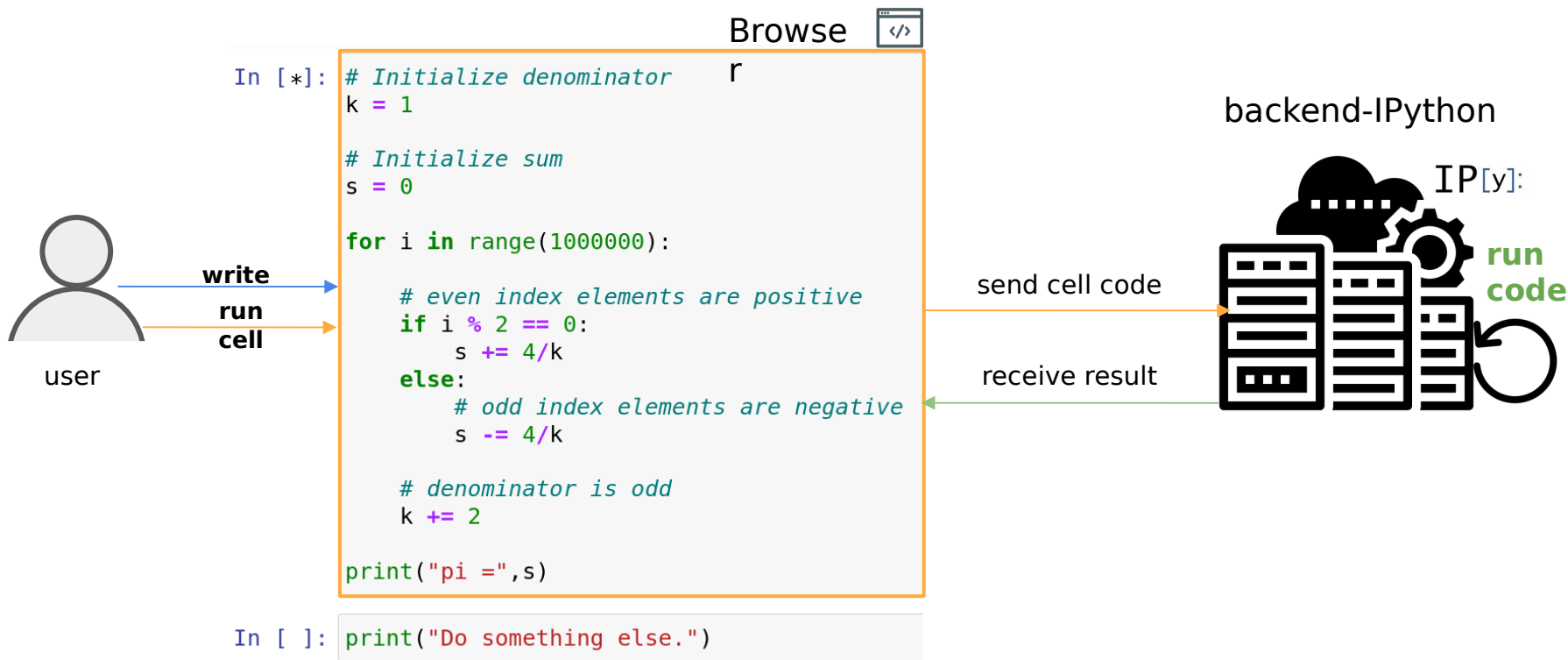
# Interactive environment





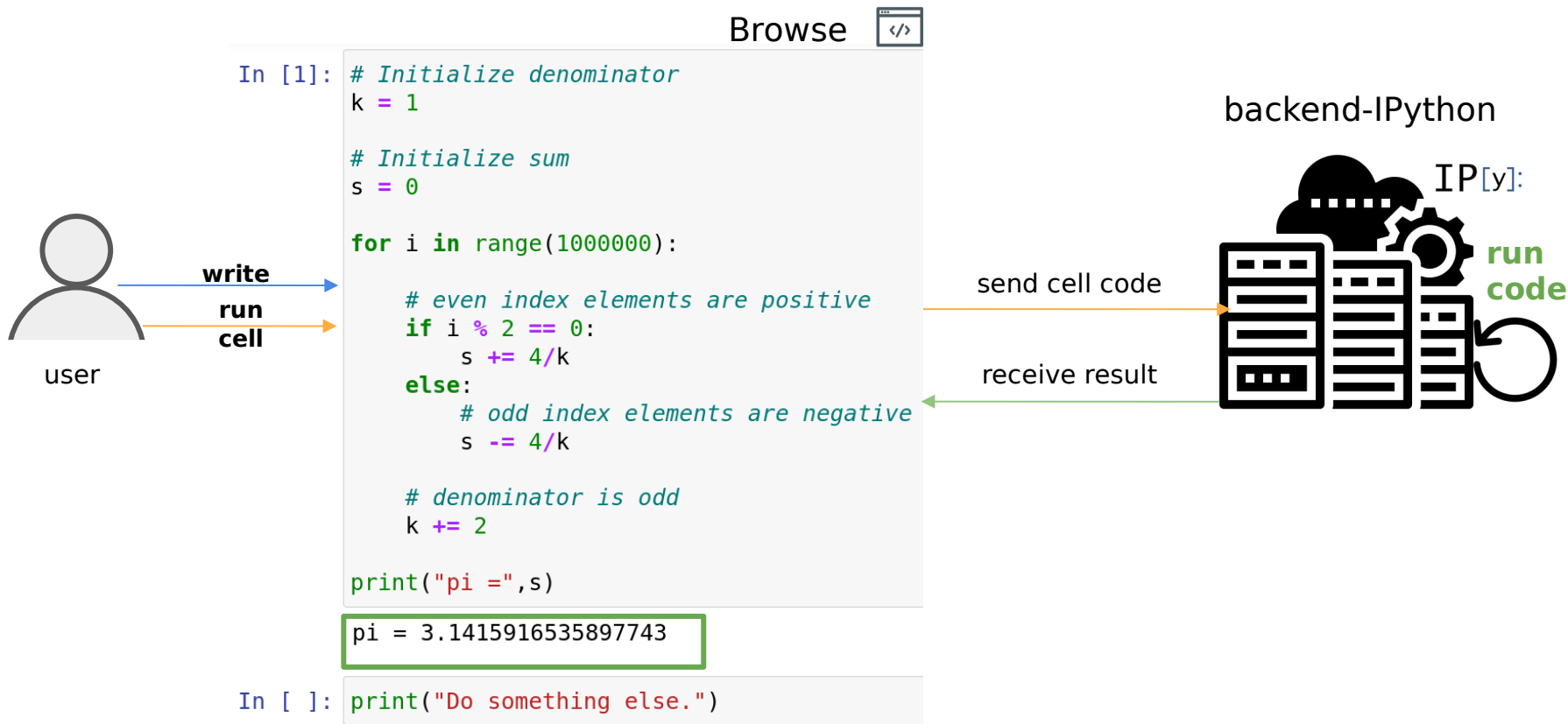


# Interactive environment





# Interactive environment





# Interactive environment

Browse



In [1]: `# Initialize denominator`

```
k = 1
```

```
# Initialize sum
```

```
s = 0
```

```
for i in range(1000000):
```

```
    # even index elements are positive
```

```
    if i % 2 == 0:
```

```
        s += 4/k
```

```
    else:
```

```
        # odd index elements are negative
```

```
        s -= 4/k
```

```
    # denominator is odd
```

```
    k += 2
```

```
print("pi =", s)
```

```
pi = 3.1415916535897743
```

In [ ]: `print("Do something else.")`

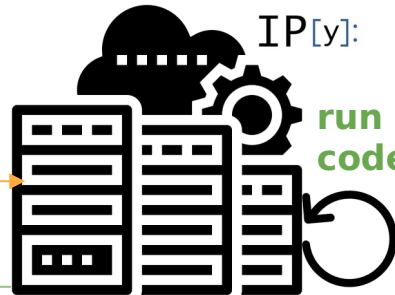
backend-IPython

IP[y]:

run  
code

send cell code

receive result





# Interactive environment



user

write

run

cell

run  
cell

Browse



```
In [1]: # Initialize denominator
```

```
k = 1
```

```
# Initialize sum
```

```
s = 0
```

```
for i in range(1000000):
```

```
    # even index elements are positive
```

```
    if i % 2 == 0:
```

```
        s += 4/k
```

```
    else:
```

```
        # odd index elements are negative
```

```
        s -= 4/k
```

```
    # denominator is odd
```

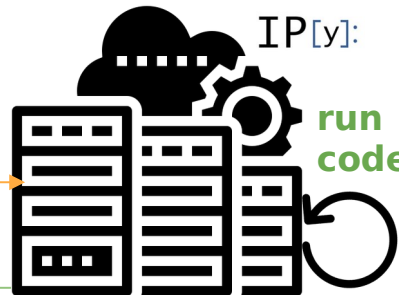
```
    k += 2
```

```
print("pi =", s)
```

```
pi = 3.1415916535897743
```

```
In [ ]: print("Do something else.")
```

backend-IPython

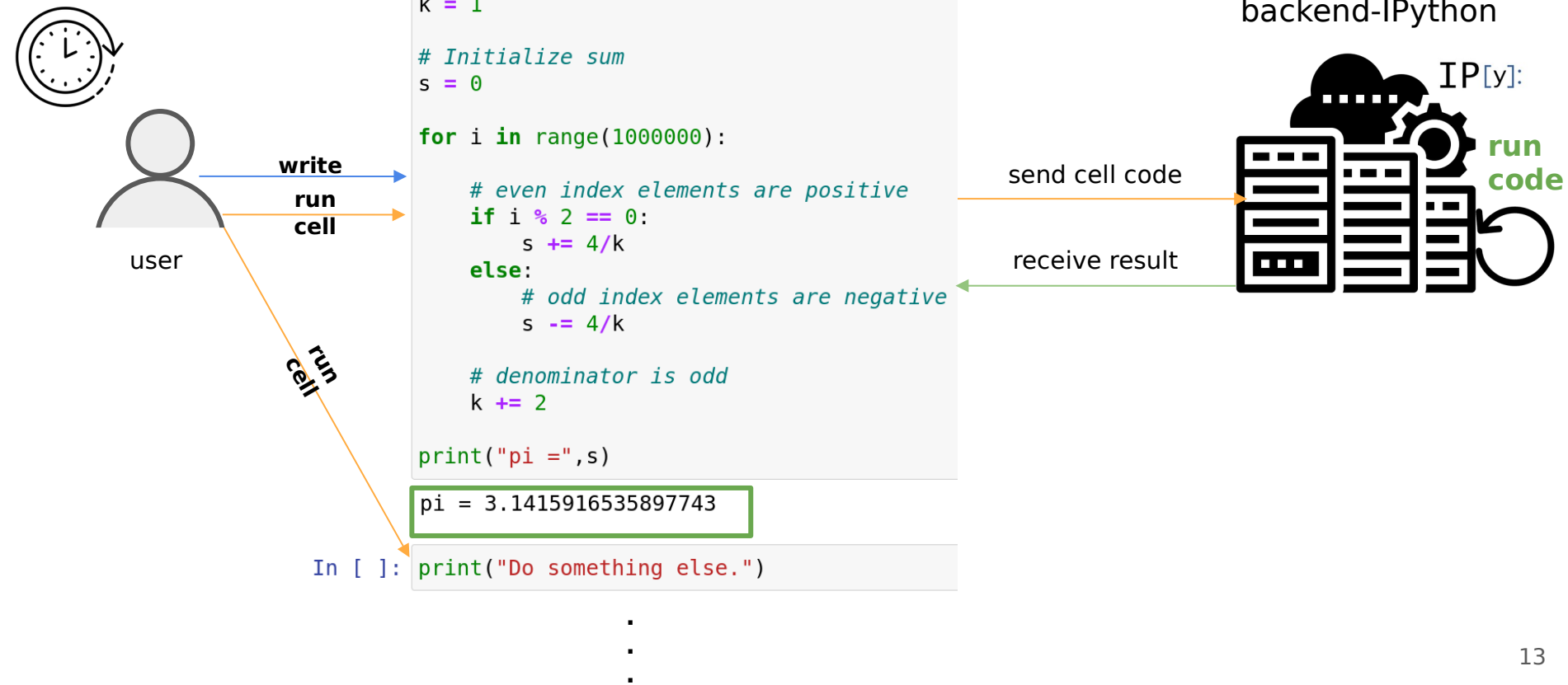


send cell code

receive result



# Interactive environment



# What is Kubeflow?

# What is Kubeflow?



# What is Kubeflow?





# What is Kubeflow?

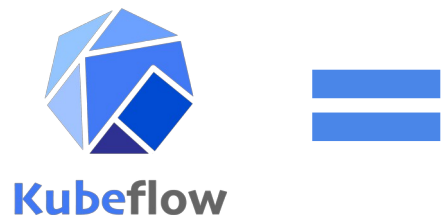


# What is Kubeflow?

web  
app



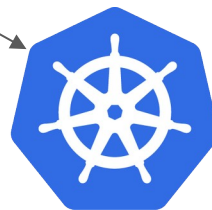
# What is Kubeflow?



web  
app



run as Pod

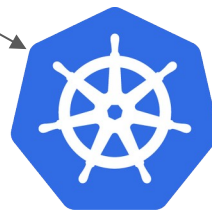


# What is Kubeflow?

web  
app



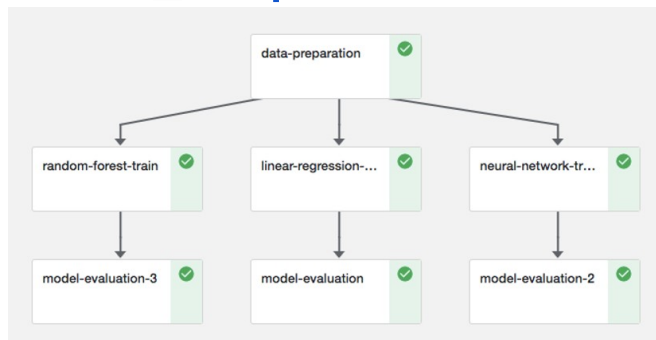
run as Pod



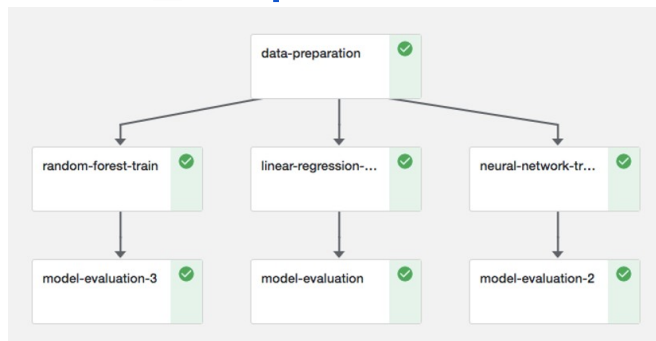
Kubeflow



Pipelines



# What is Kubeflow?

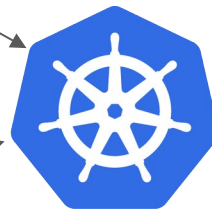


web  
app

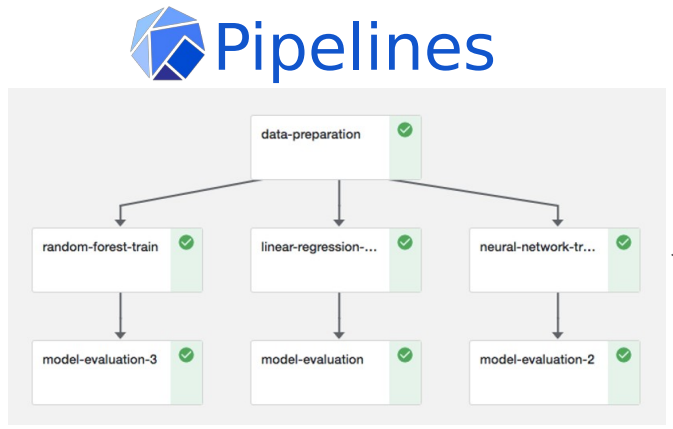
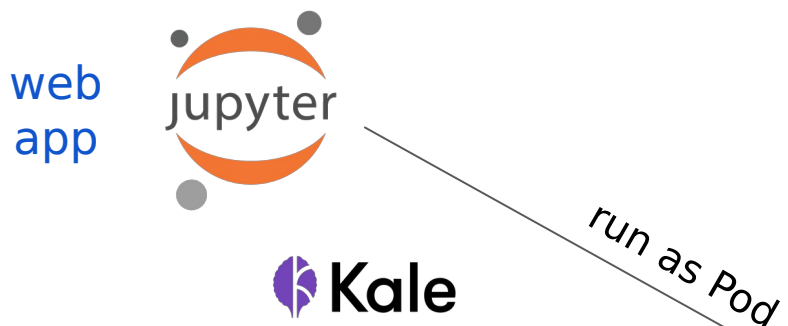
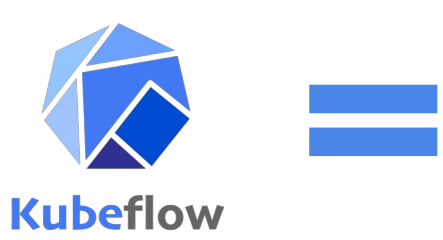


run as Pod

Pod per step

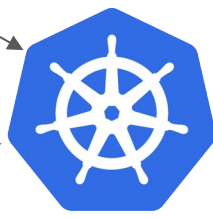


# What is Kubeflow?

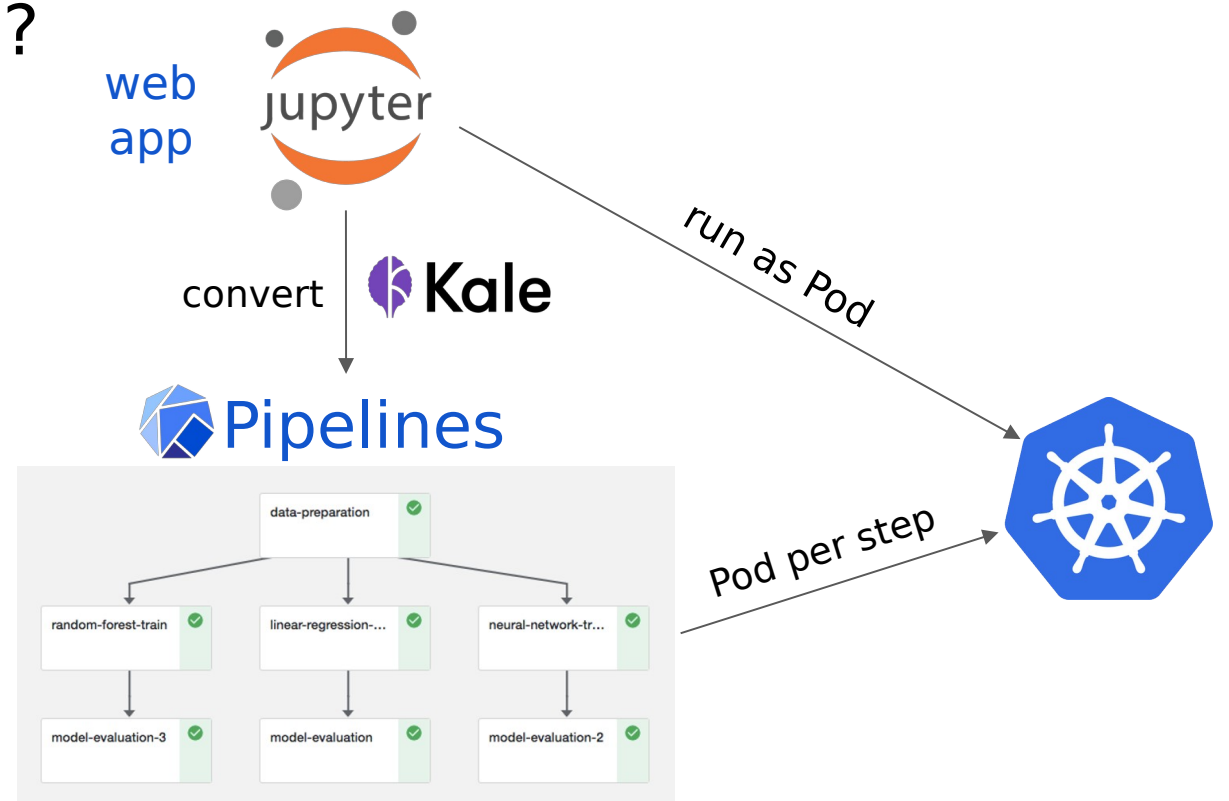
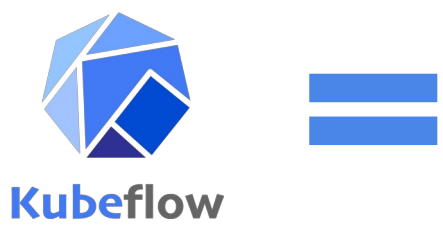


run as Pod

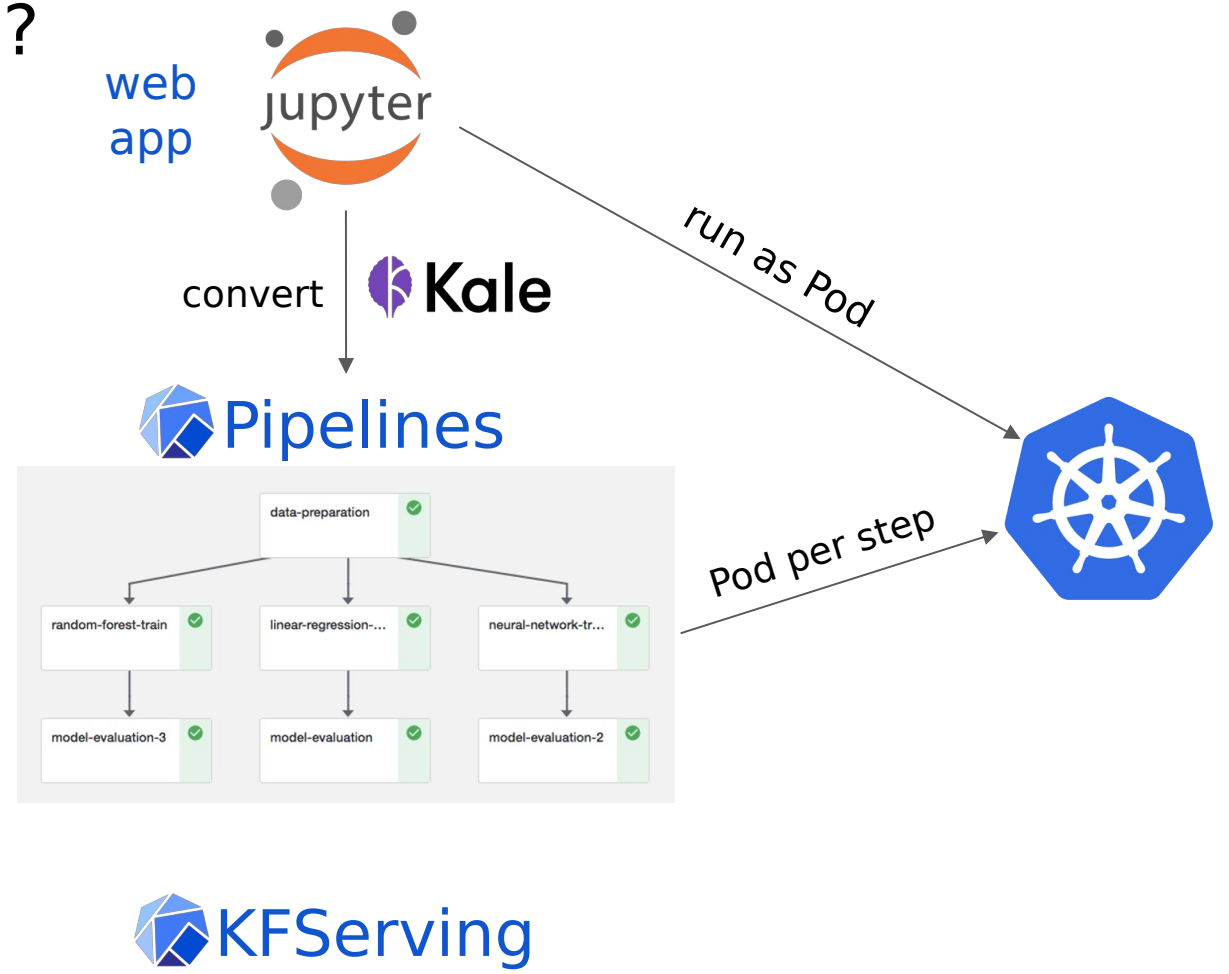
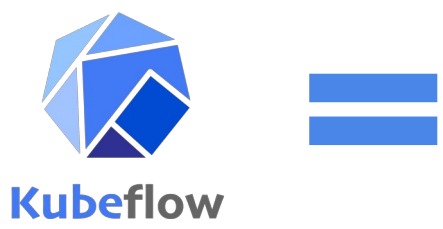
Pod per step



# What is Kubeflow?

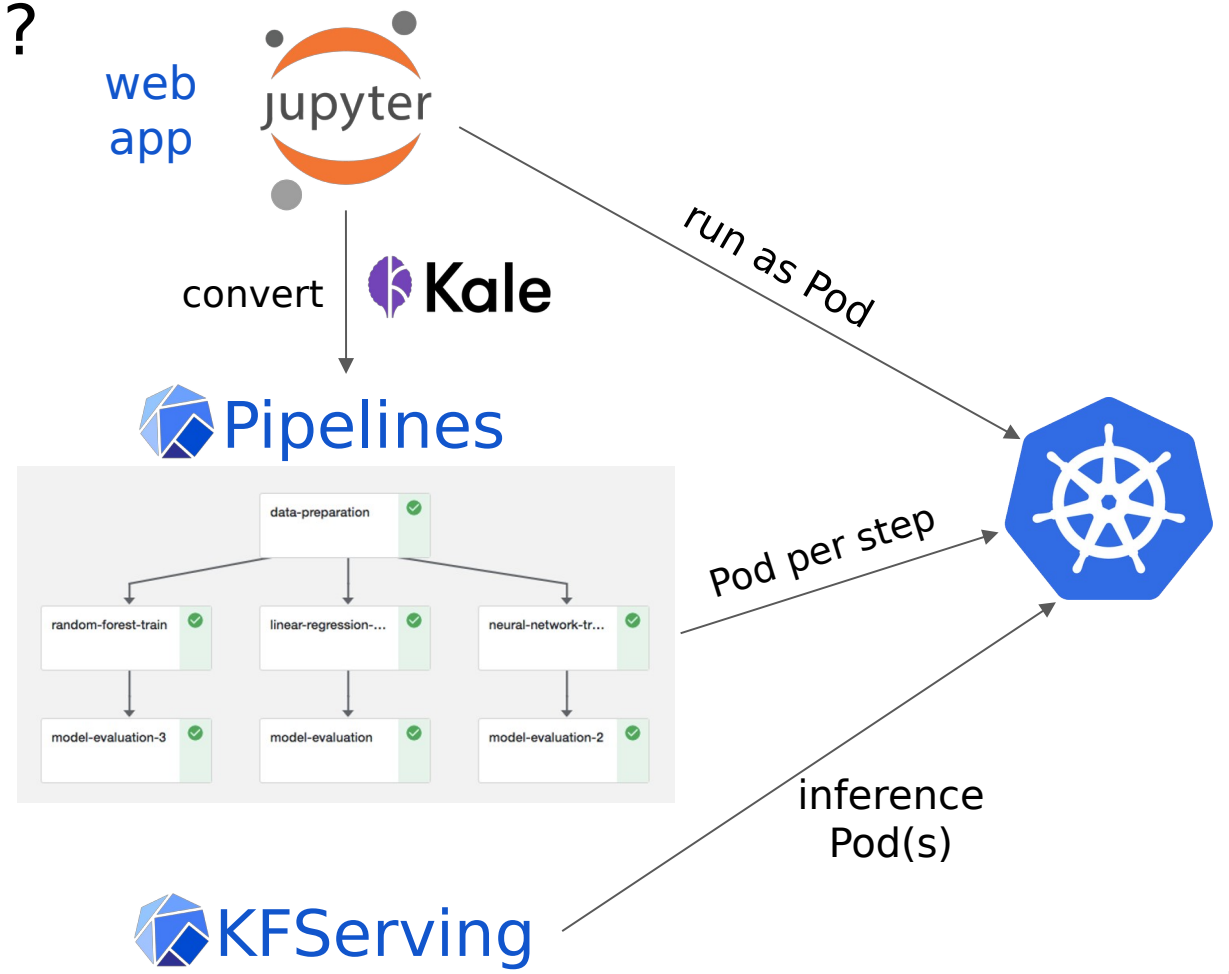
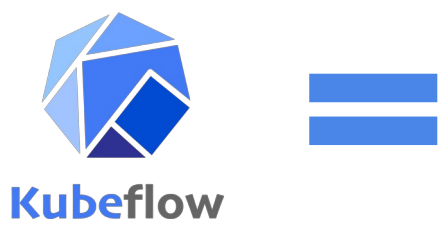


# What is Kubeflow?





# What is Kubeflow?



# The Problem



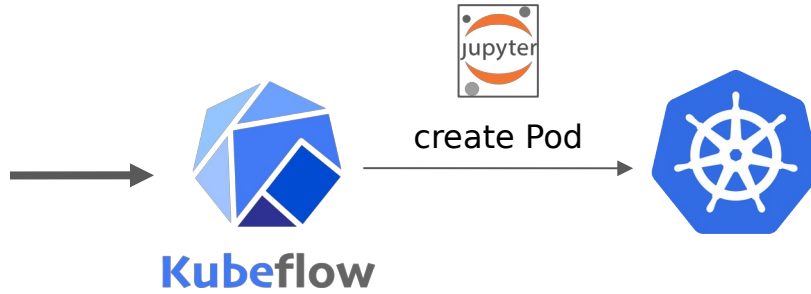
# The Problem



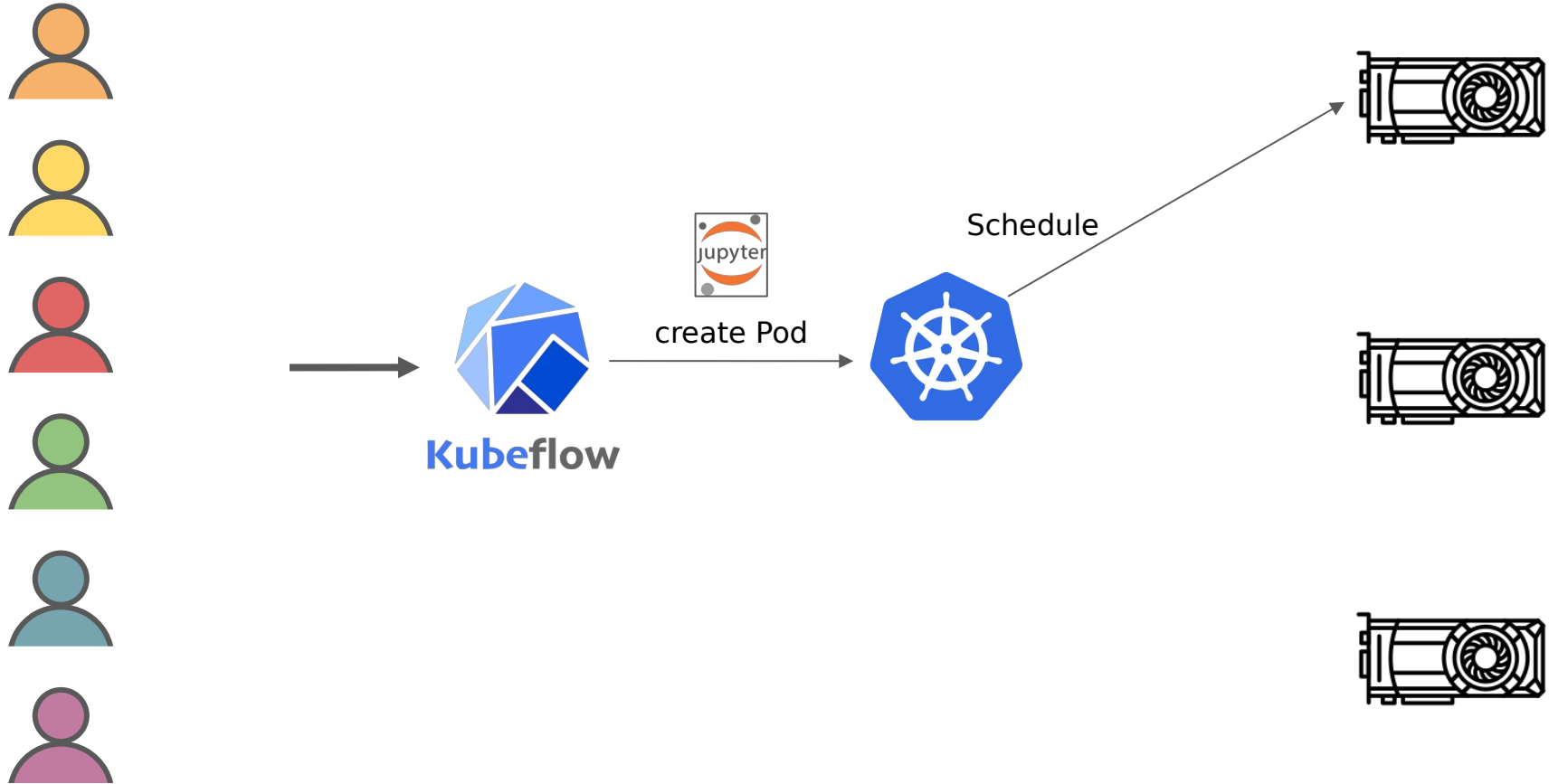
Kubeflow



# The Problem



# The Problem



# The Problem



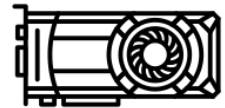
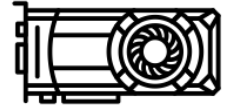
Kubeflow



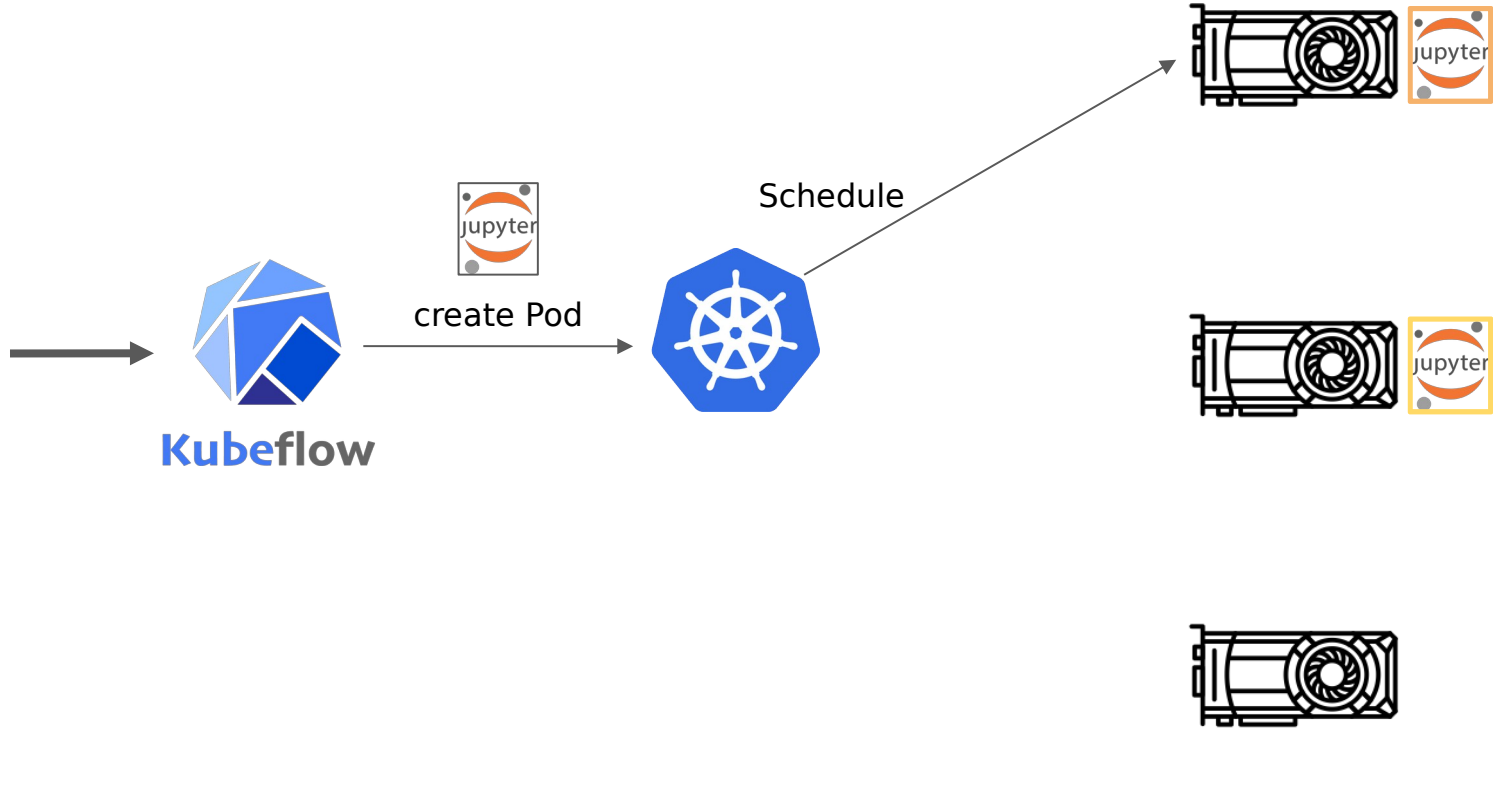
create Pod



Schedule



# The Problem



# The Problem



Kubeflow



create Pod



Schedule





# The Problem



?



?



?



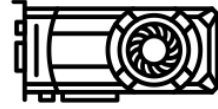
Kubeflow



create Pod



Schedule



# The Problem



Kubeflow



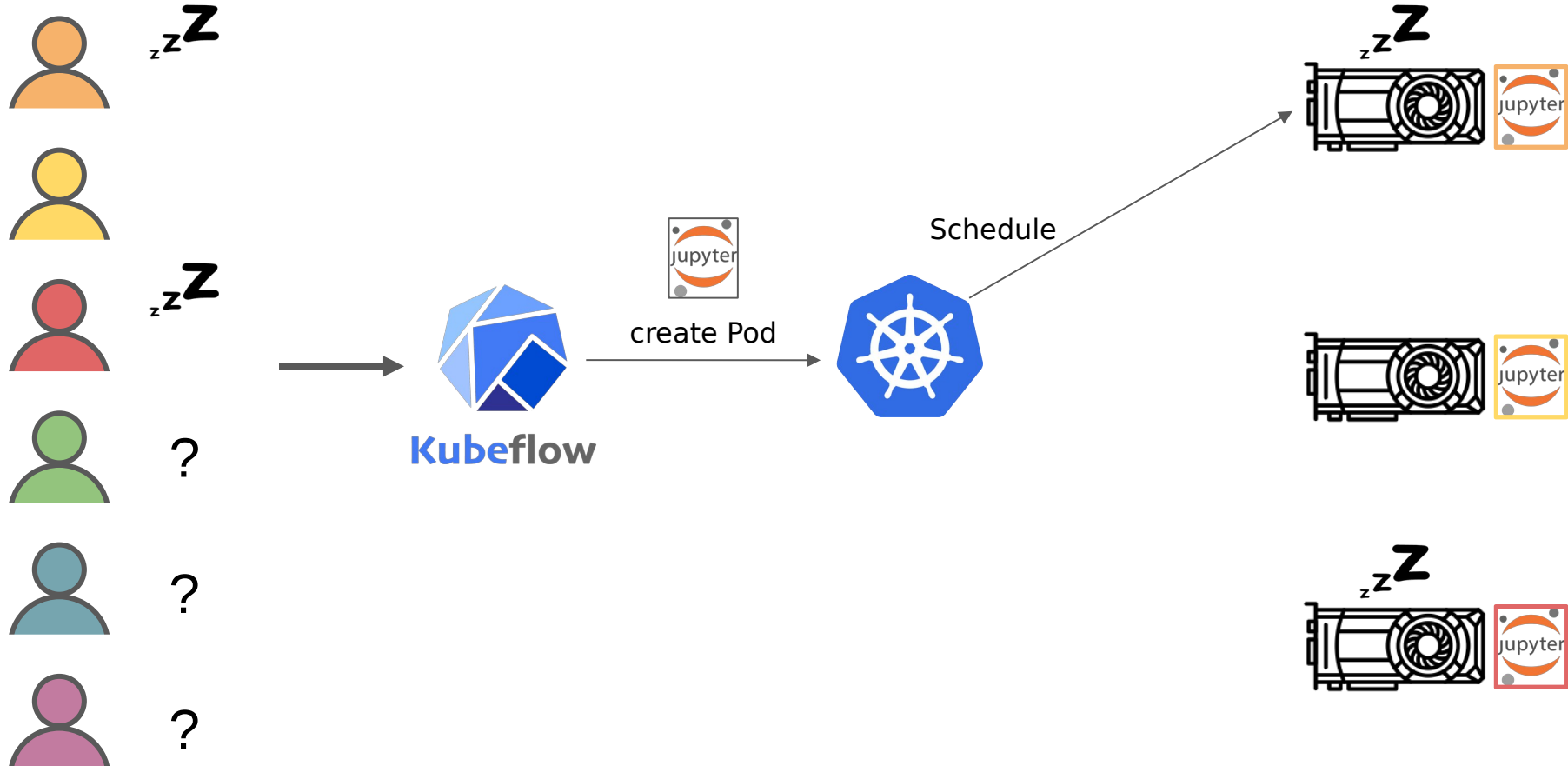
create Pod



Schedule

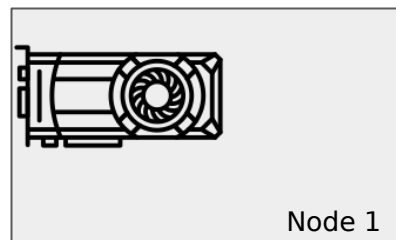


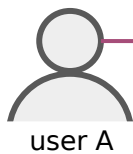
# The Problem





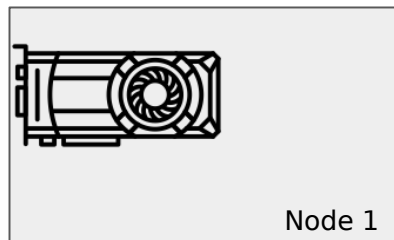
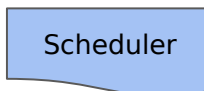
Scheduler





write

```
kind: Pod
metadata:
  name: jupyter-1
spec:
  image: jupyter:latest
  resources:
    limits:
      nvidia.com/gpu: 1
  nodeName:
```

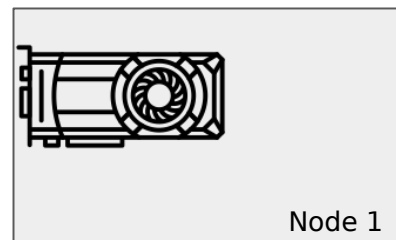
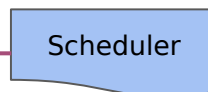




write

```
kind: Pod
metadata:
  name: jupyter-1
spec:
  image: jupyter:latest
  resources:
    limits:
      nvidia.com/gpu: 1
  nodeName:
```

schedule

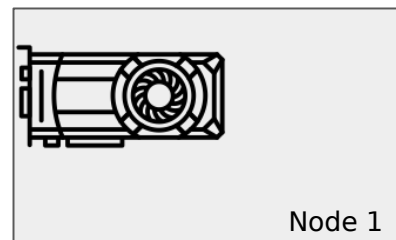
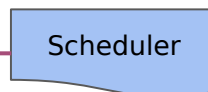




write

```
kind: Pod
metadata:
  name: jupyter-1
spec:
  image: jupyter:latest
  resources:
    limits:
      nvidia.com/gpu: 1
  nodeName: Node 1
```

schedule

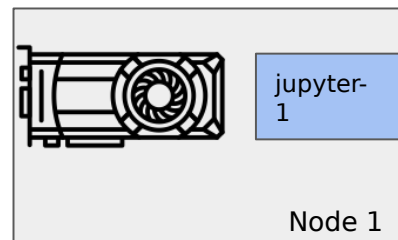
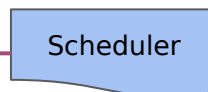




write

```
kind: Pod
metadata:
  name: jupyter-1
spec:
  image: jupyter:latest
  resources:
    limits:
      nvidia.com/gpu: 1
  nodeName: Node 1
```

schedule







doing work...



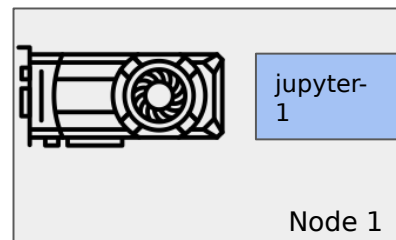
user A

write

```
kind: Pod
metadata:
  name: jupyter-1
spec:
  image: jupyter:latest
  resources:
    limits:
      nvidia.com/gpu: 1
  nodeName: Node 1
```

schedule

Scheduler





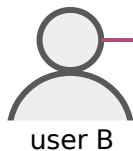
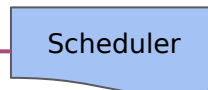
doing work...



write

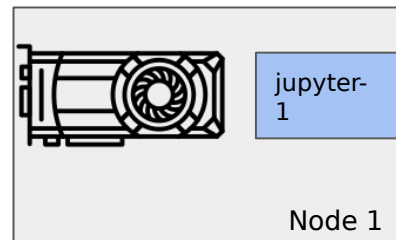
```
kind: Pod
metadata:
  name: jupyter-1
spec:
  image: jupyter:latest
  resources:
    limits:
      nvidia.com/gpu: 1
  nodeName: Node 1
```

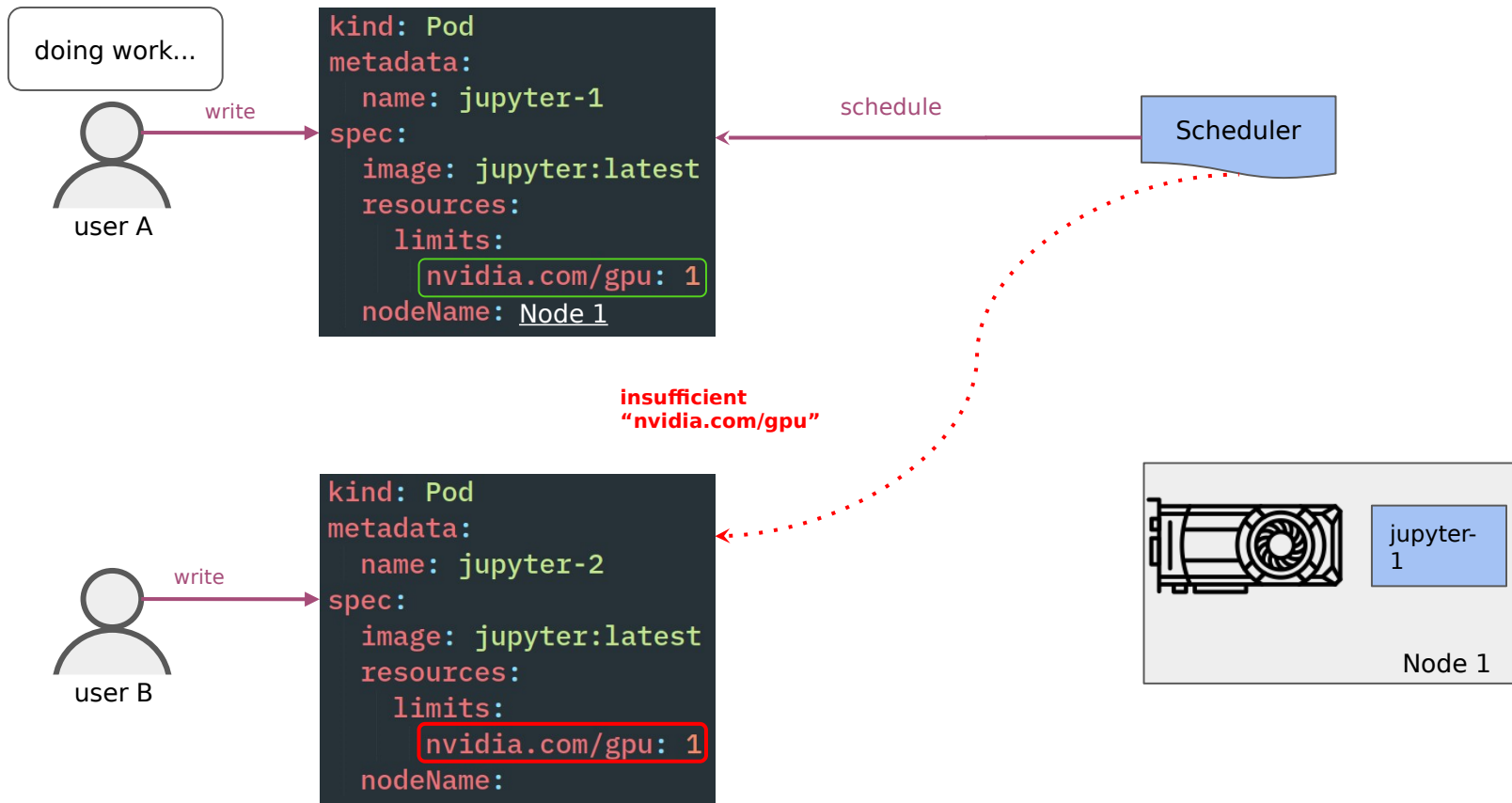
schedule

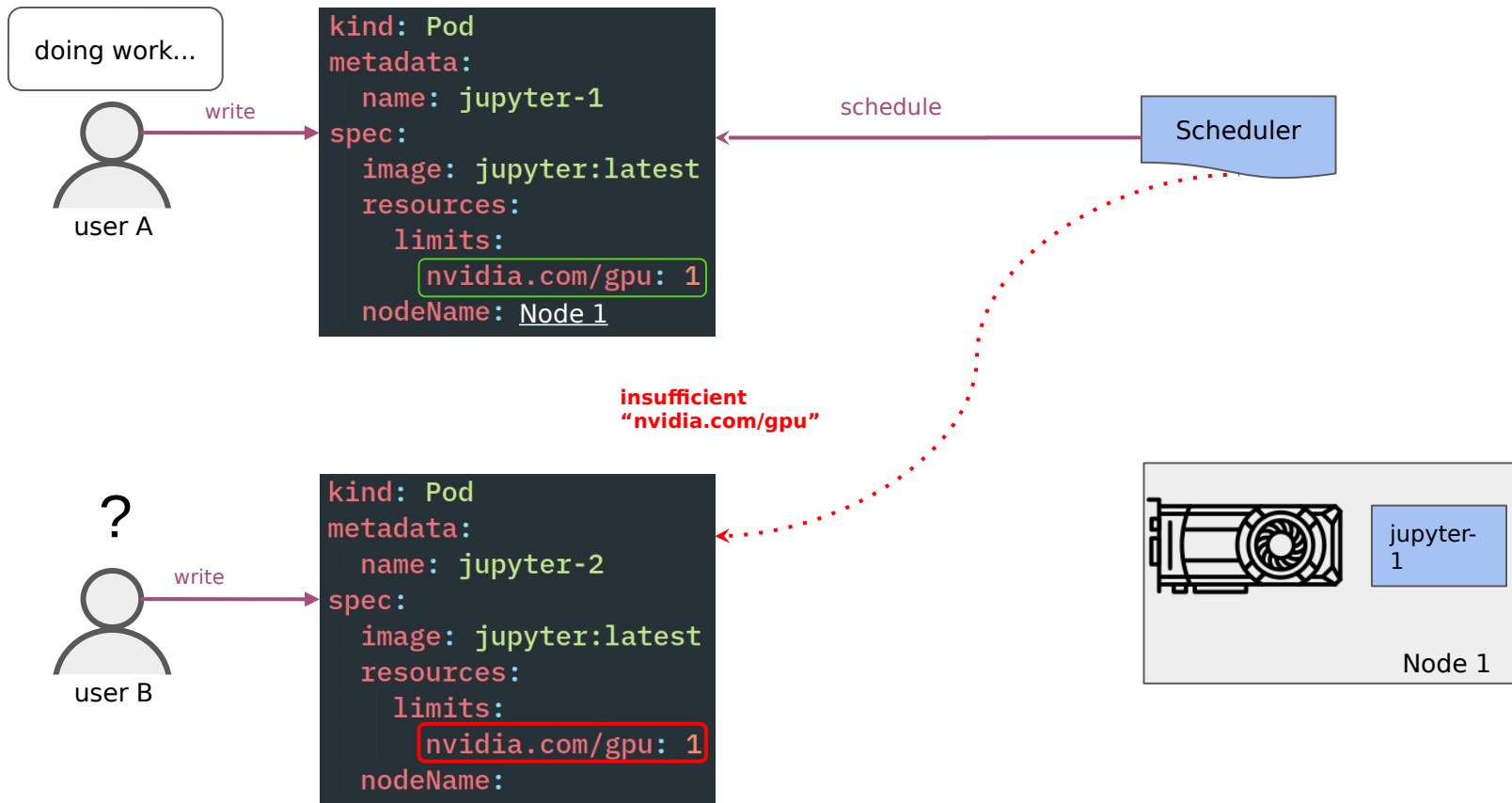


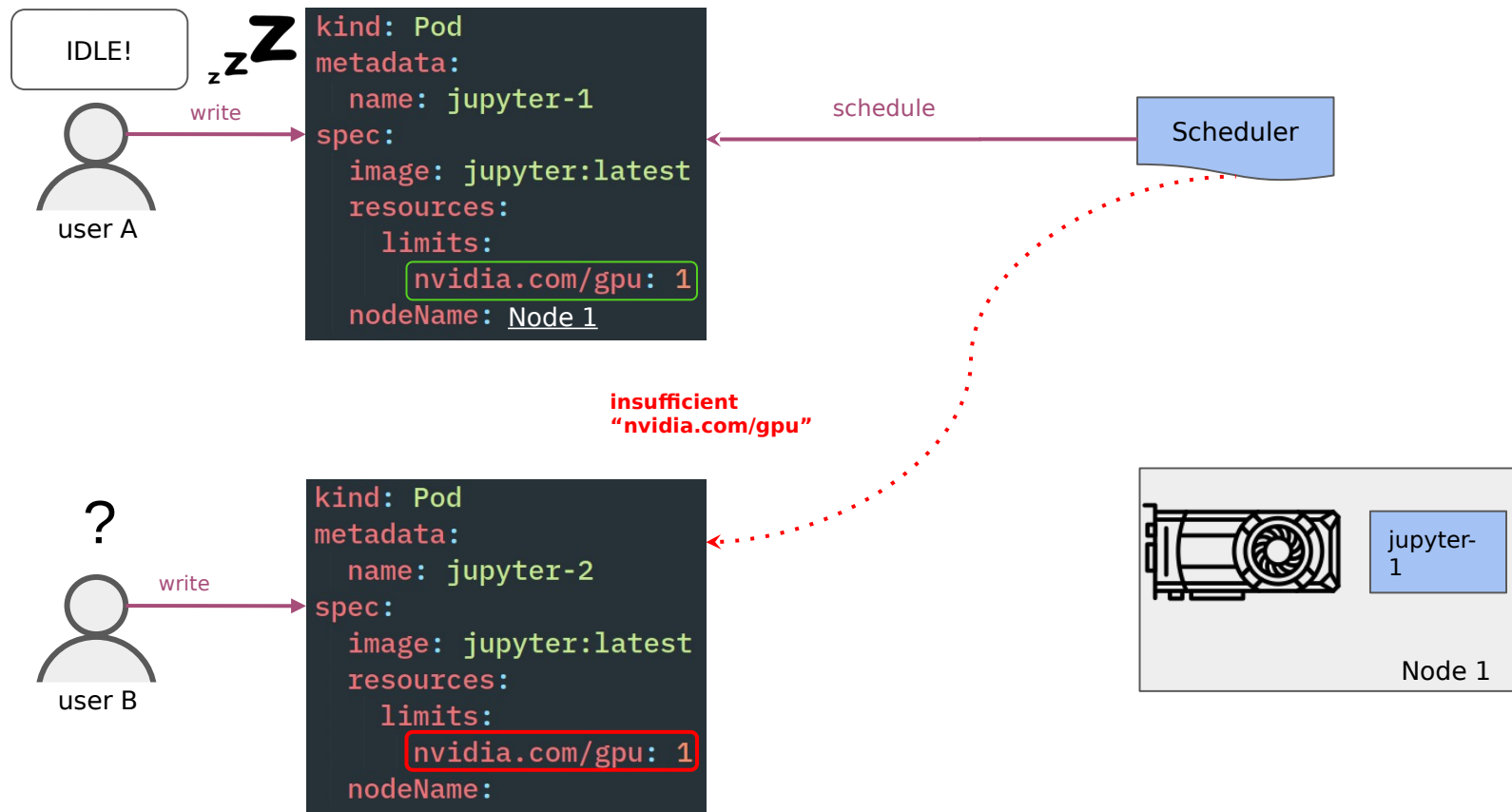
write

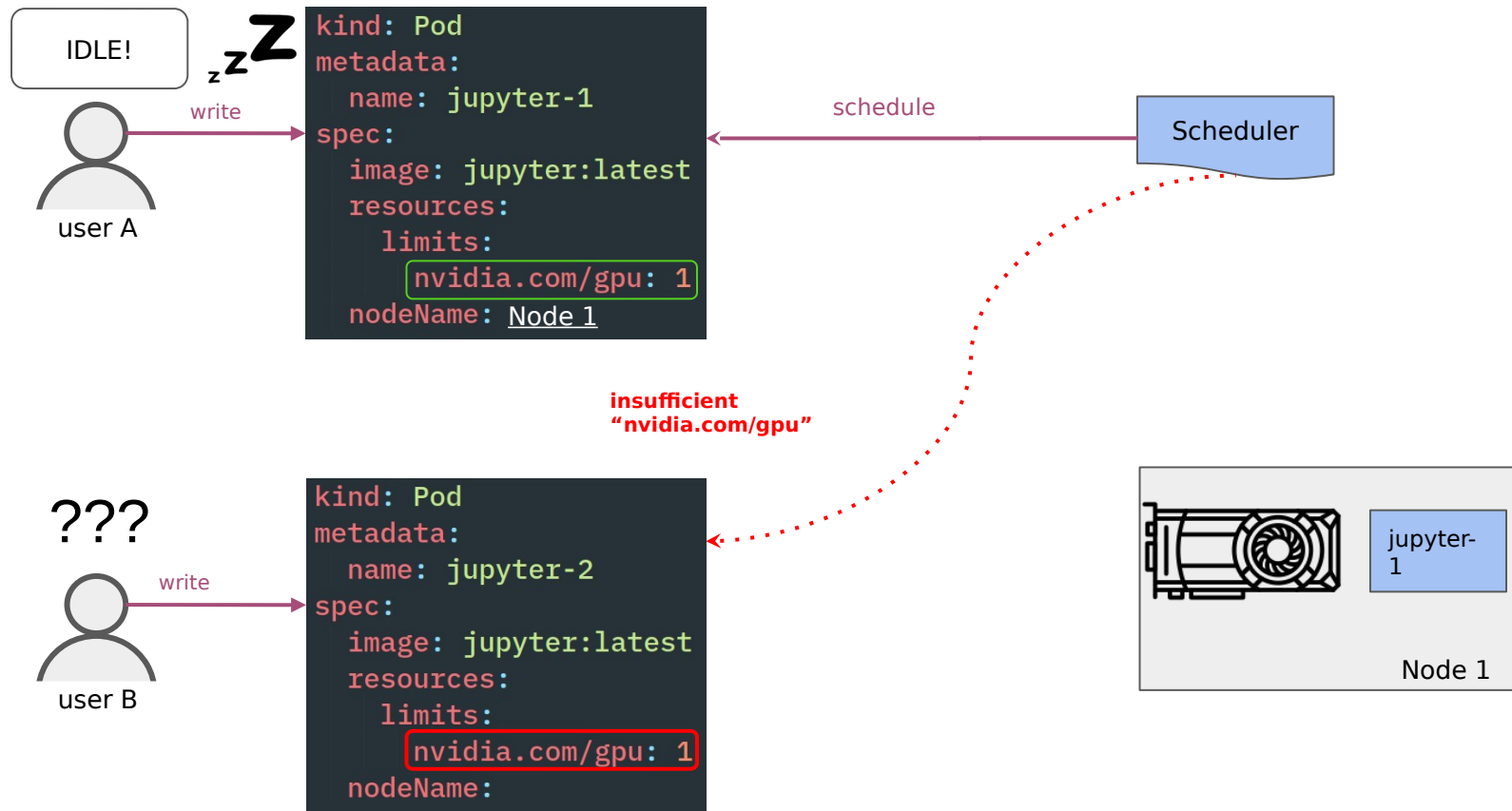
```
kind: Pod
metadata:
  name: jupyter-2
spec:
  image: jupyter:latest
  resources:
    limits:
      nvidia.com/gpu: 1
  nodeName:
```











# Existing Approaches (Aliyun, Kubeshare)

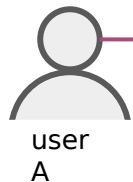
Scheduler

8 GB VRAM



Node 1

# Existing Approaches (Aliyun, Kubeshare)

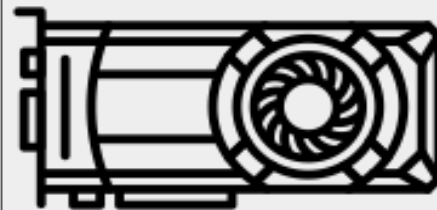


write

```
kind: Pod
metadata:
  name: jupyter-1
spec:
  image: jupyter:latest
  resources:
    limits:
      aliyun.com/gpu-mem: 4
  nodeName:
```

Scheduler

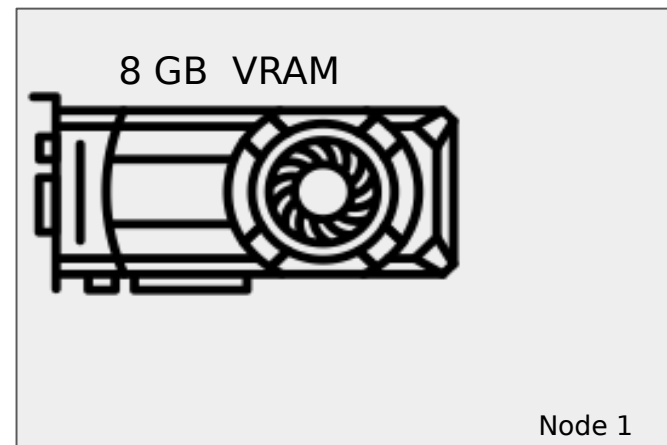
8 GB VRAM



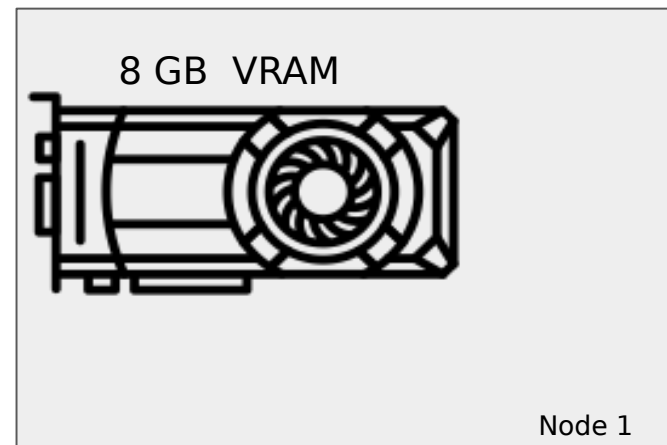
Node 1



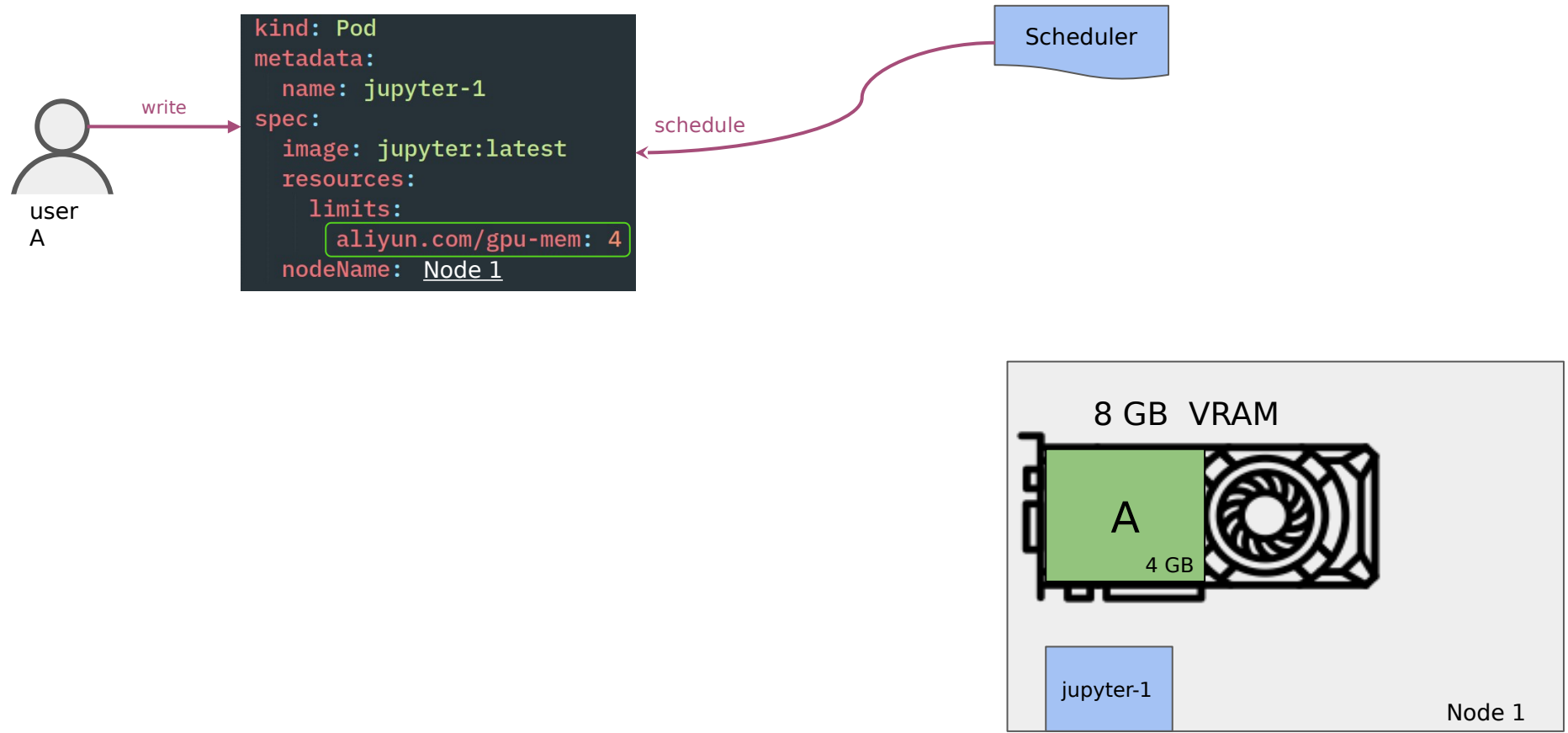
# Existing Approaches (Aliyun, Kubeshare)



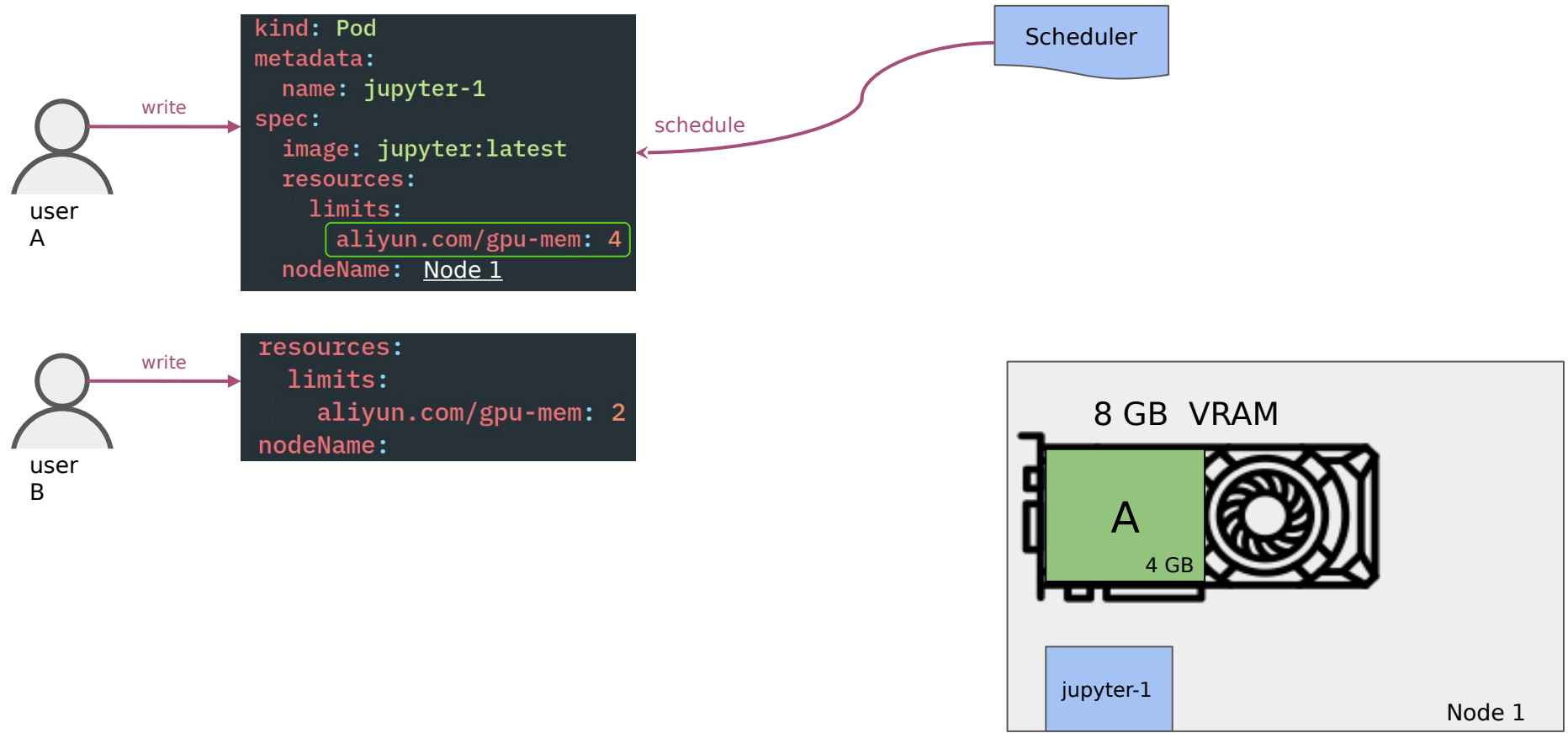
# Existing Approaches (Aliyun, Kubeshare)



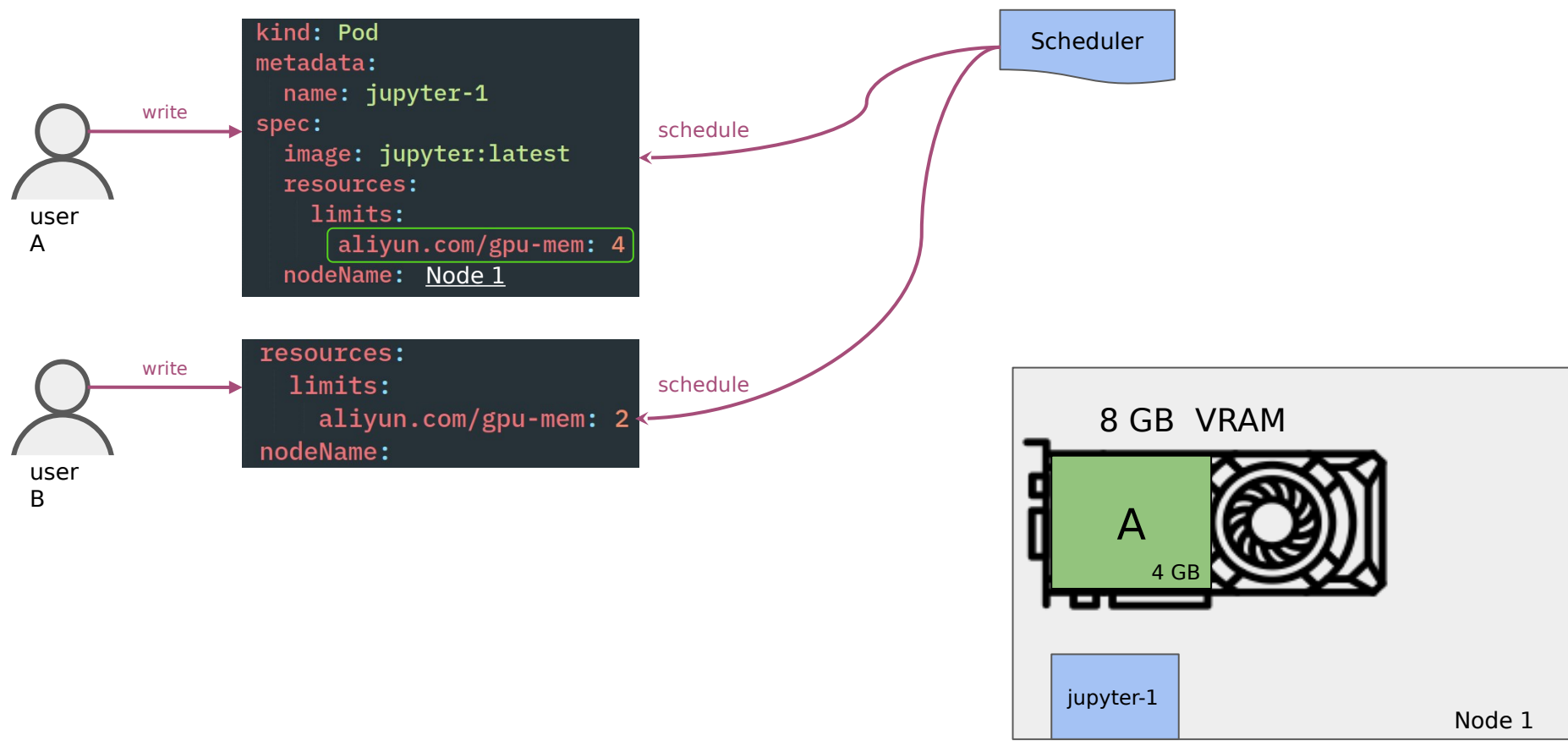
# Existing Approaches (Aliyun, Kubeshare)



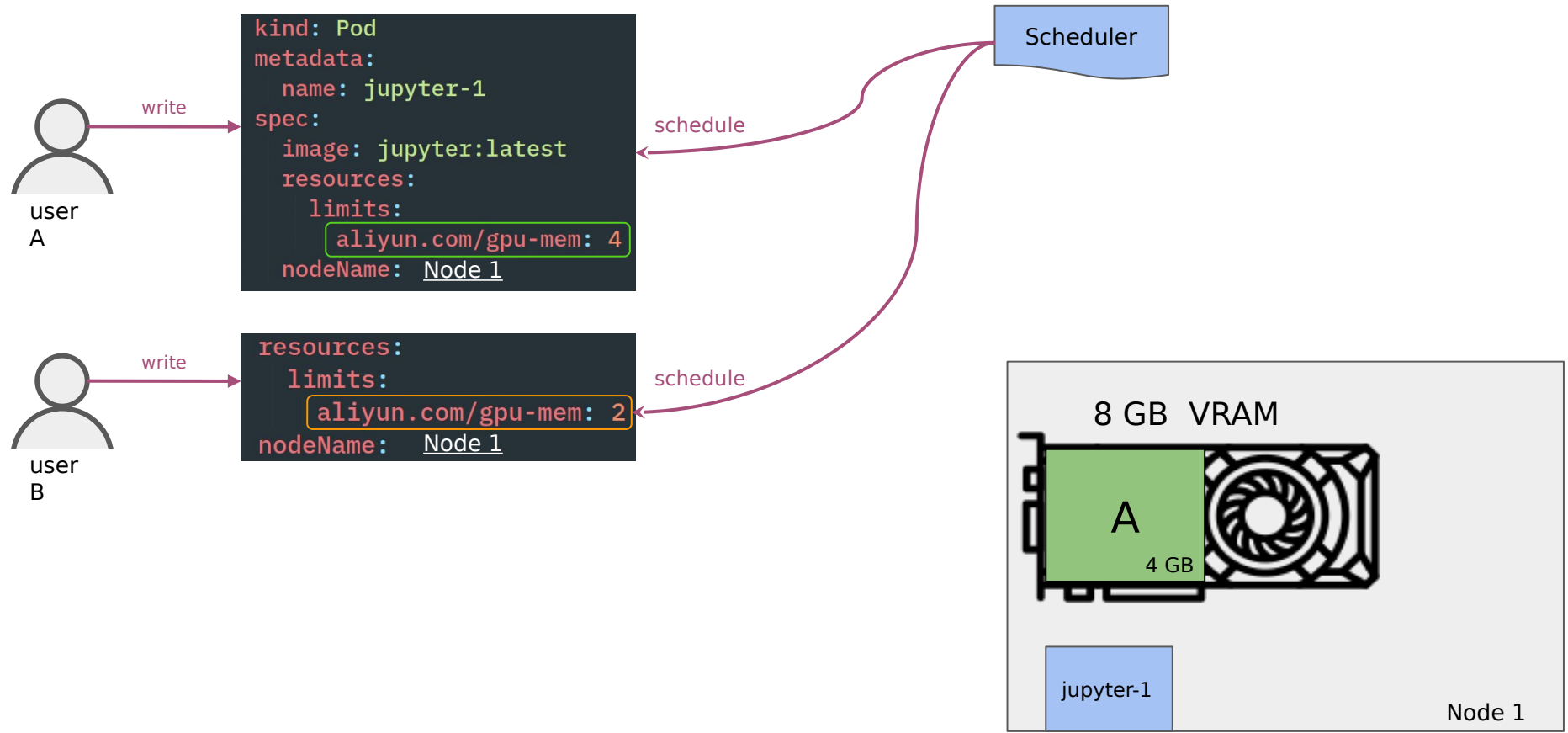
# Existing Approaches (Aliyun, Kubeshare)



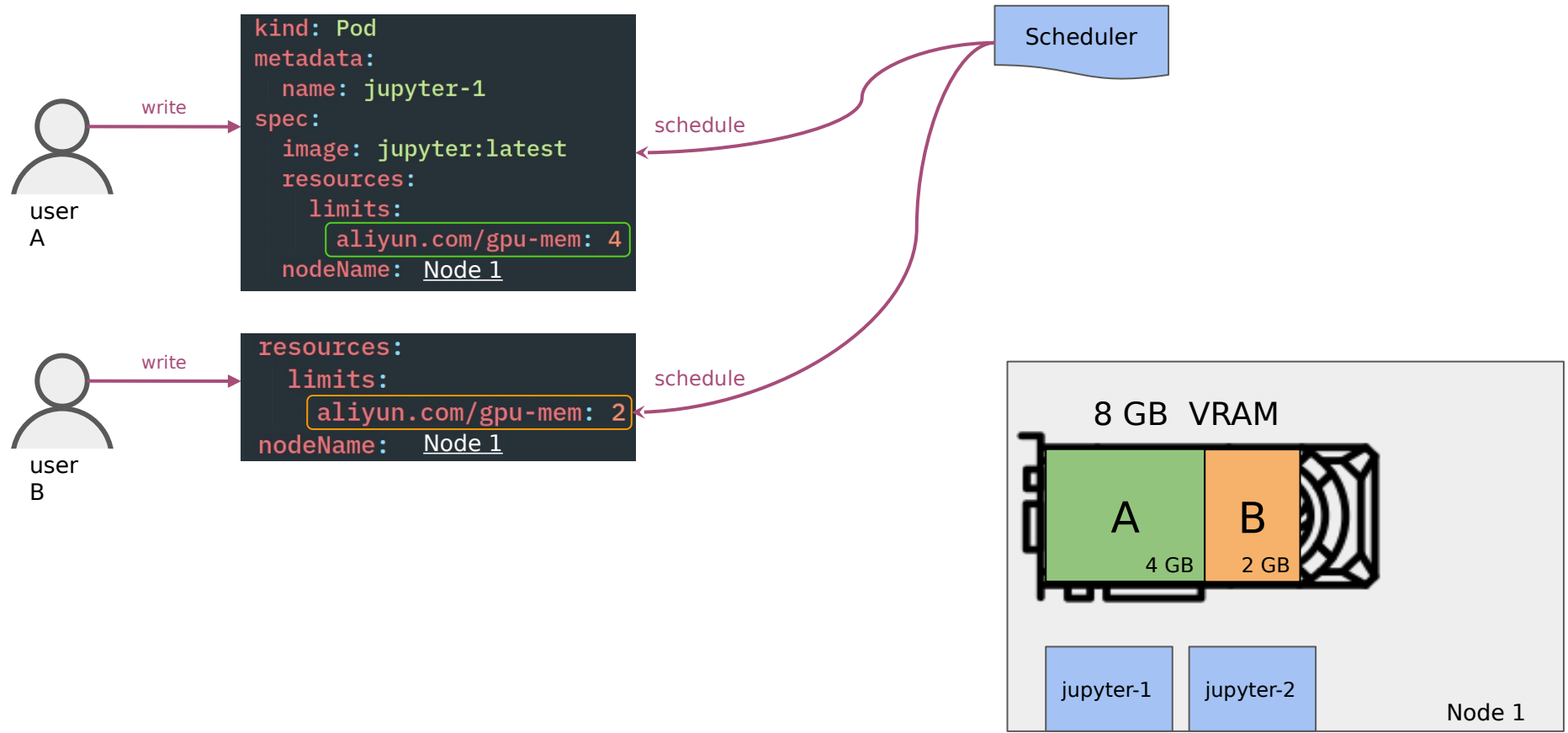
# Existing Approaches (Aliyun, Kubeshare)



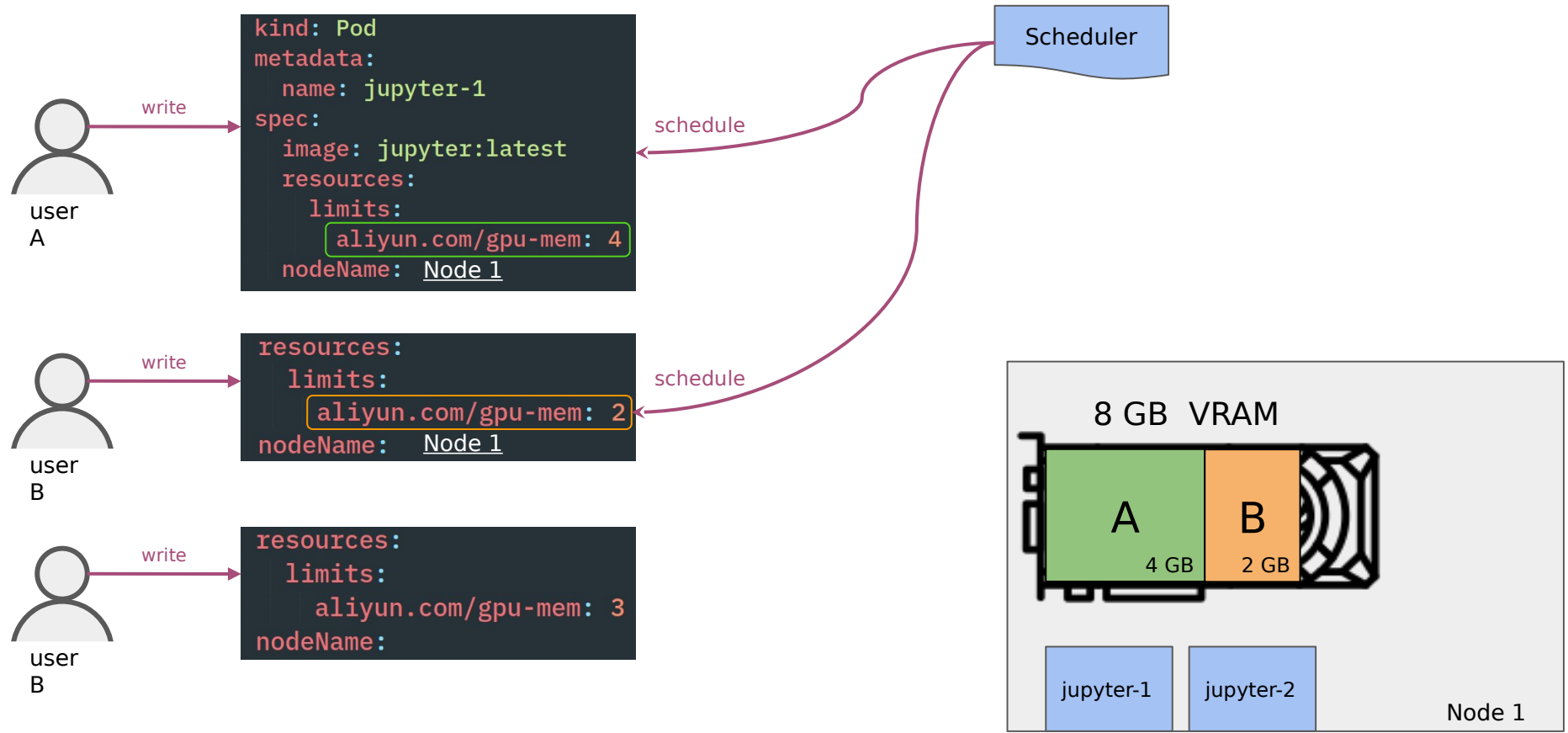
# Existing Approaches (Aliyun, Kubeshare)



# Existing Approaches (Aliyun, Kubeshare)

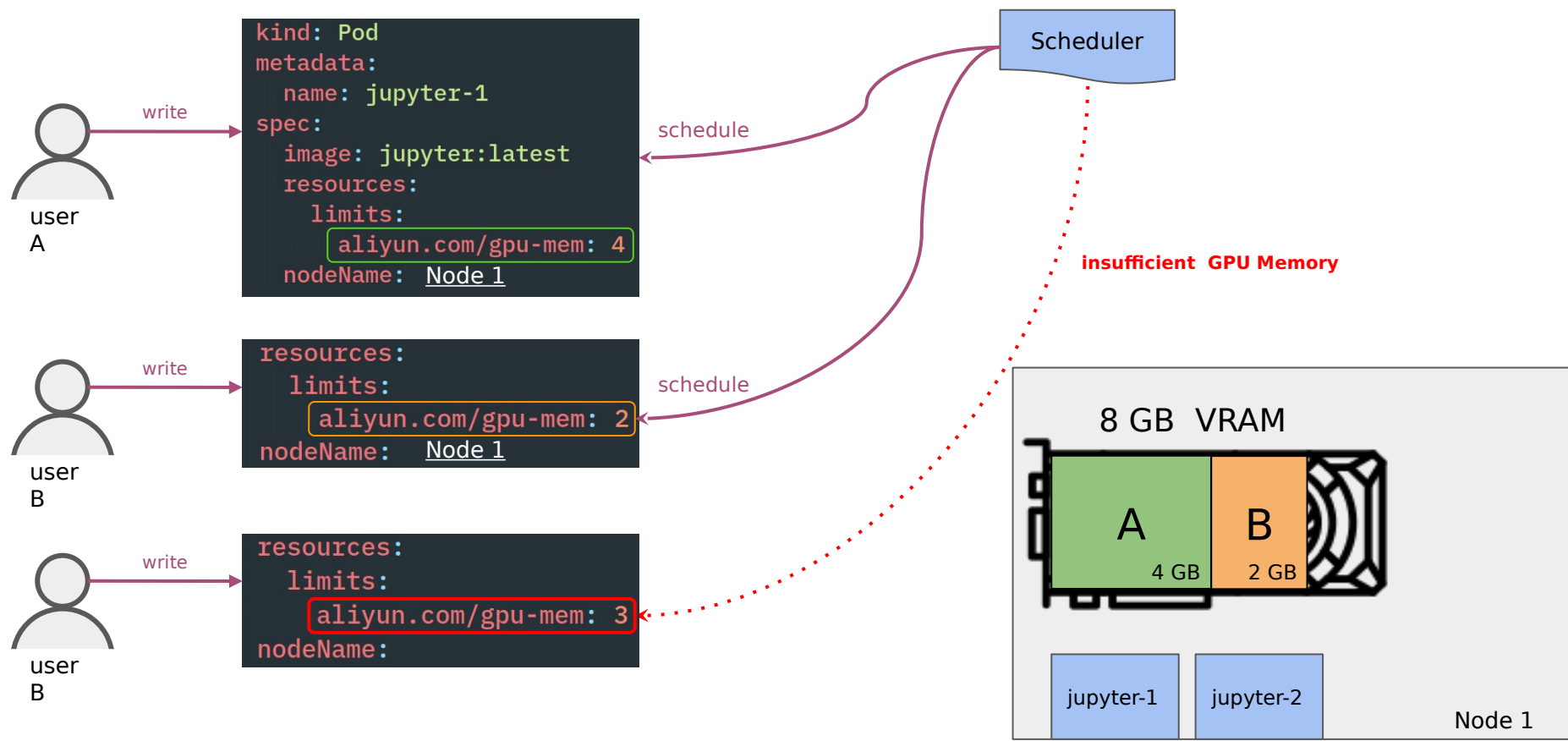


# Existing Approaches (Aliyun, Kubeshare)

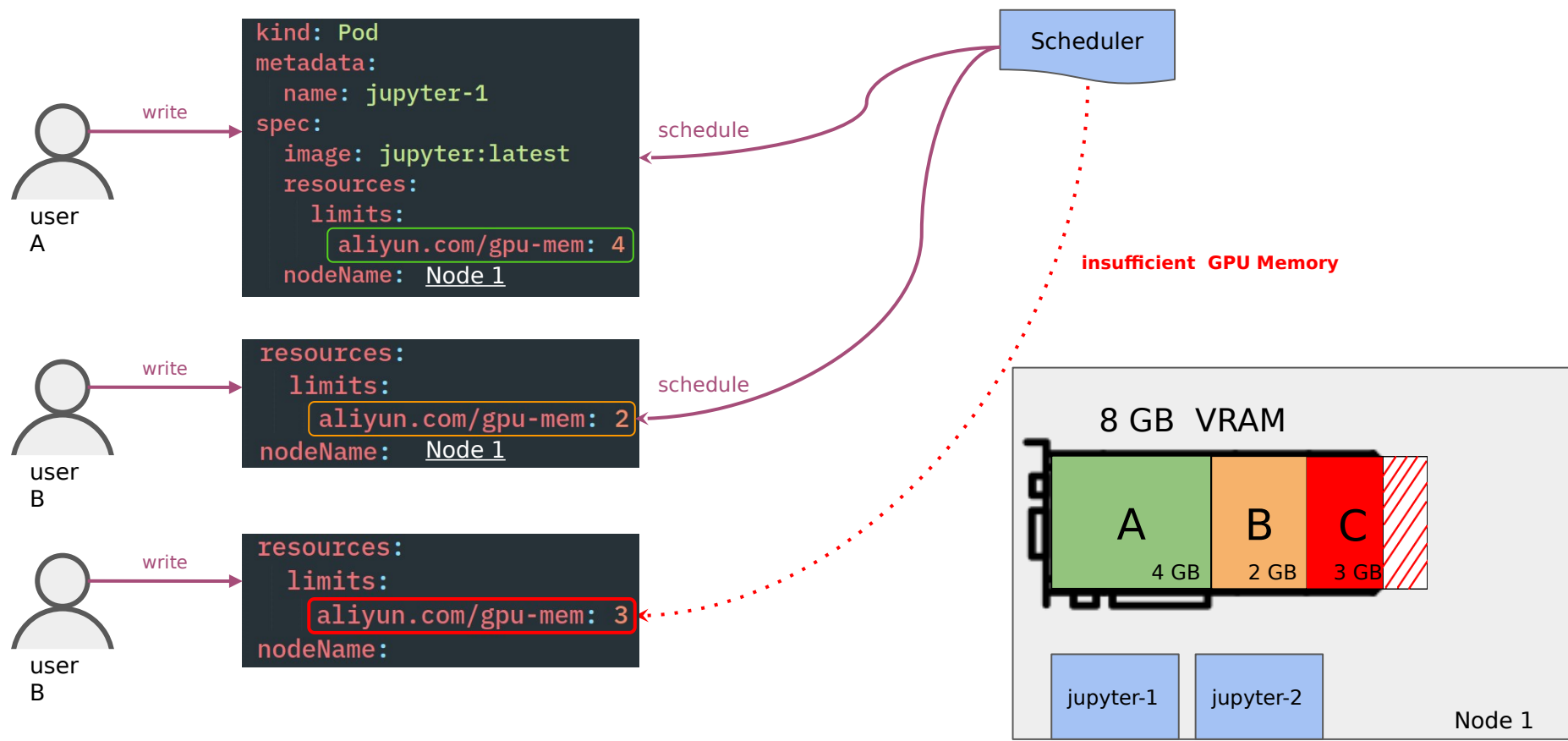




# Existing Approaches (Aliyun, Kubeshare)



# Existing Approaches (Aliyun, Kubeshare)



# Key Idea

# Key Idea



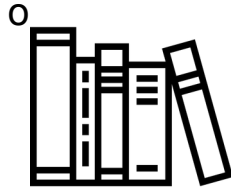
user

# Key Idea



user

libnvshare.s

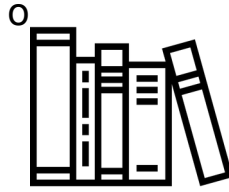


# Key Idea

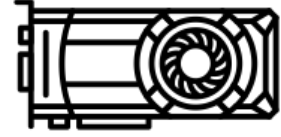


user

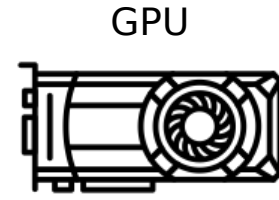
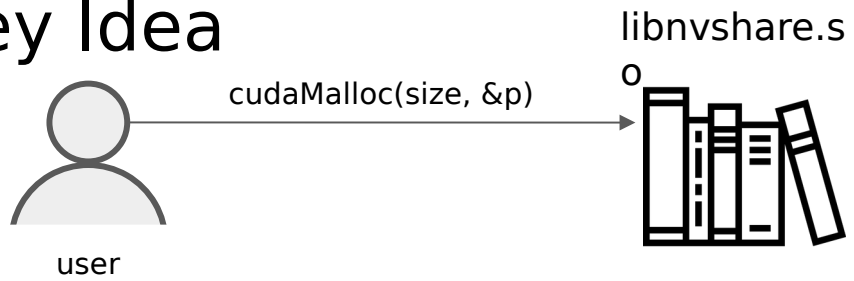
libnvshare.s



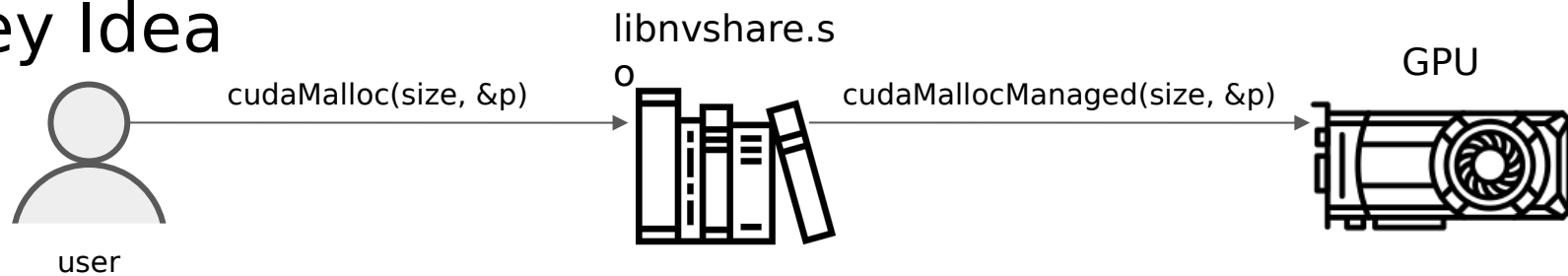
GPU



# Key Idea

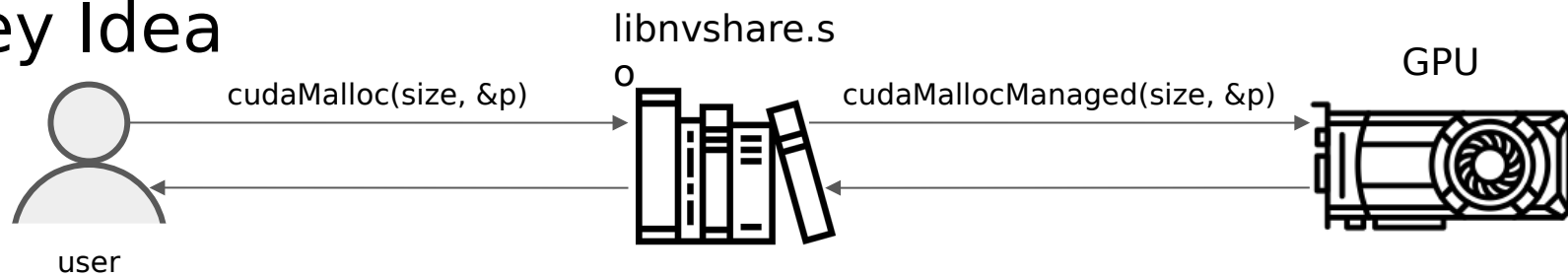


# Key Idea

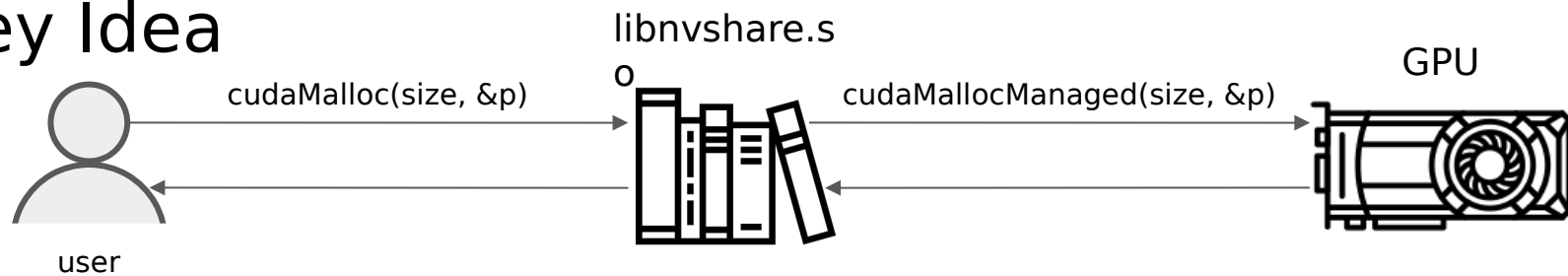




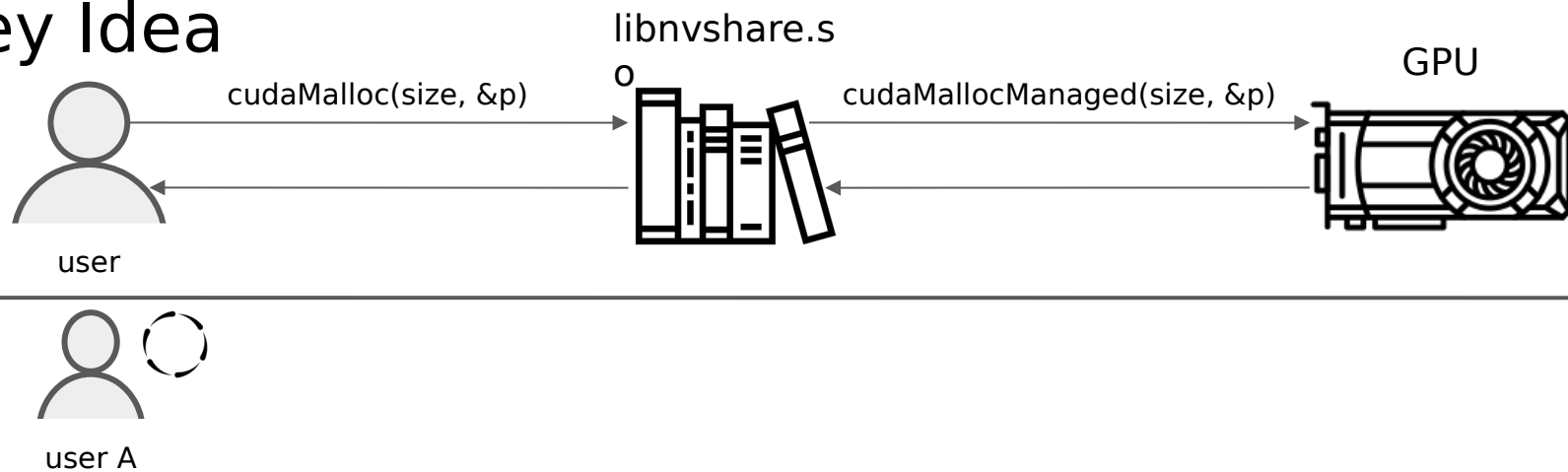
# Key Idea



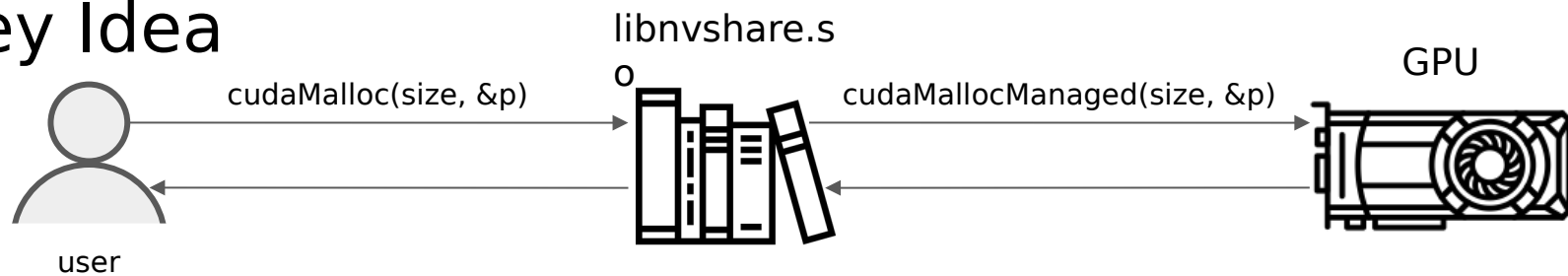
# Key Idea



# Key Idea



# Key Idea



user A

GPU

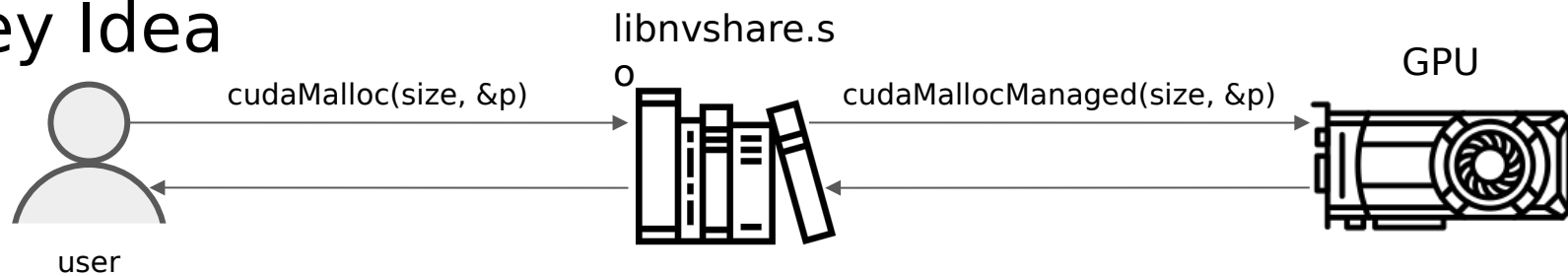
A0

A1

A2

A3

# Key Idea



user A

GPU

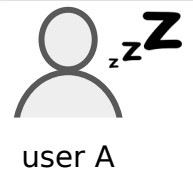
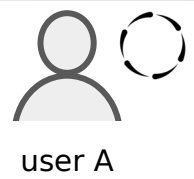
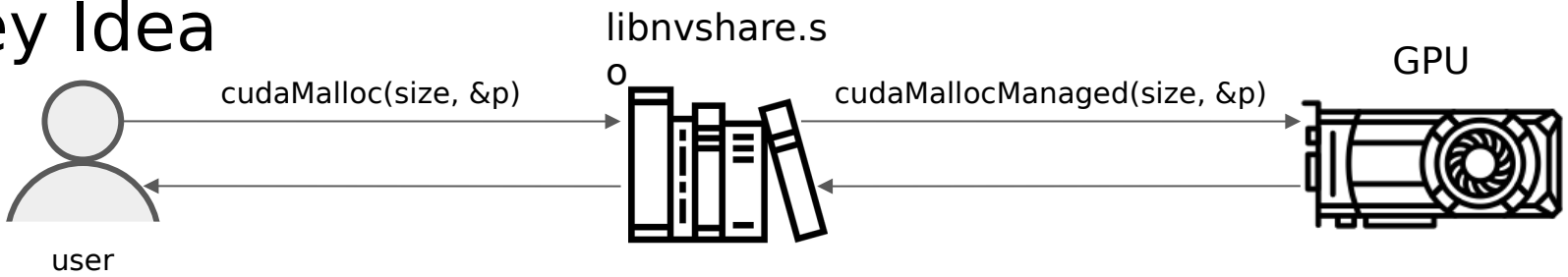
A0

A1

A2

A3

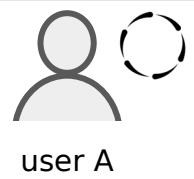
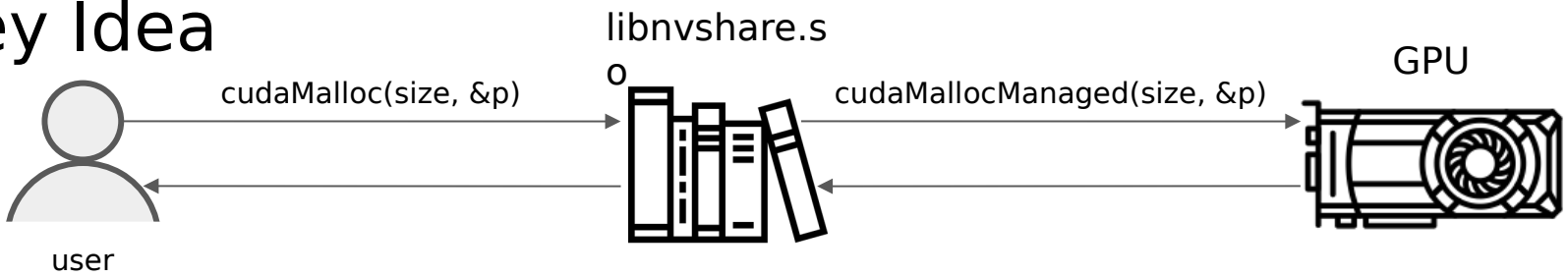
# Key Idea



GPU

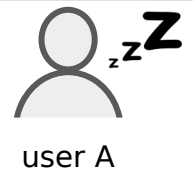
A0
A1
A2
A3

# Key Idea



GPU

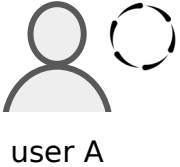
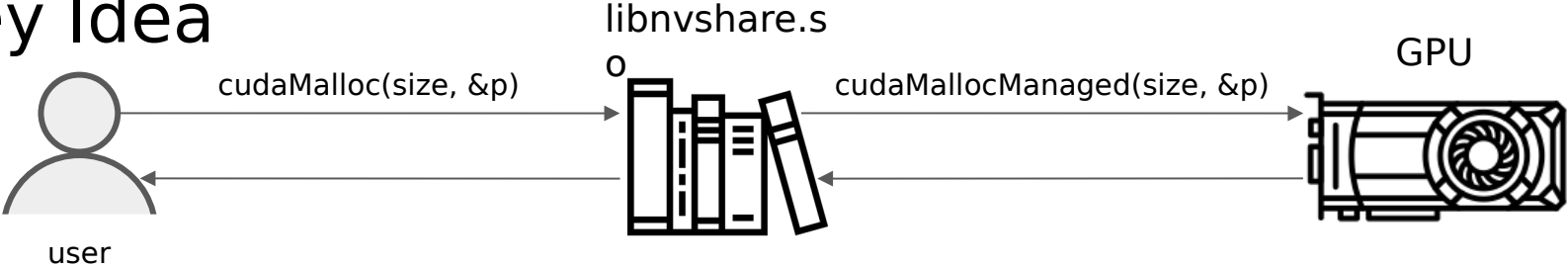
A0
A1
A2
A3



GPU

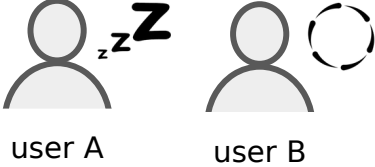
A0
A1
A2
A3

# Key Idea



GPU

A0
A1
A2
A3

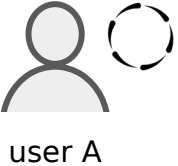
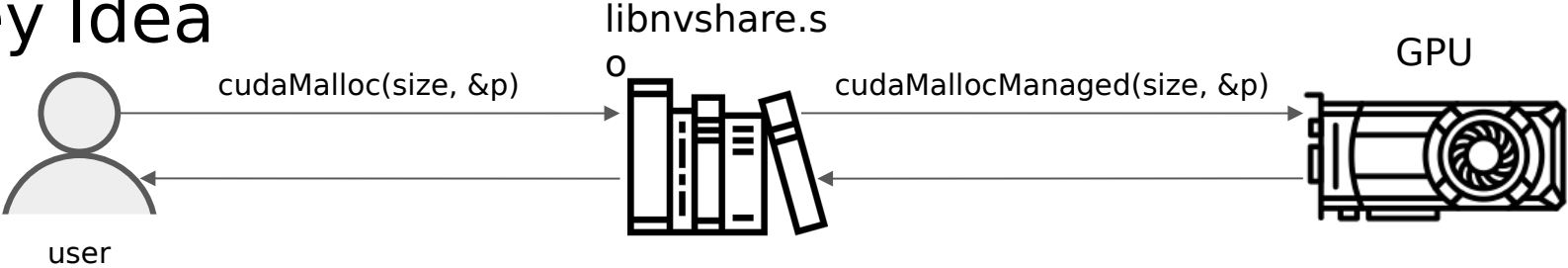


GPU

A0
A1
A2
A3

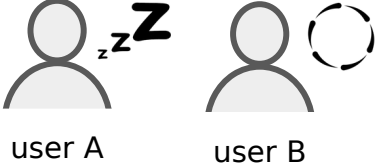


# Key Idea



GPU

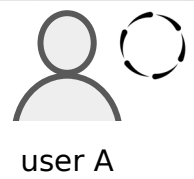
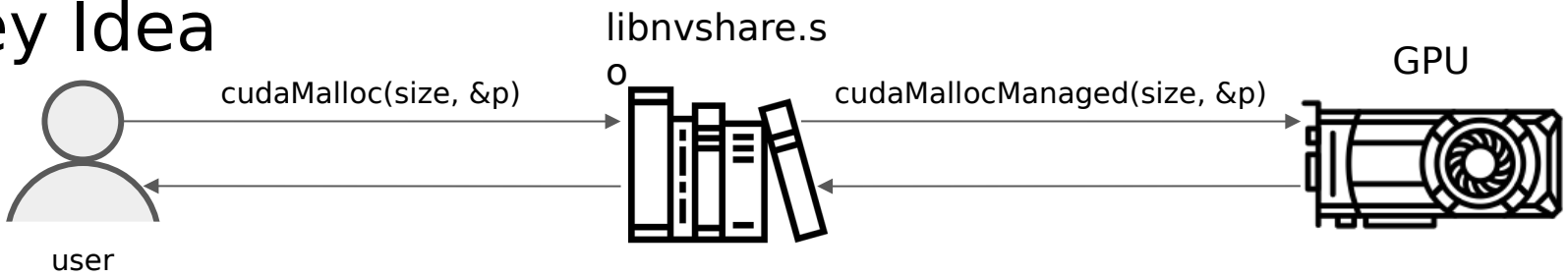
A0
A1
A2
A3



GPU

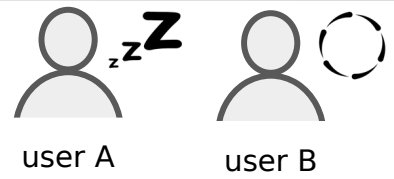
B0
A1
A2
A3

# Key Idea



GPU

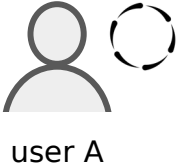
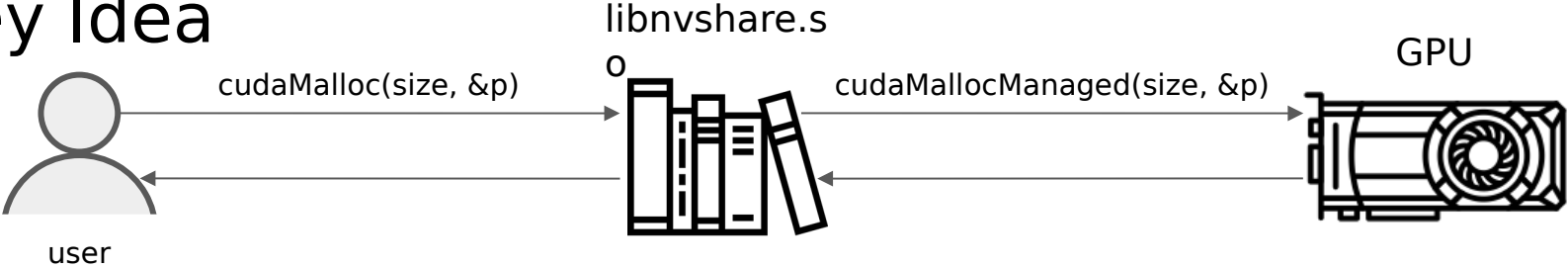
A0
A1
A2
A3



GPU

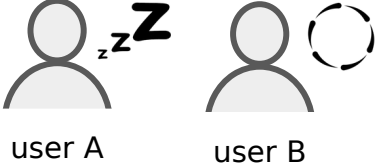
B0
B1
A2
A3

# Key Idea



GPU

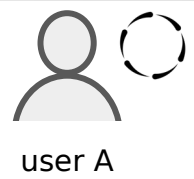
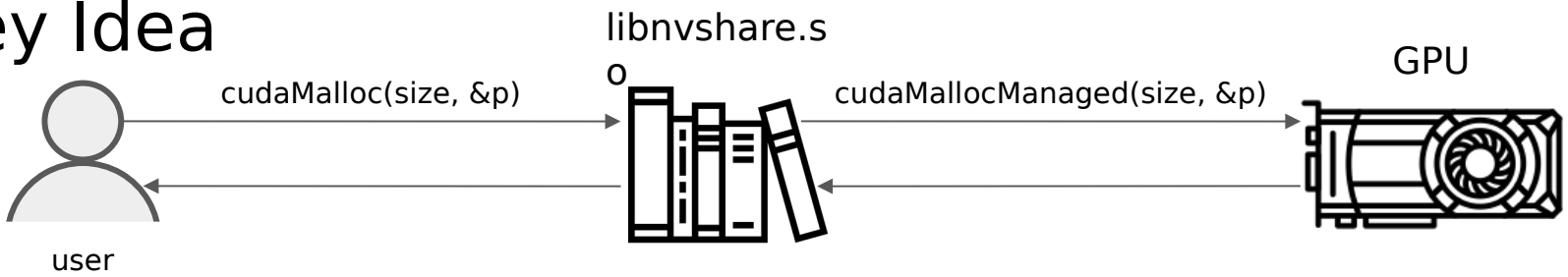
A0
A1
A2
A3



GPU

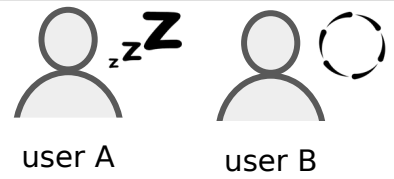
B0
B1
B2
A3

# Key Idea



GPU

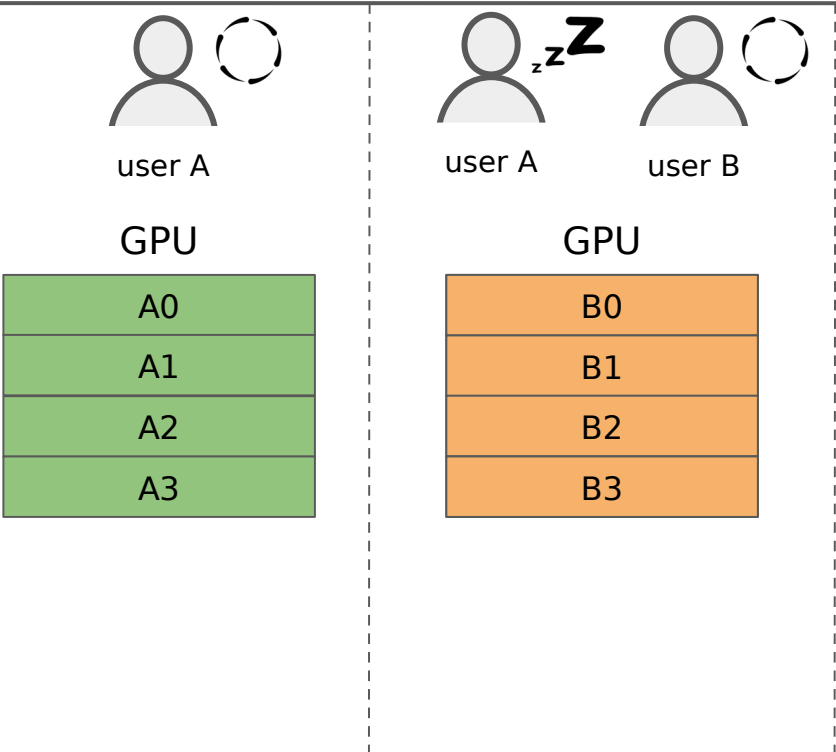
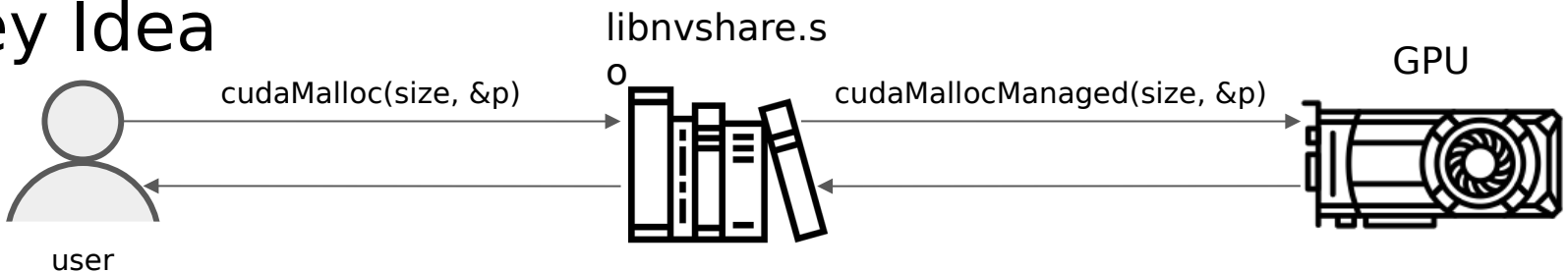
A0
A1
A2
A3



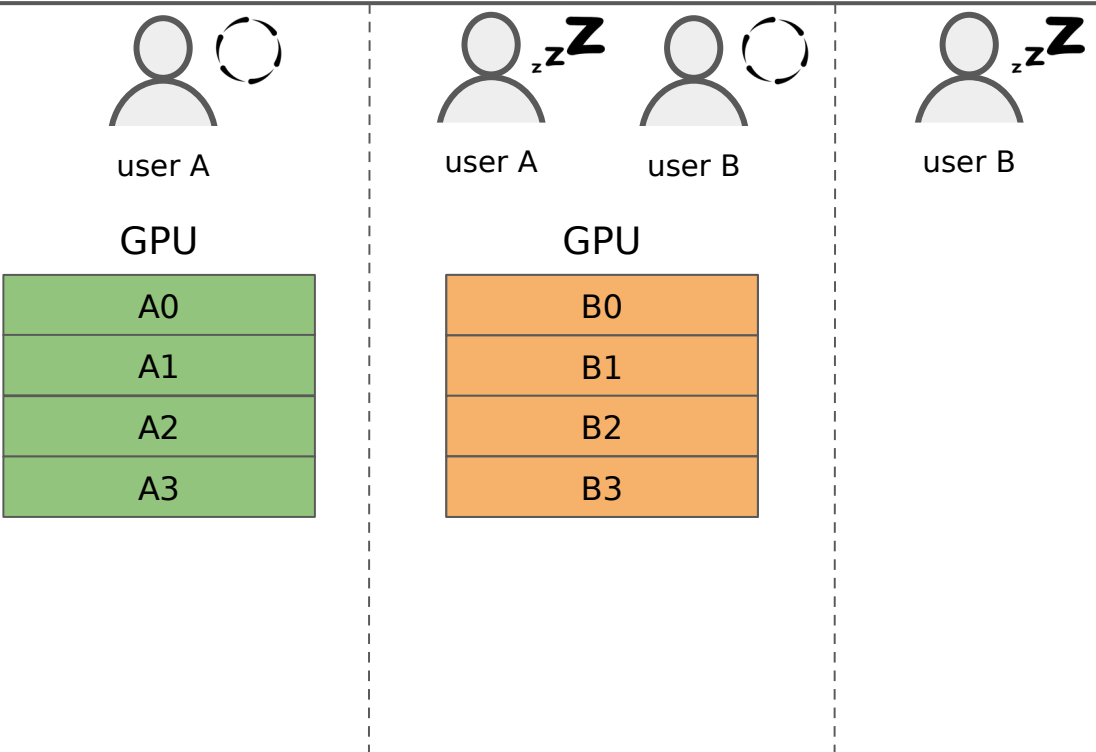
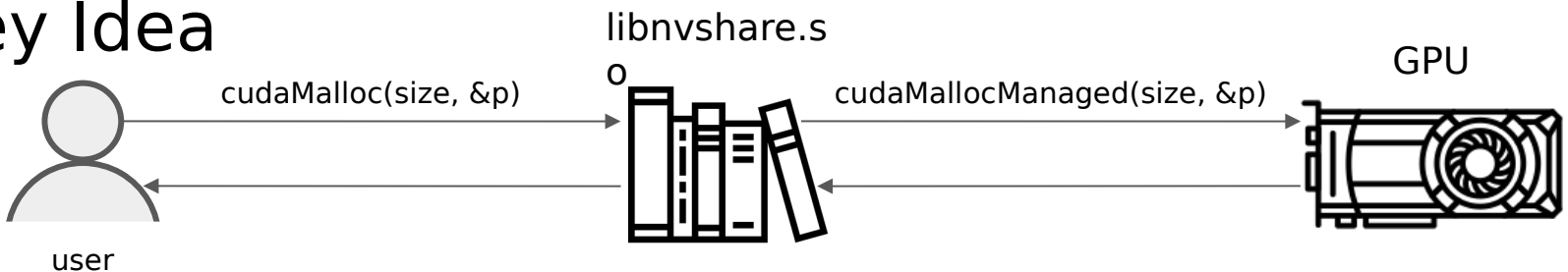
GPU

B0
B1
B2
B3

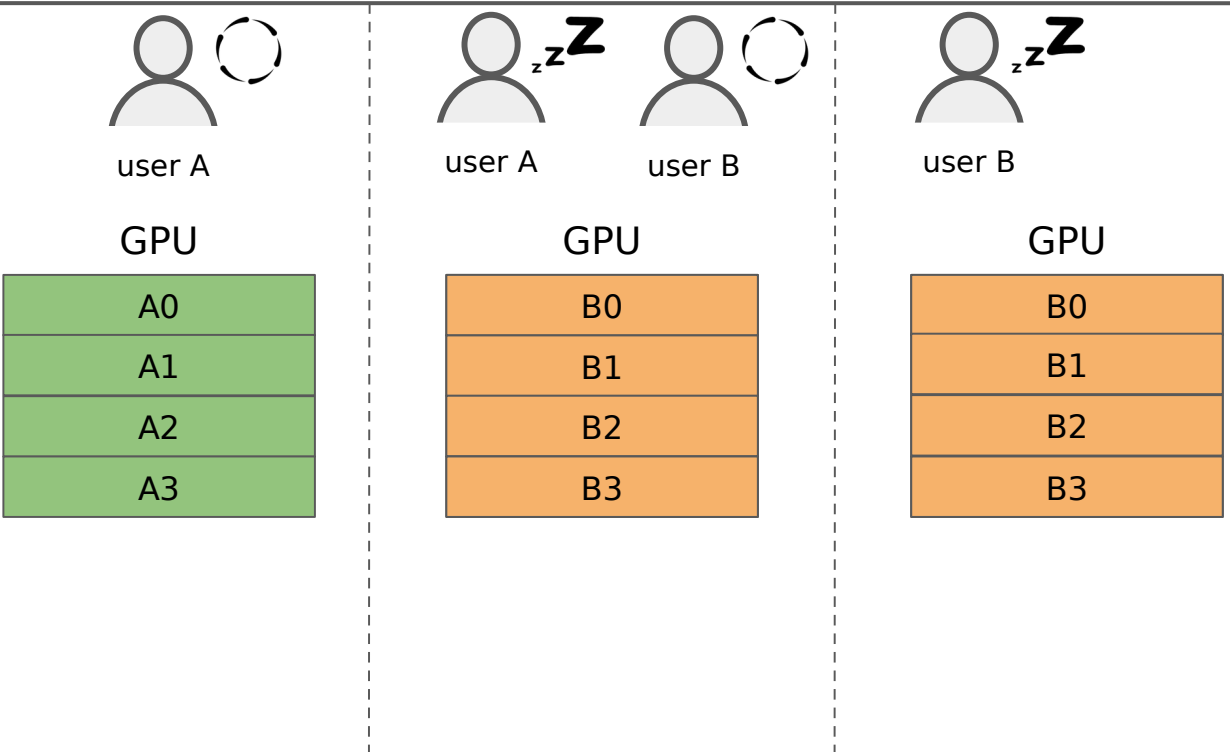
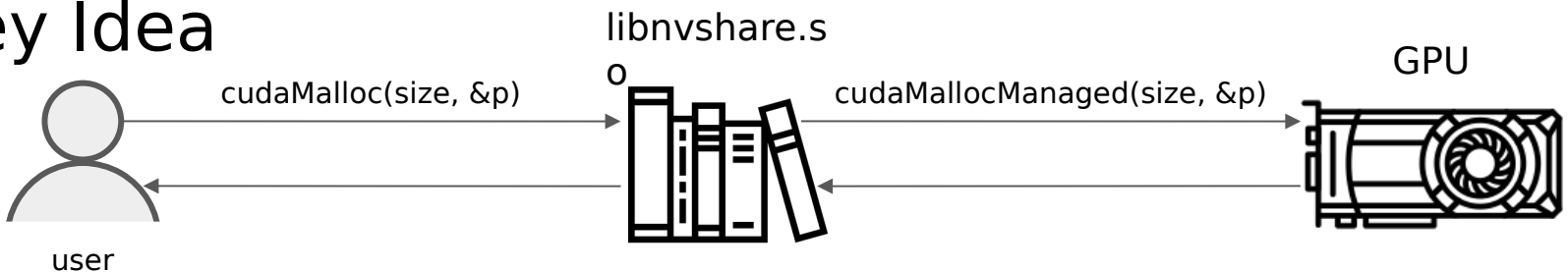
# Key Idea



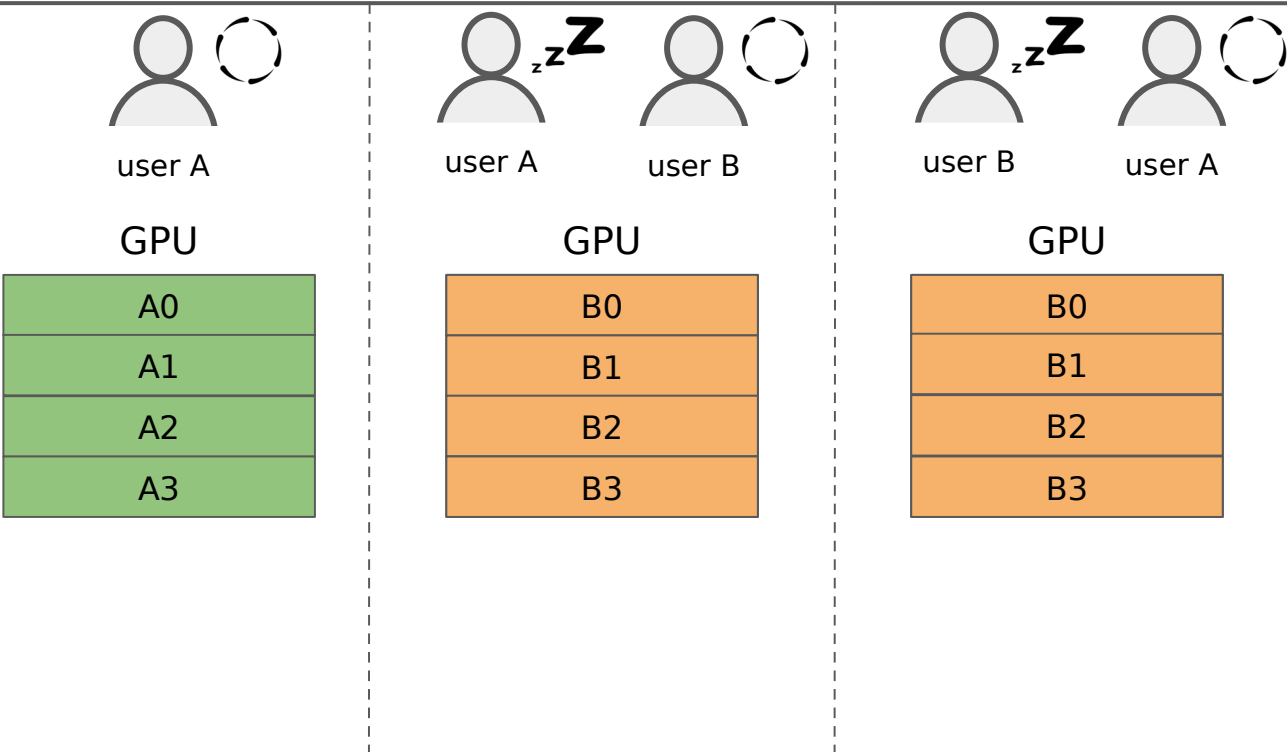
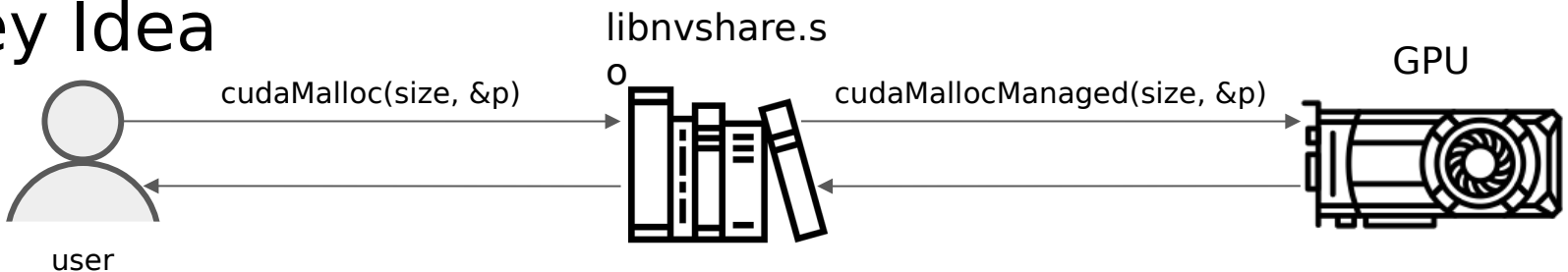
# Key Idea



# Key Idea

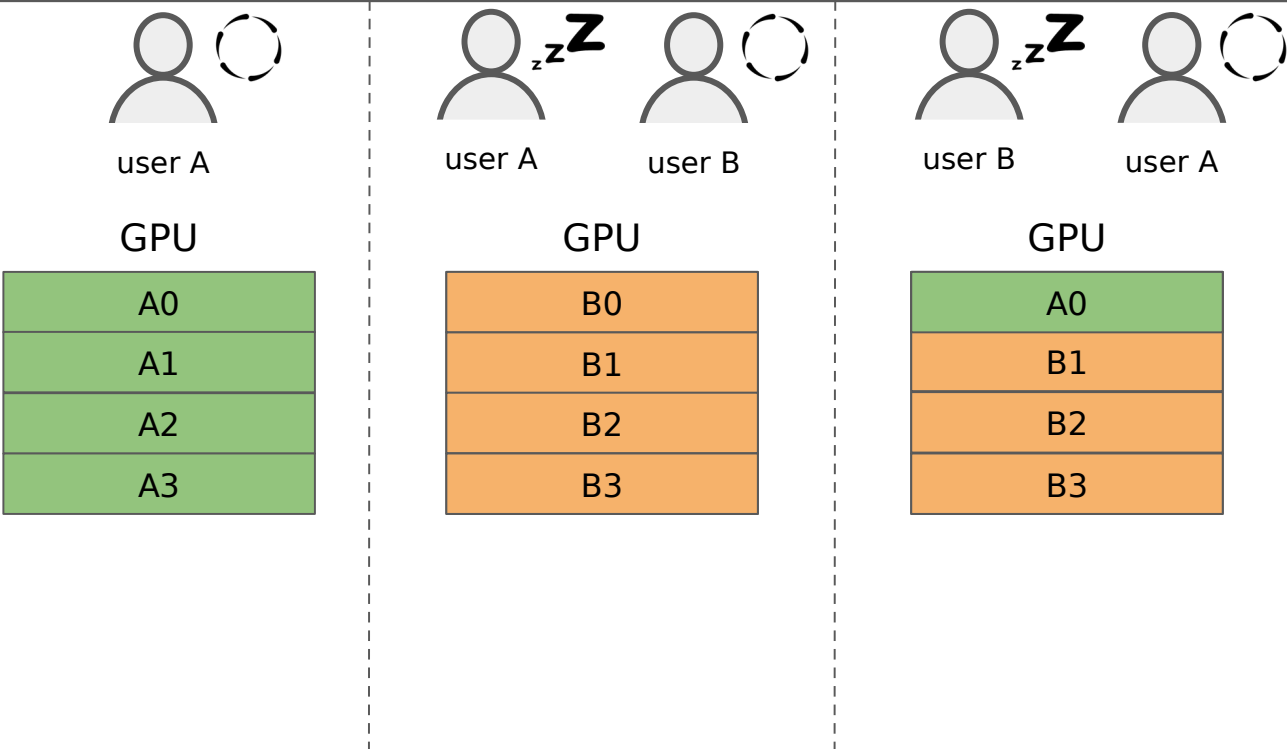
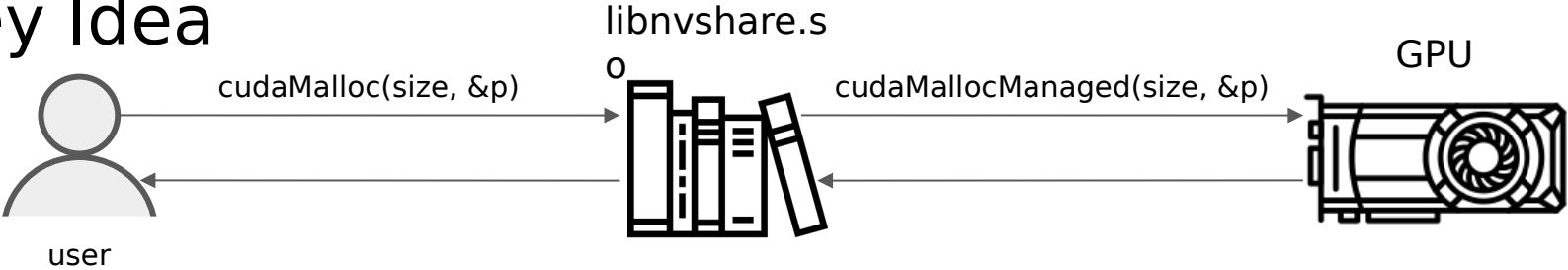


# Key Idea

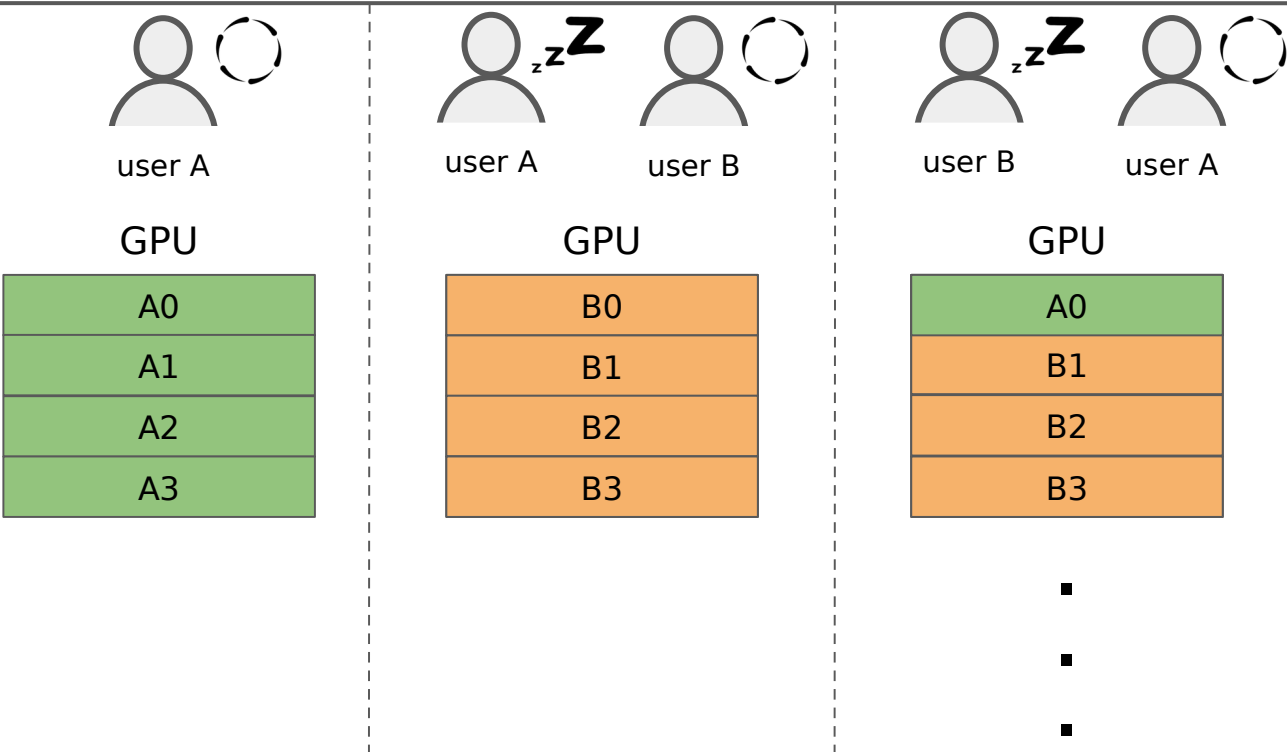
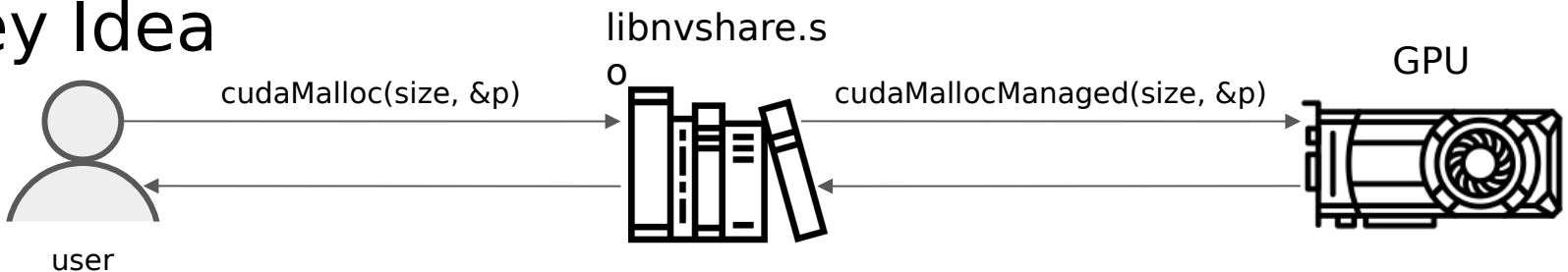




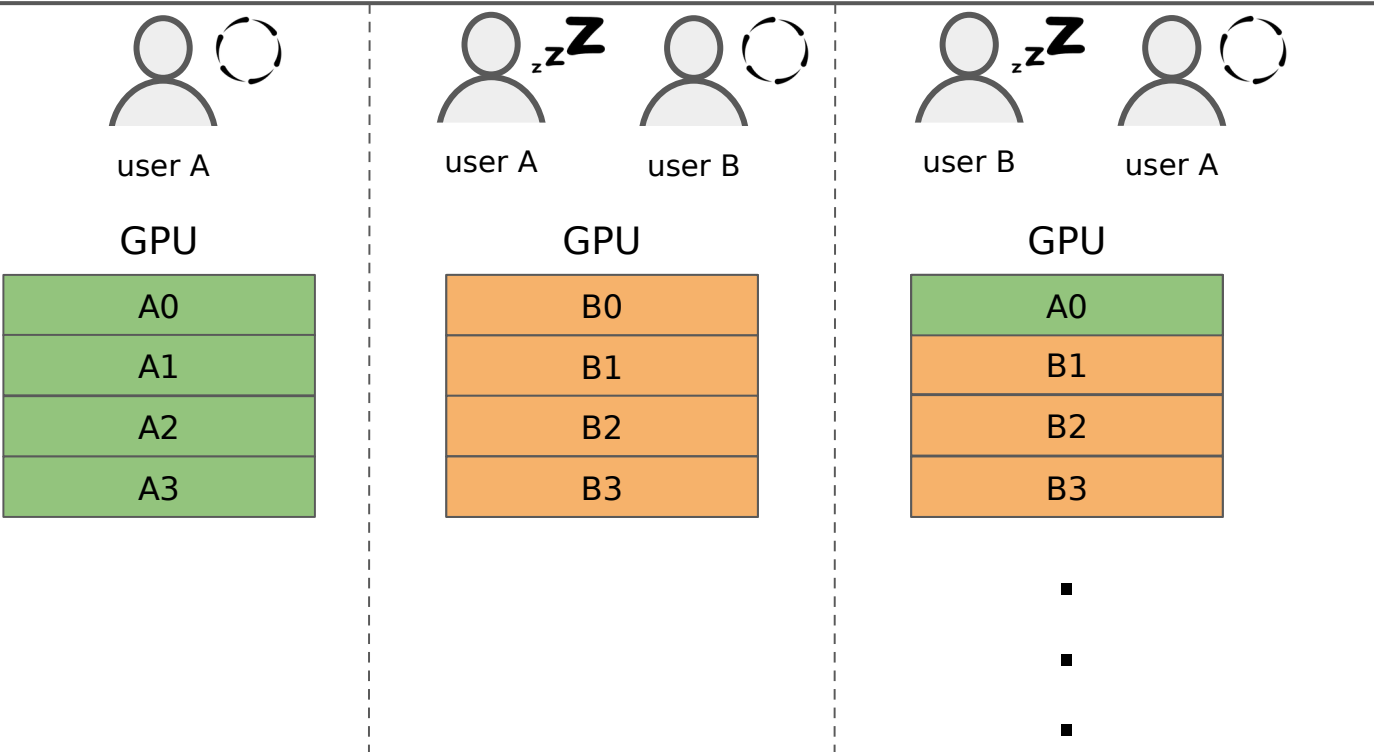
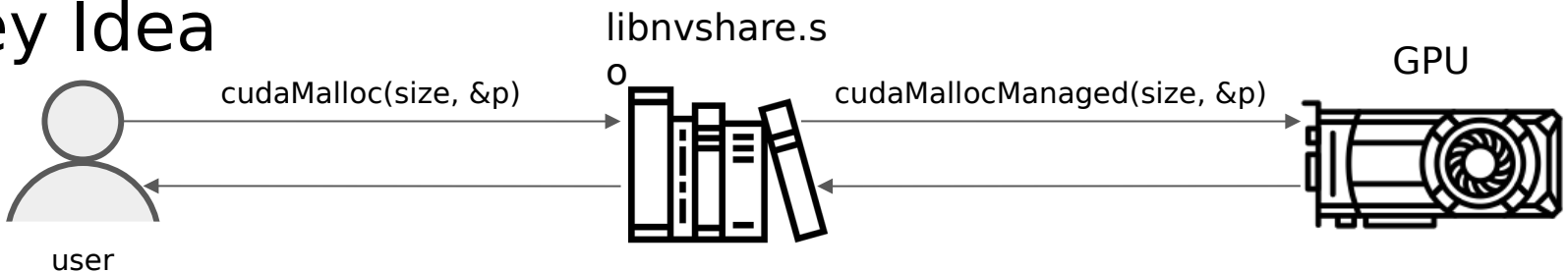
# Key Idea



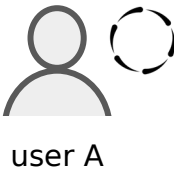
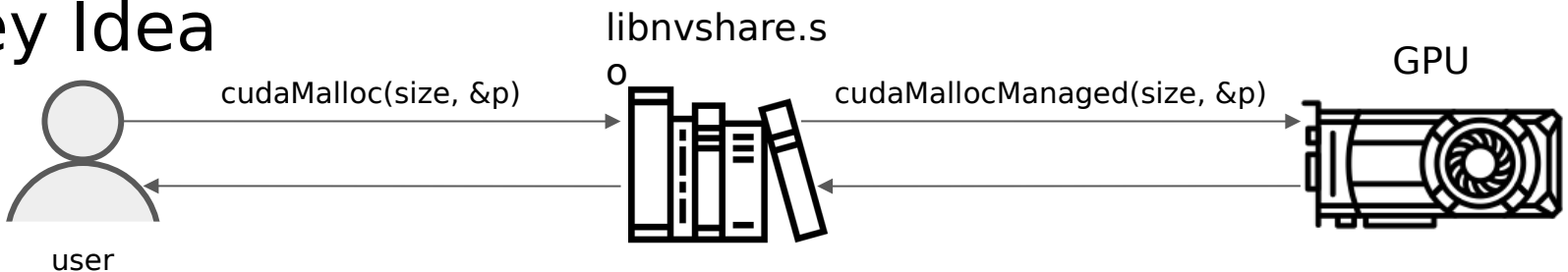
# Key Idea



# Key Idea

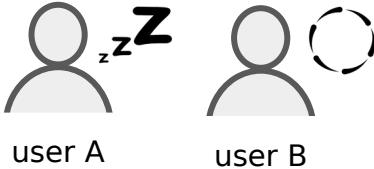


# Key Idea



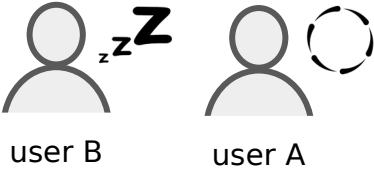
GPU

A0
A1
A2
A3



GPU

B0
B1
B2
B3



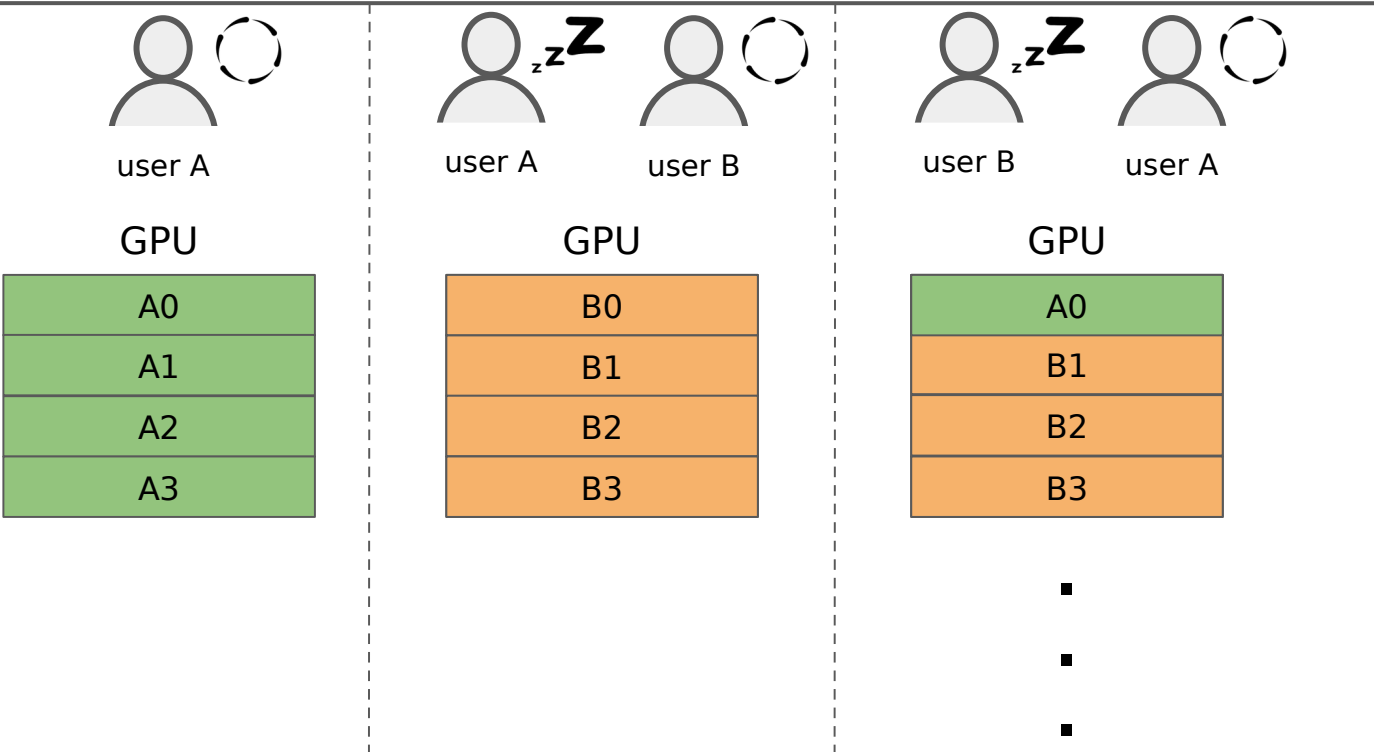
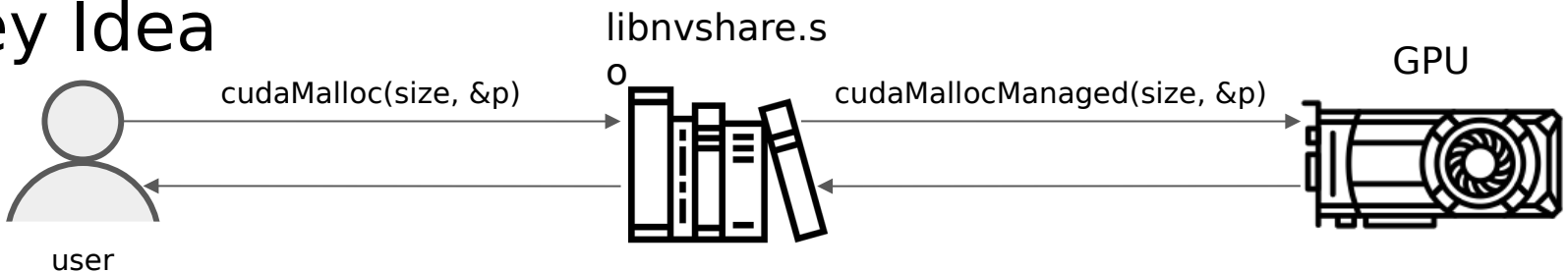
GPU

A0
B1
B2
B3

- 
- 
- 

Result

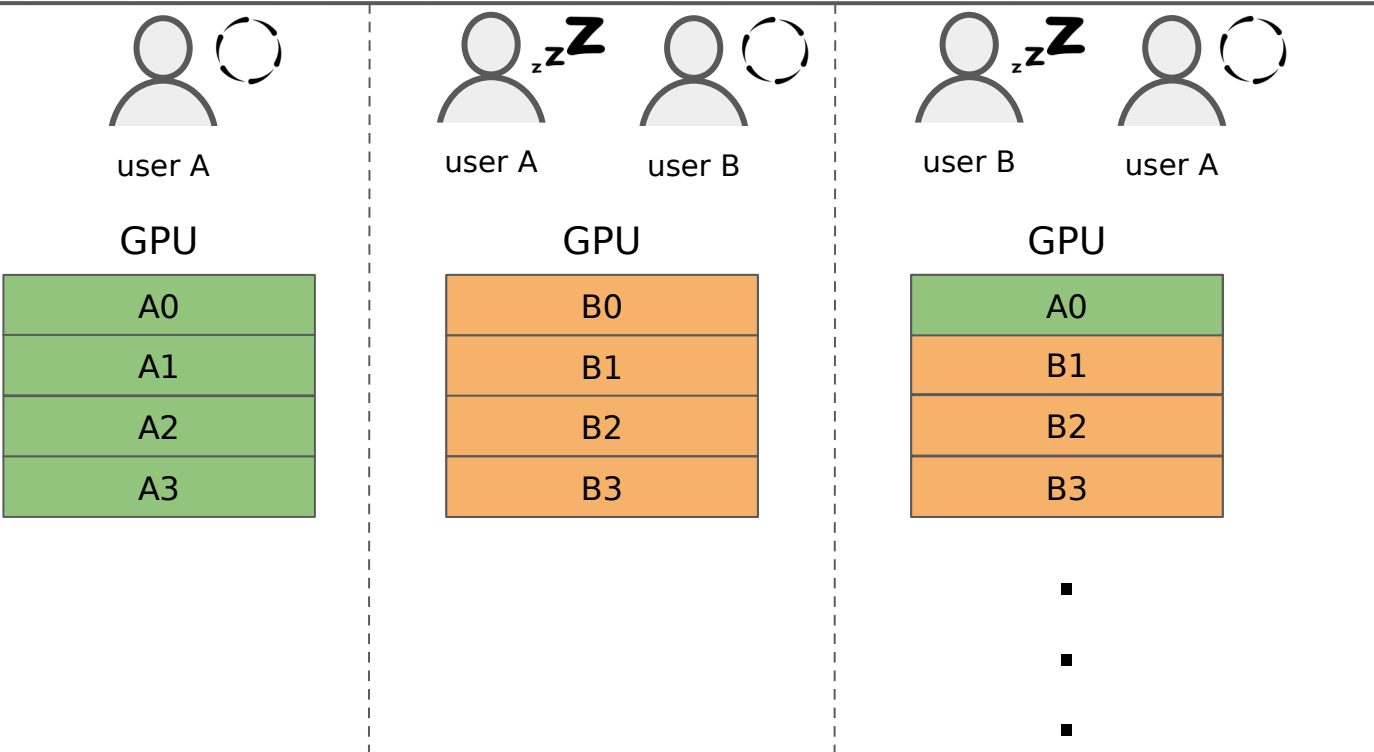
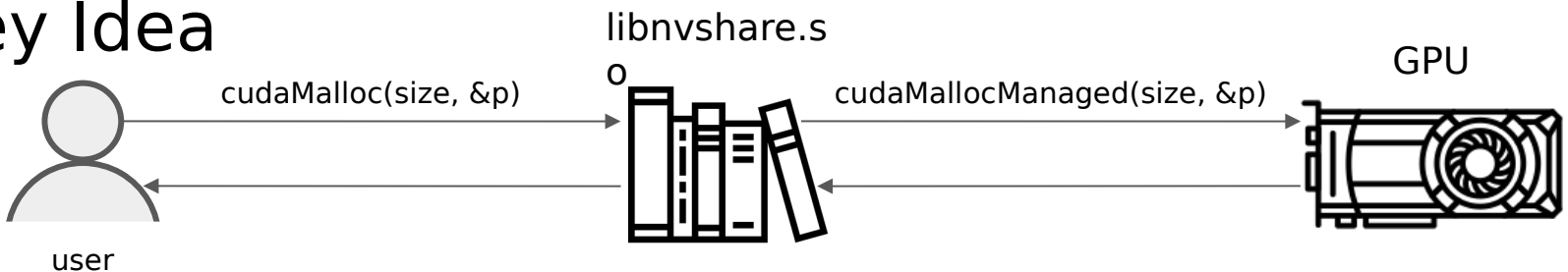
# Key Idea



## Result

- # of co-located processes limited by system RAM

# Key Idea



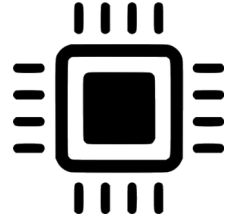
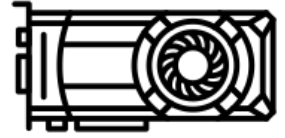
## Result

- # of co-located processes limited by system RAM
- each process can use the whole GPU memory

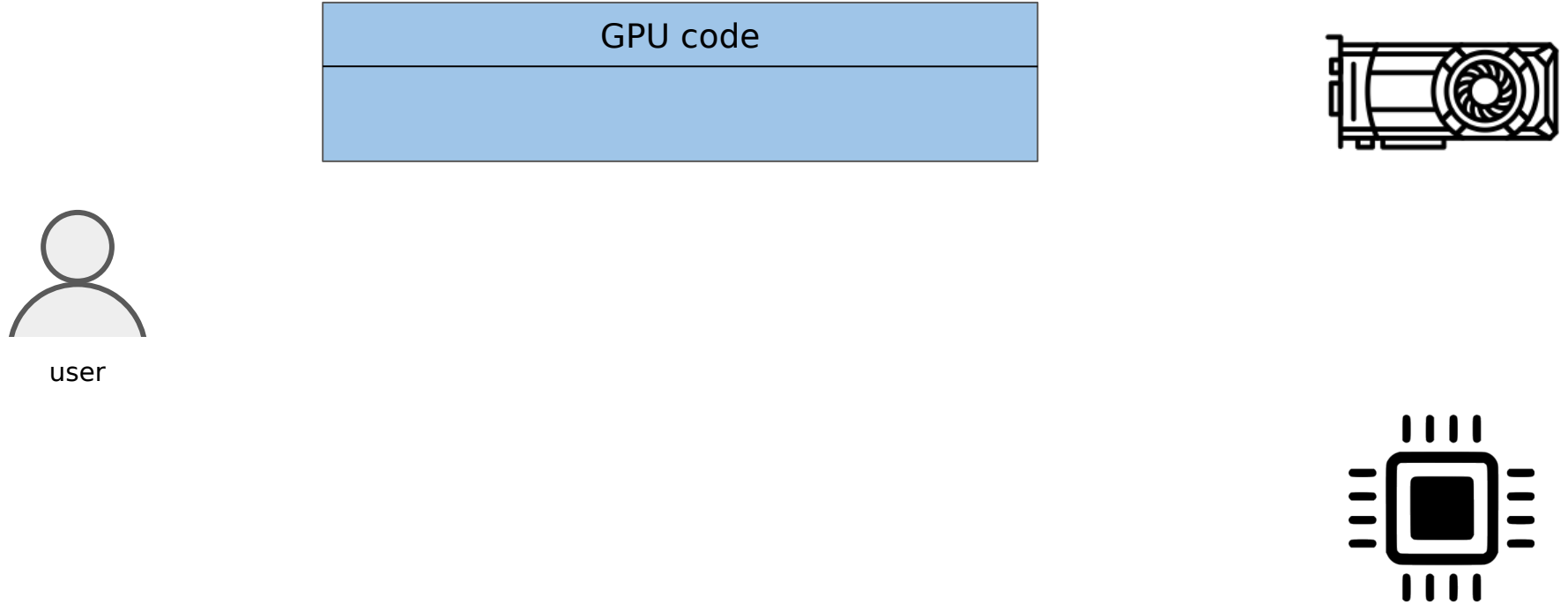
# CUDA programming model



user



# CUDA programming model

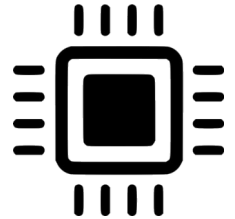
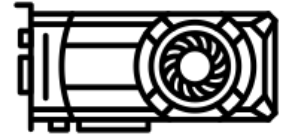
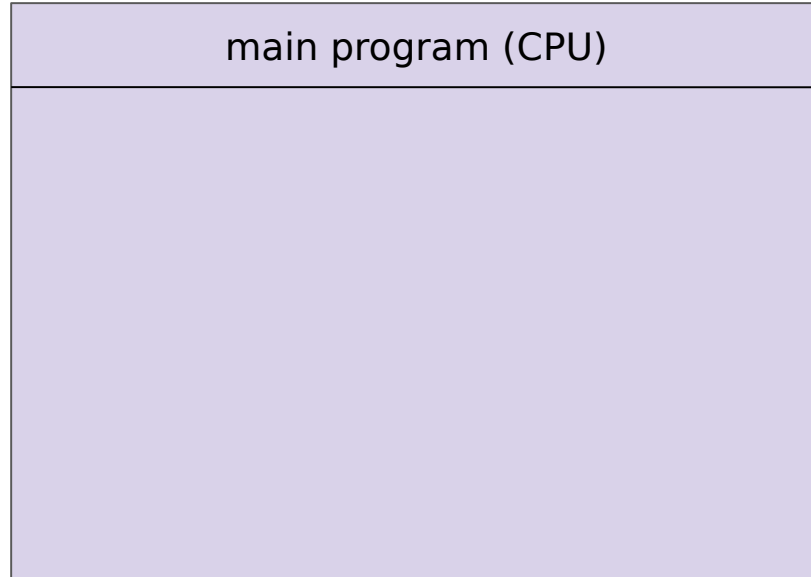
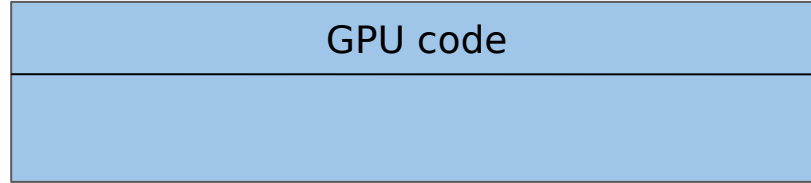




# CUDA programming model



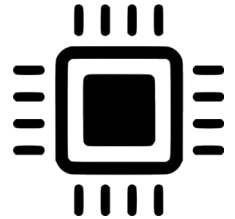
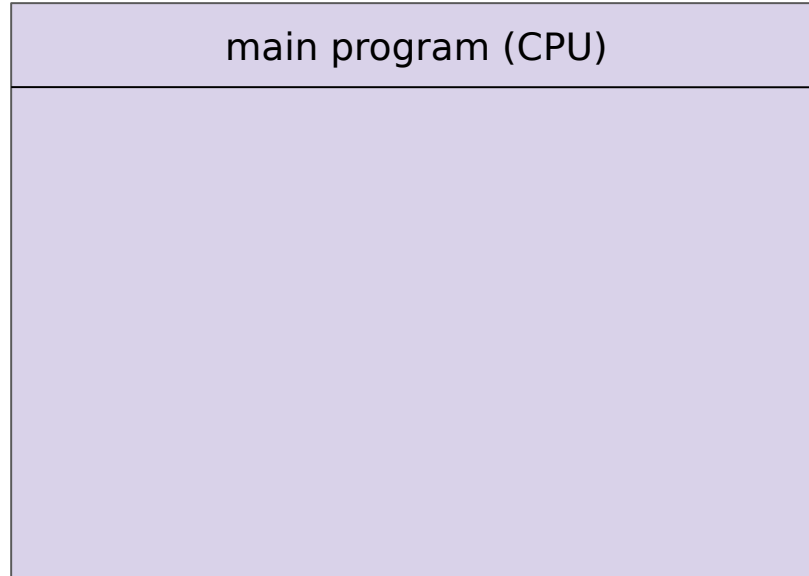
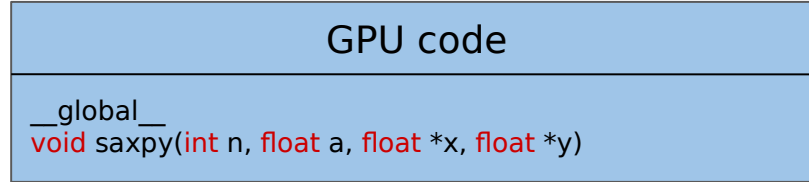
user



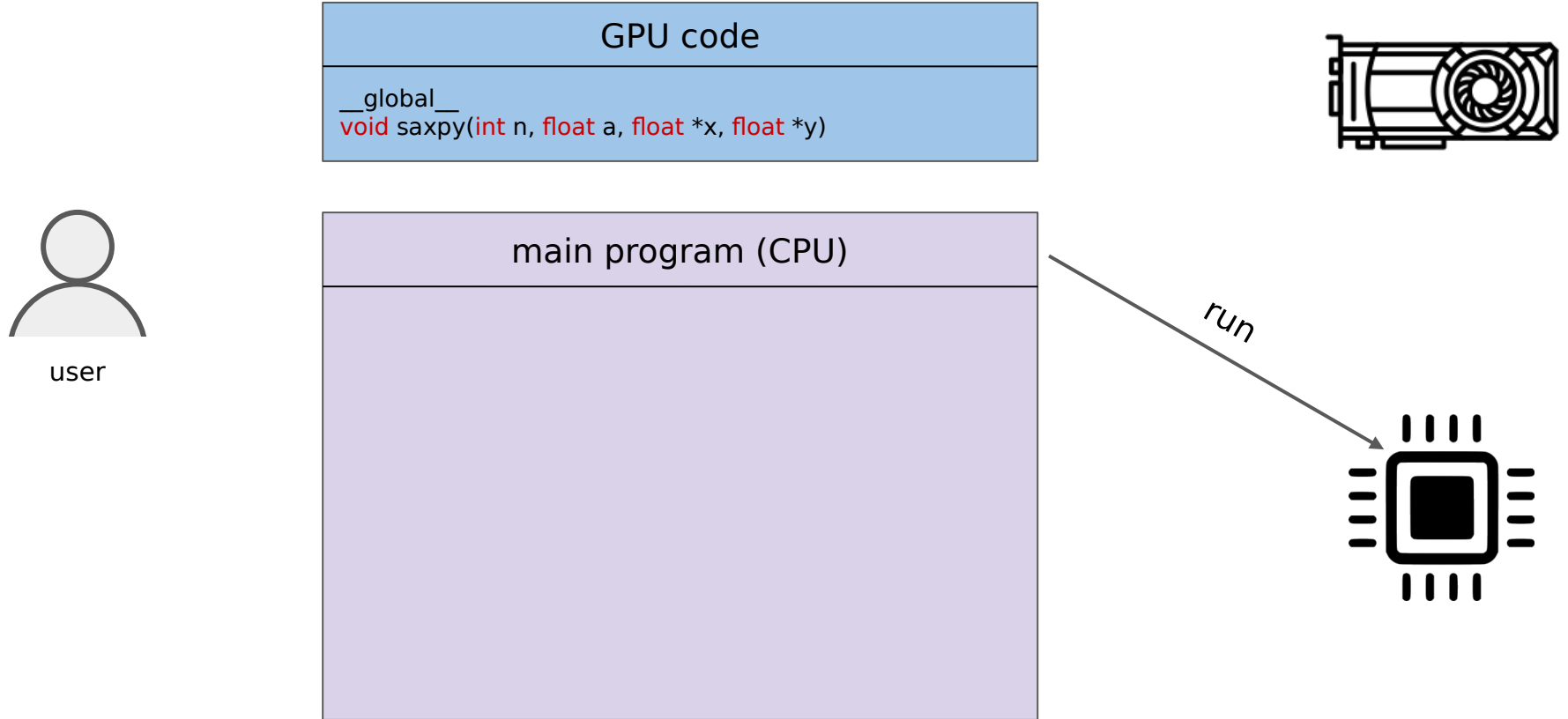
# CUDA programming model



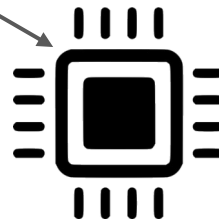
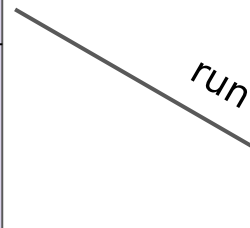
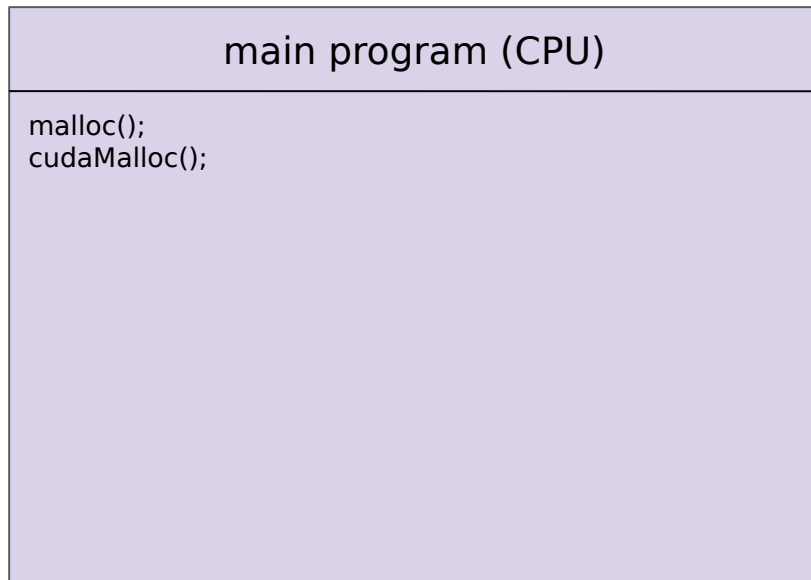
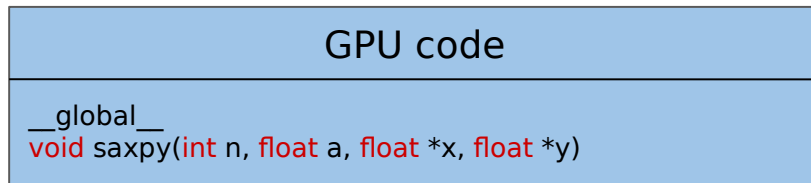
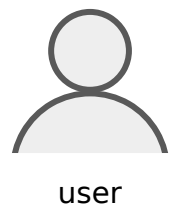
user



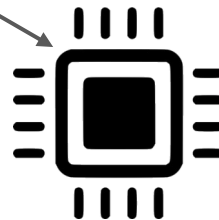
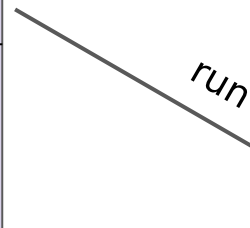
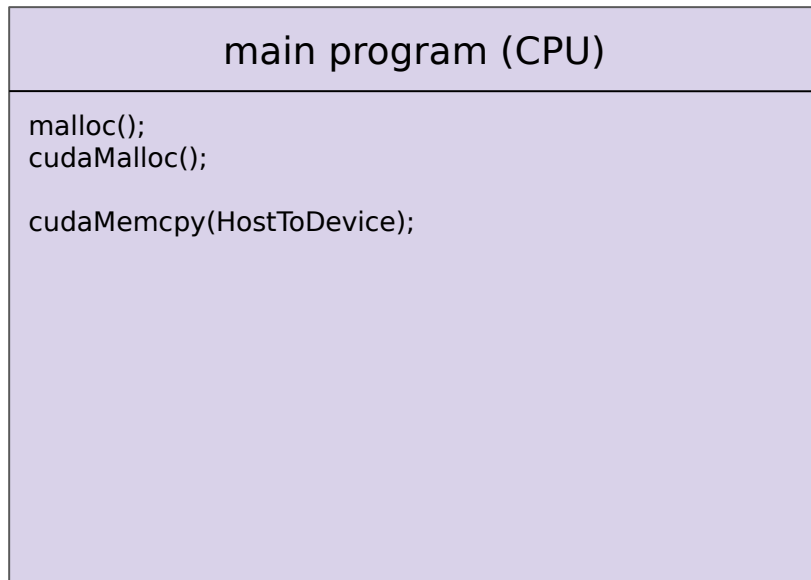
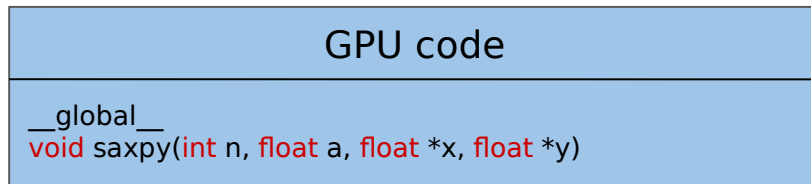
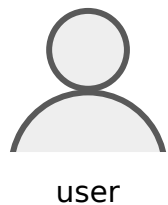
# CUDA programming model



# CUDA programming model



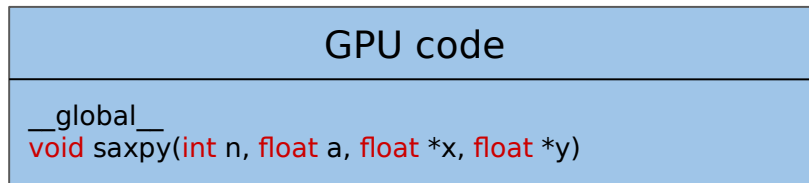
# CUDA programming model



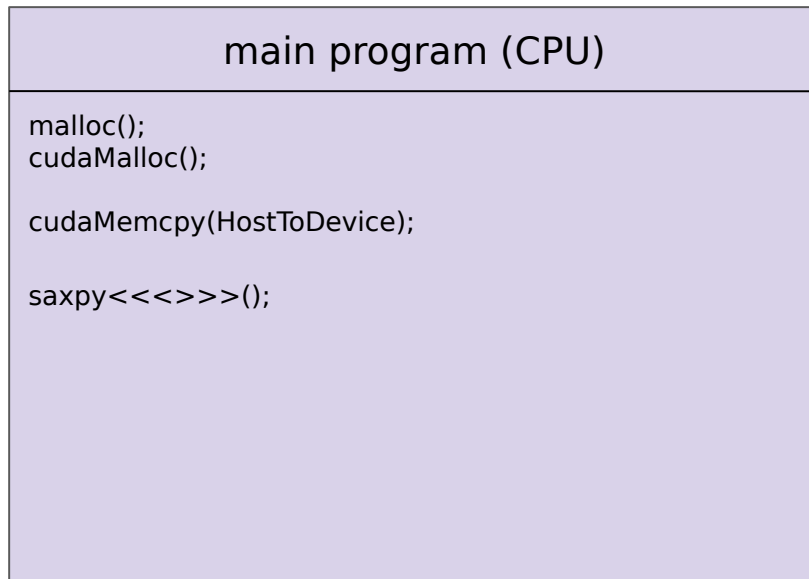
# CUDA programming model



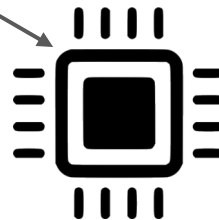
user



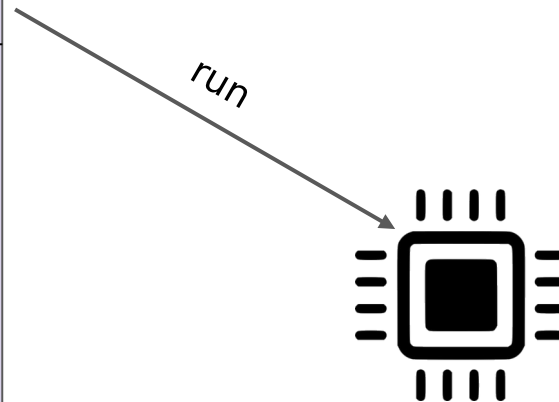
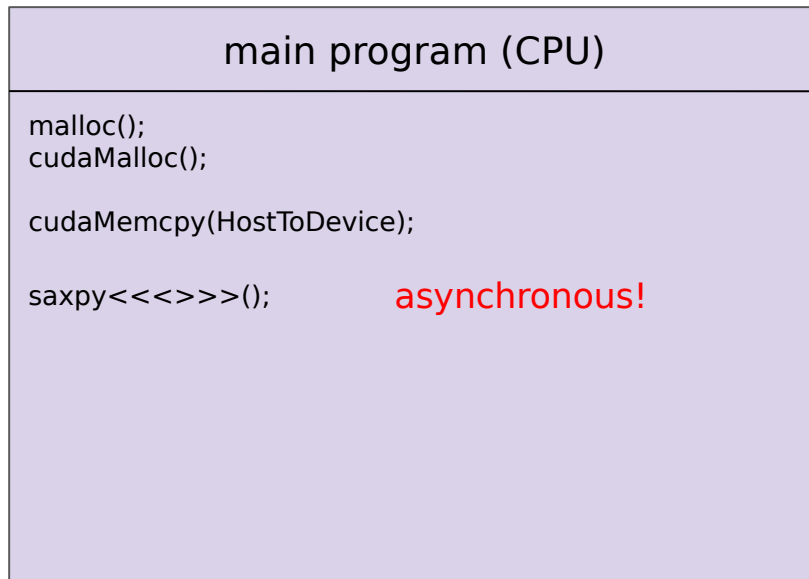
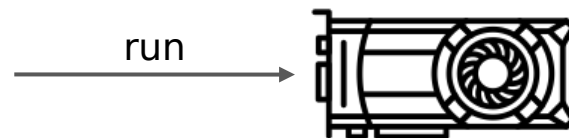
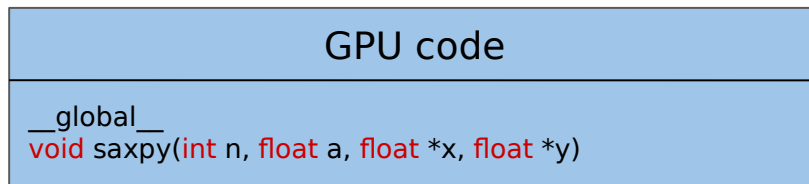
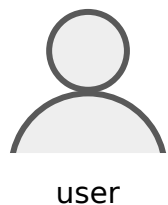
run



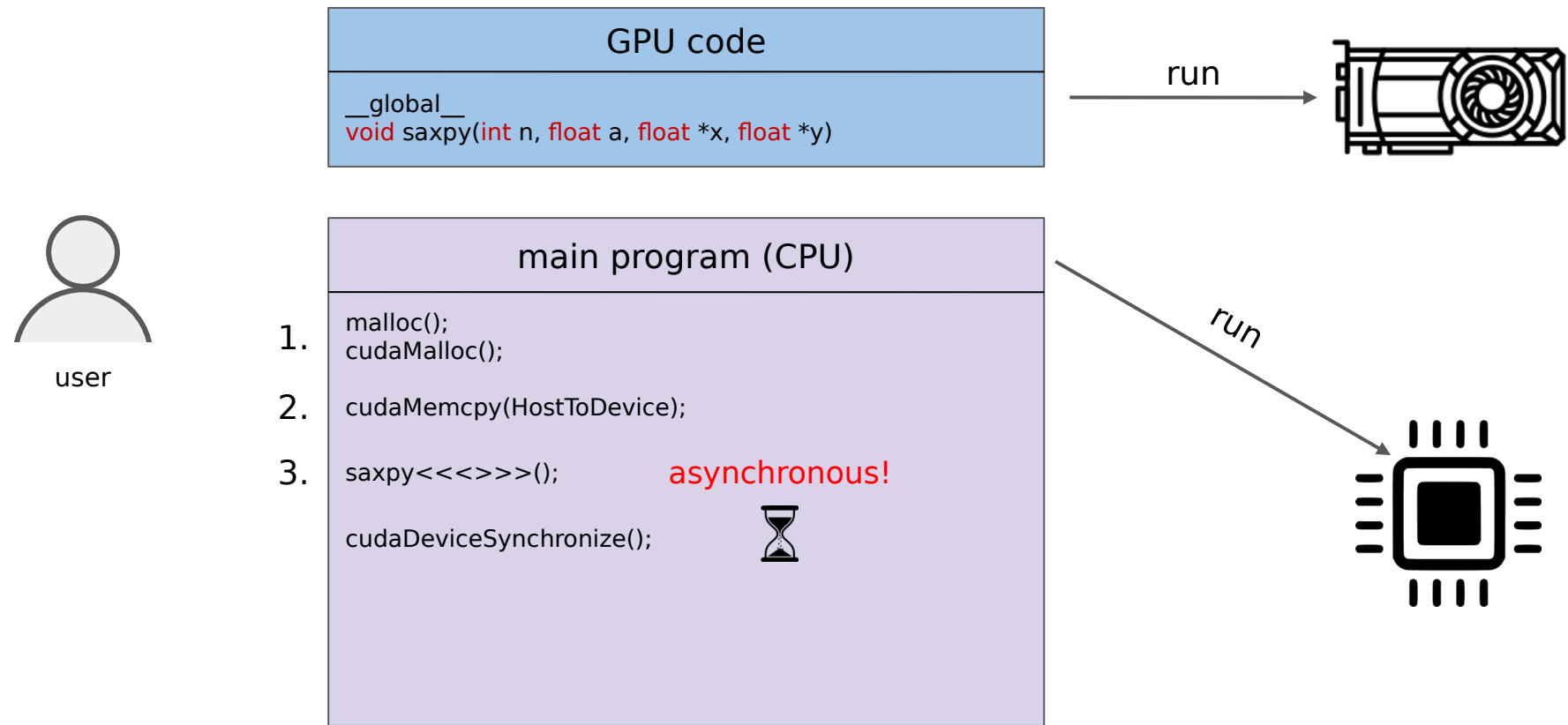
run



# CUDA programming model

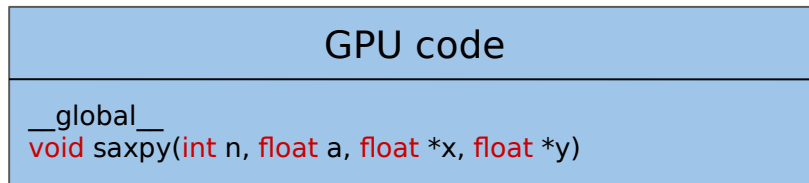
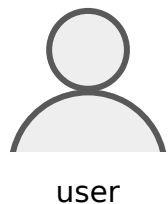


# CUDA programming model

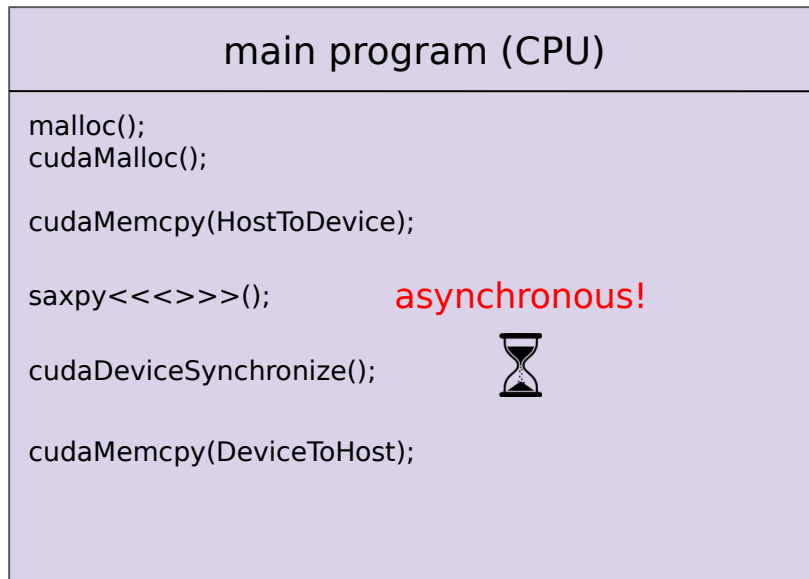




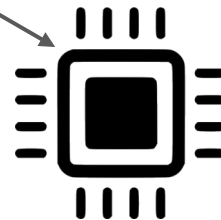
# CUDA programming model



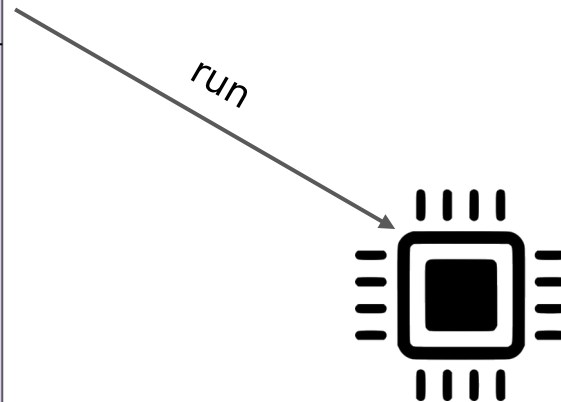
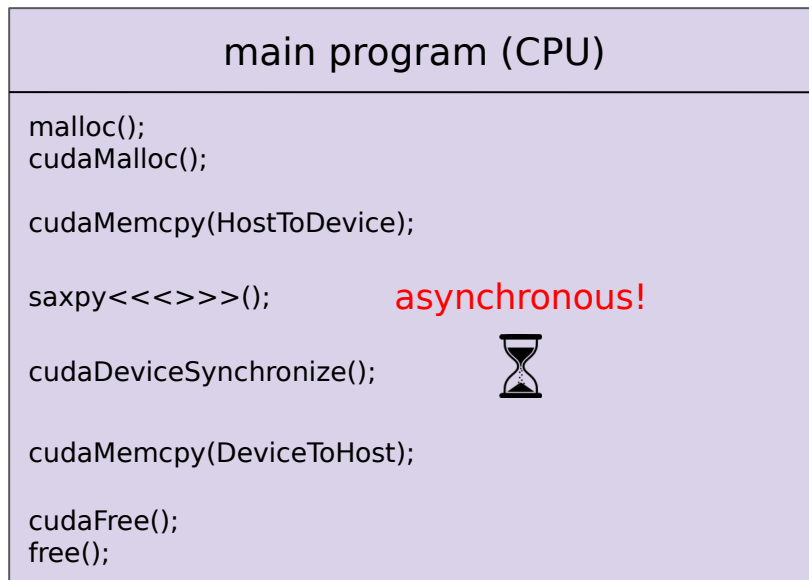
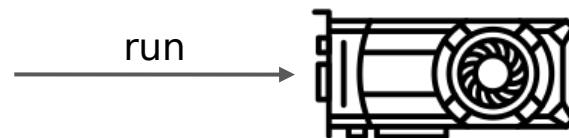
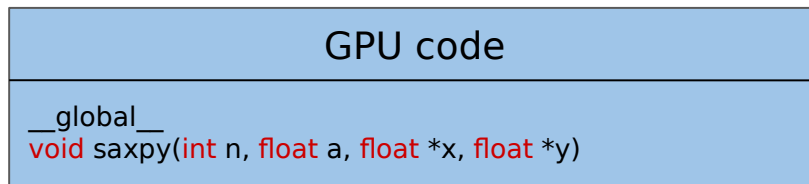
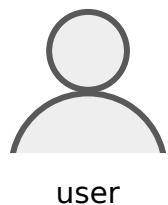
run



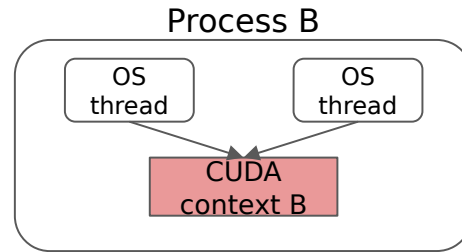
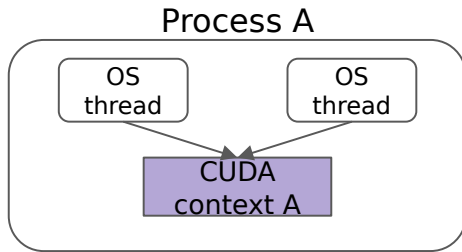
run



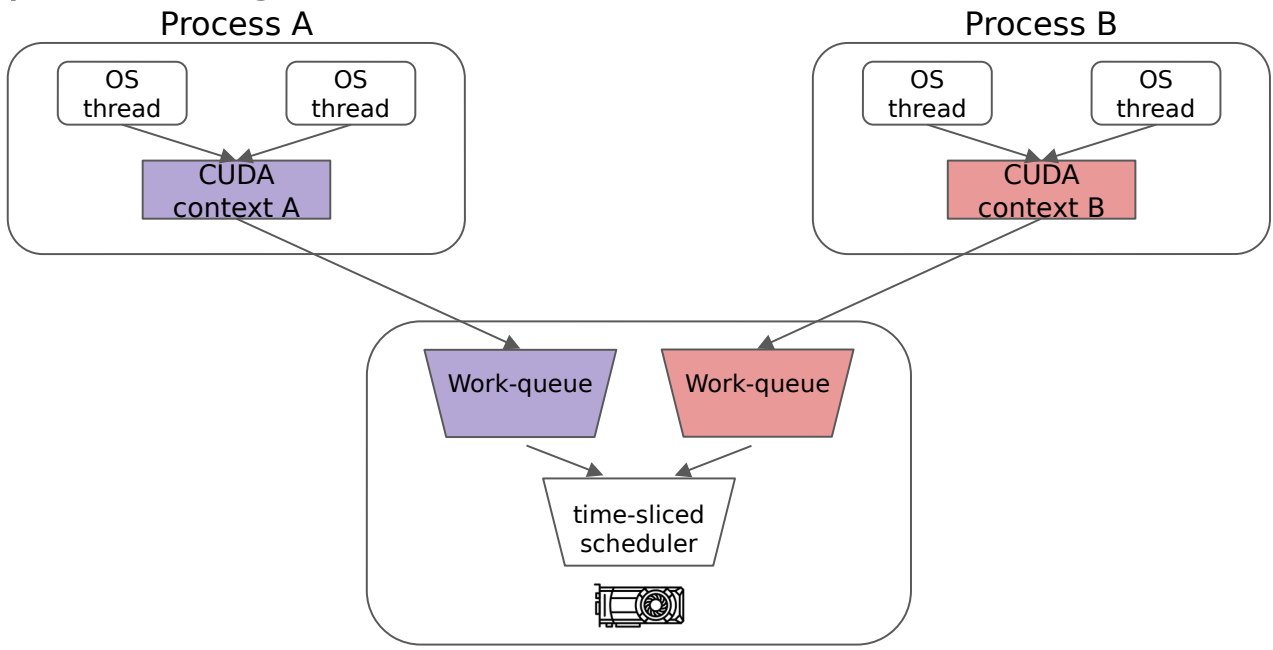
# CUDA programming model



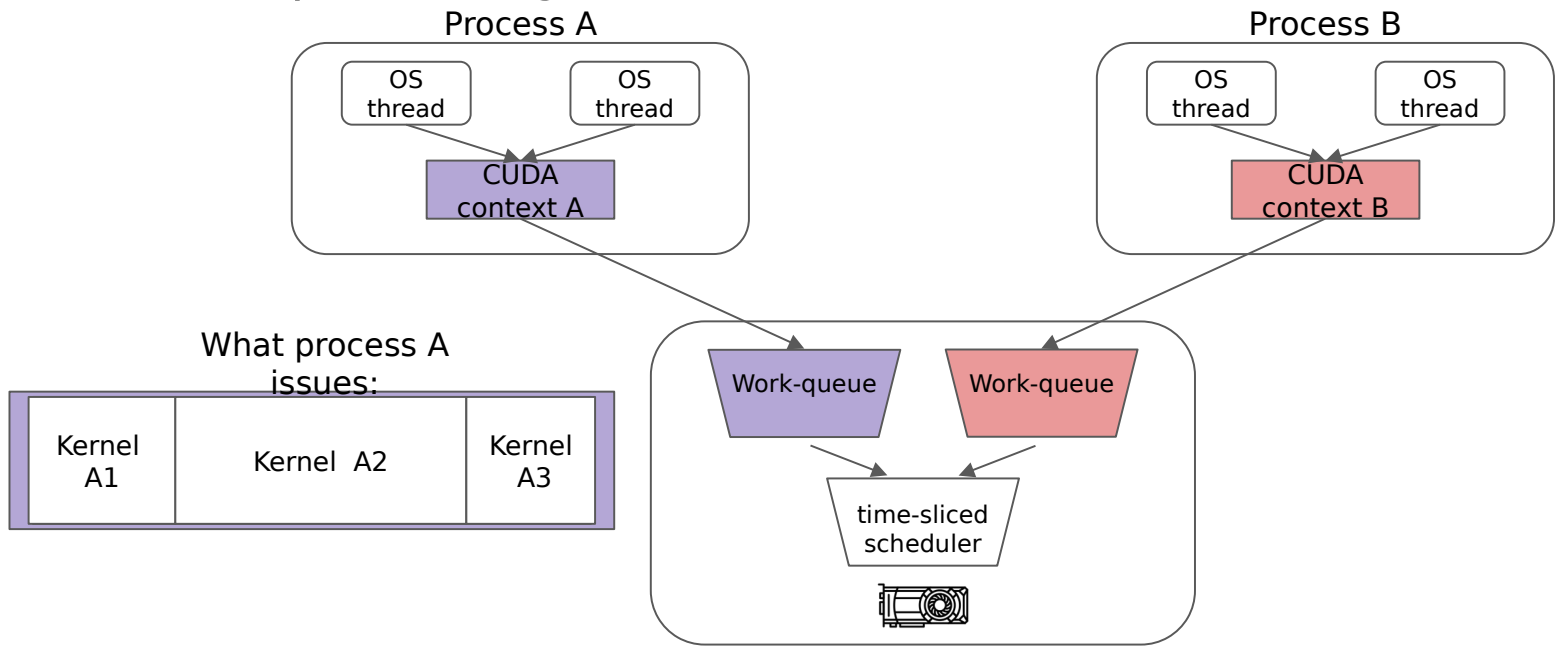
# CUDA multi-processing



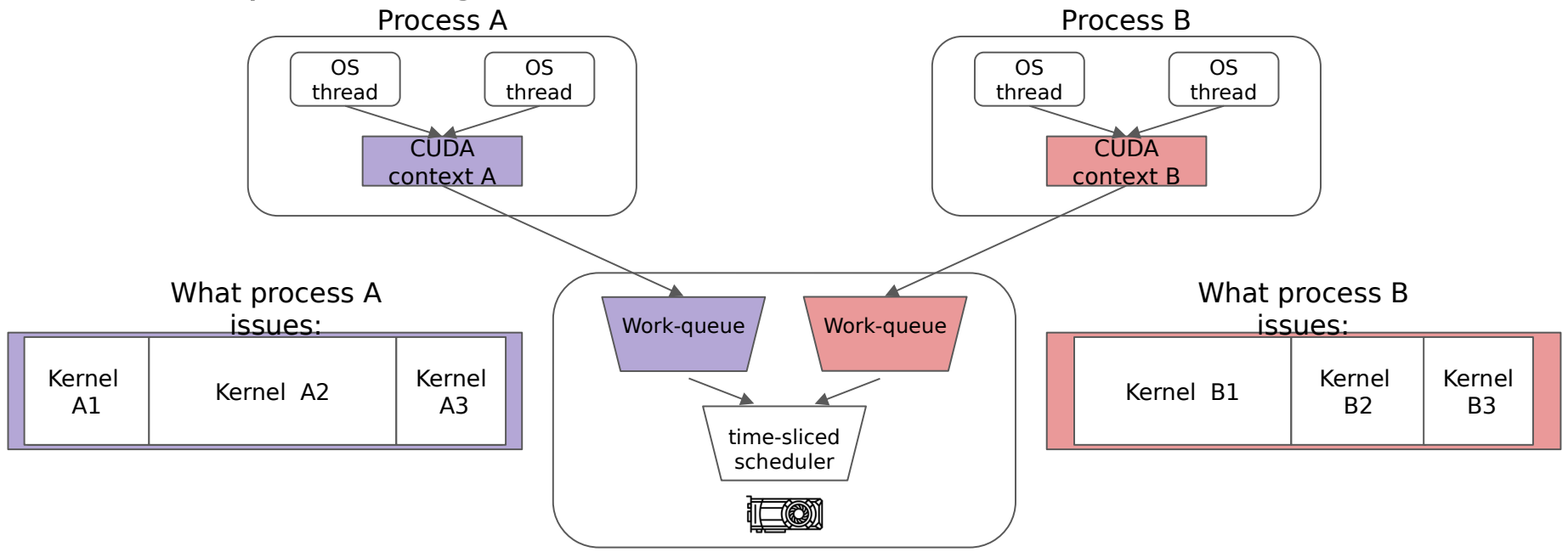
# CUDA multi-processing



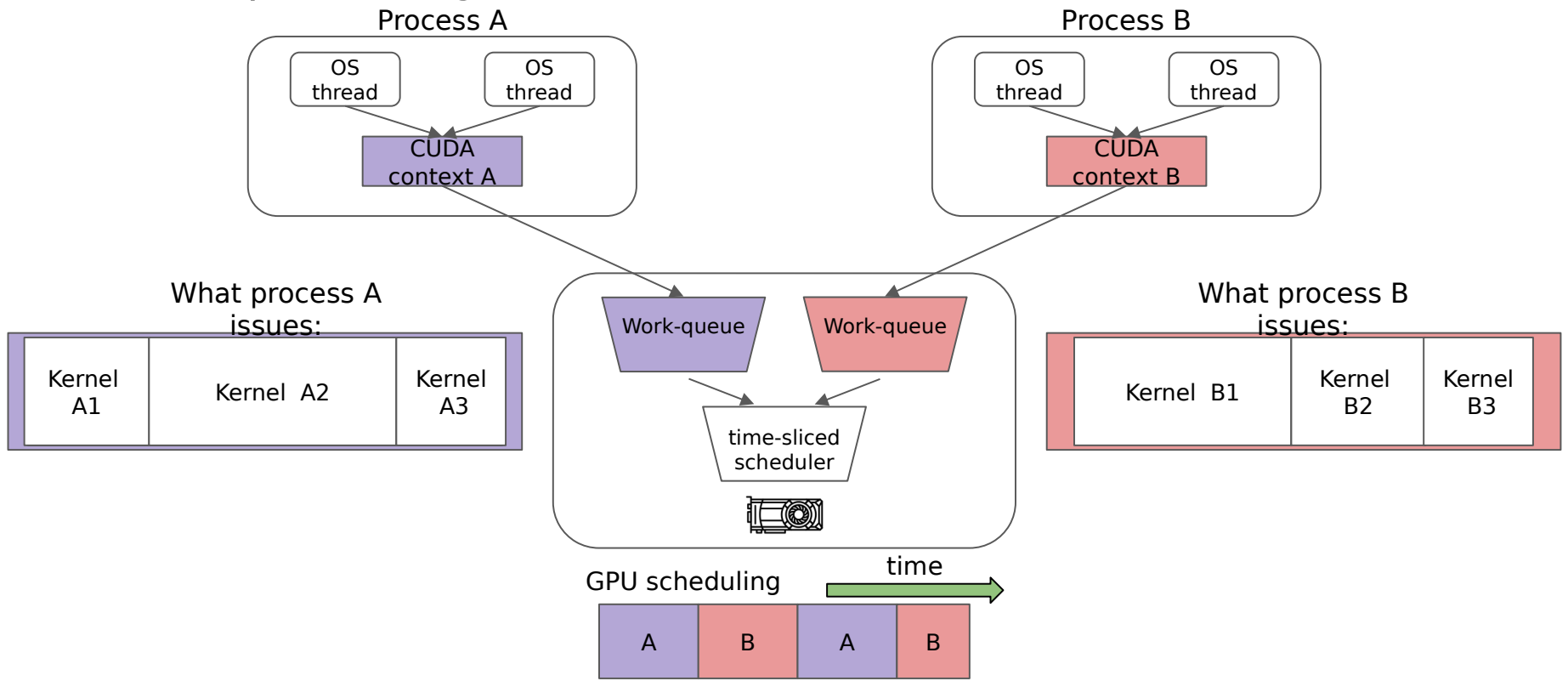
# CUDA multi-processing



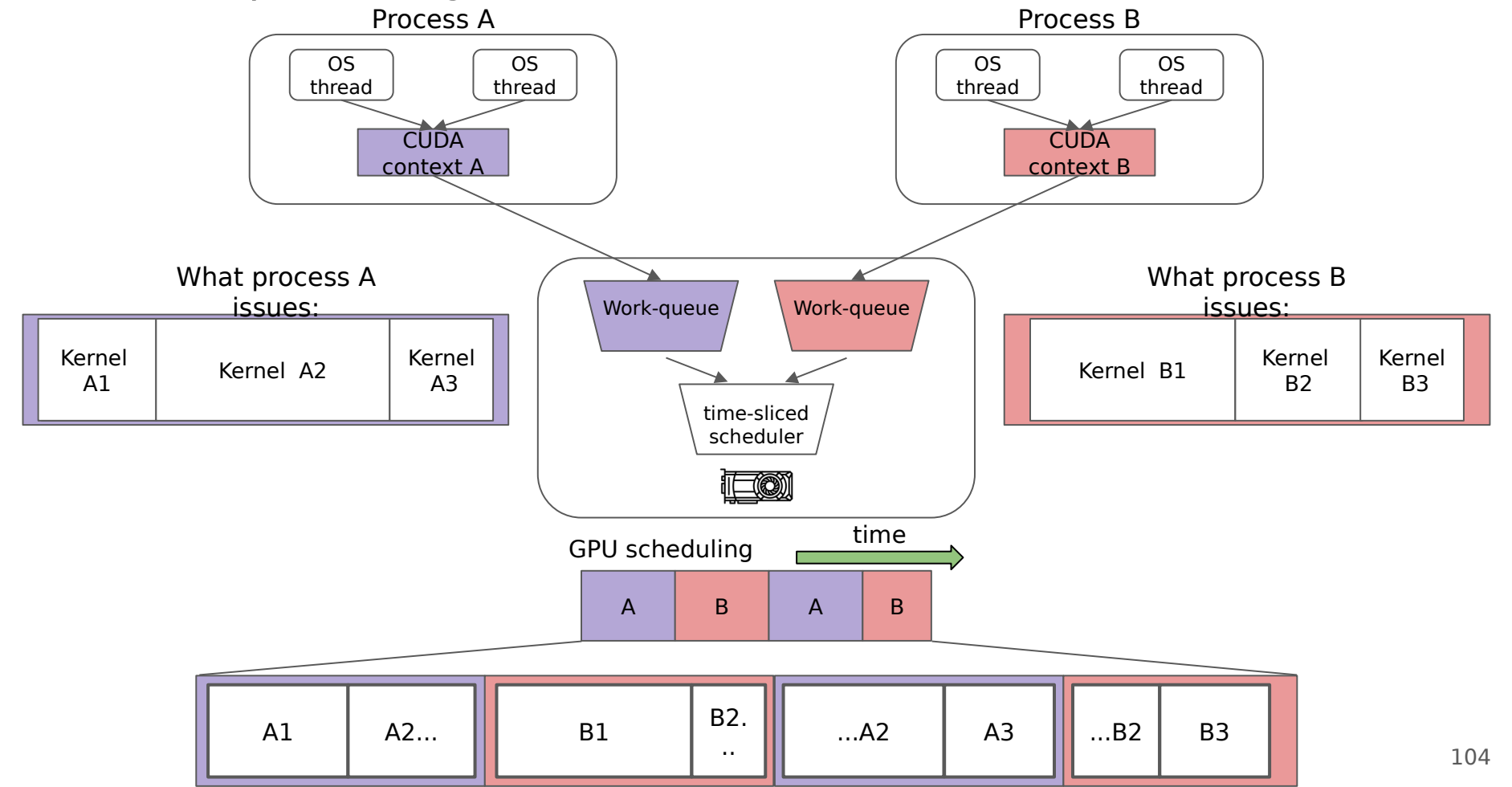
# CUDA multi-processing



# CUDA multi-processing

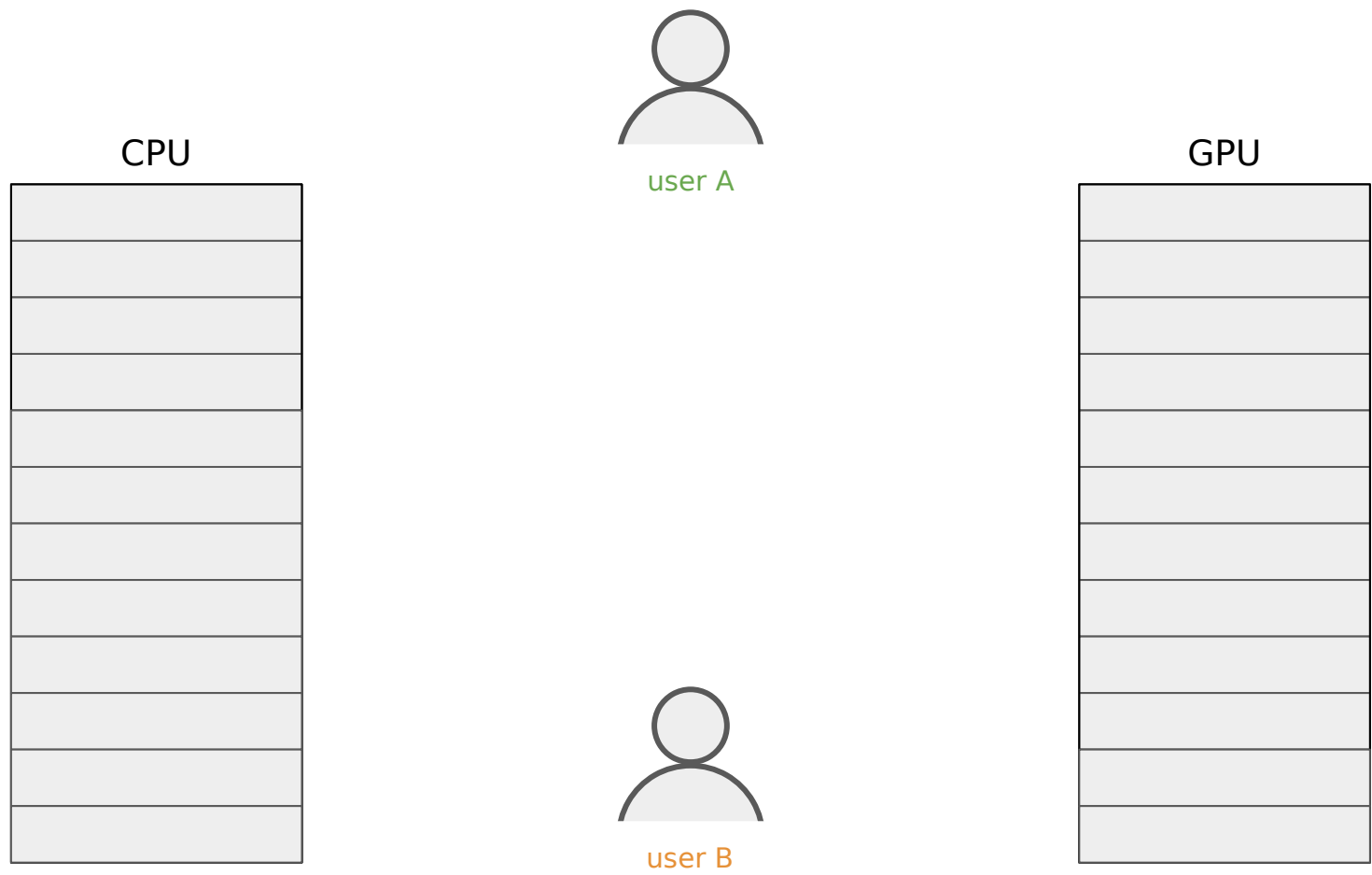


# CUDA multi-processing

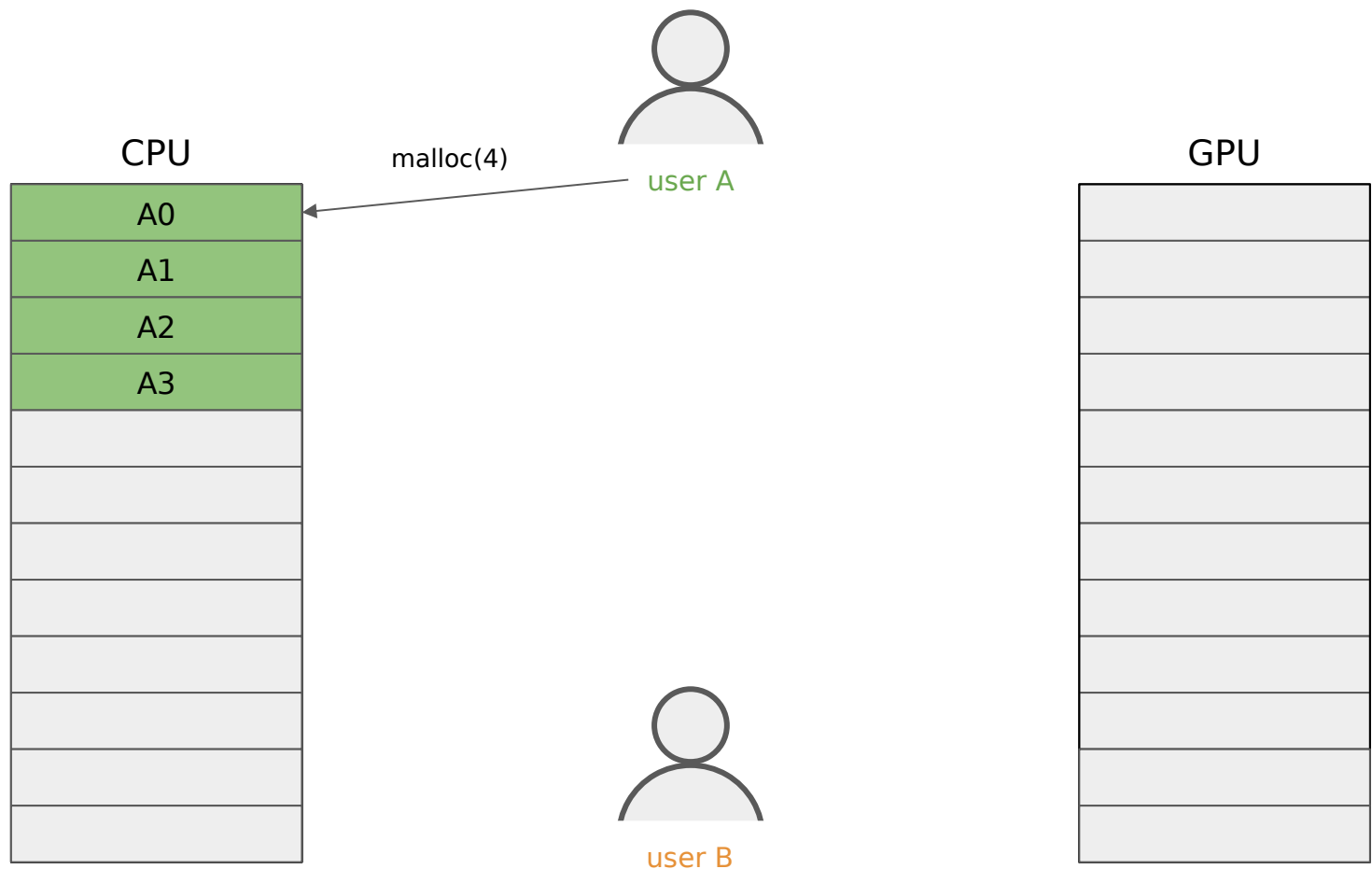




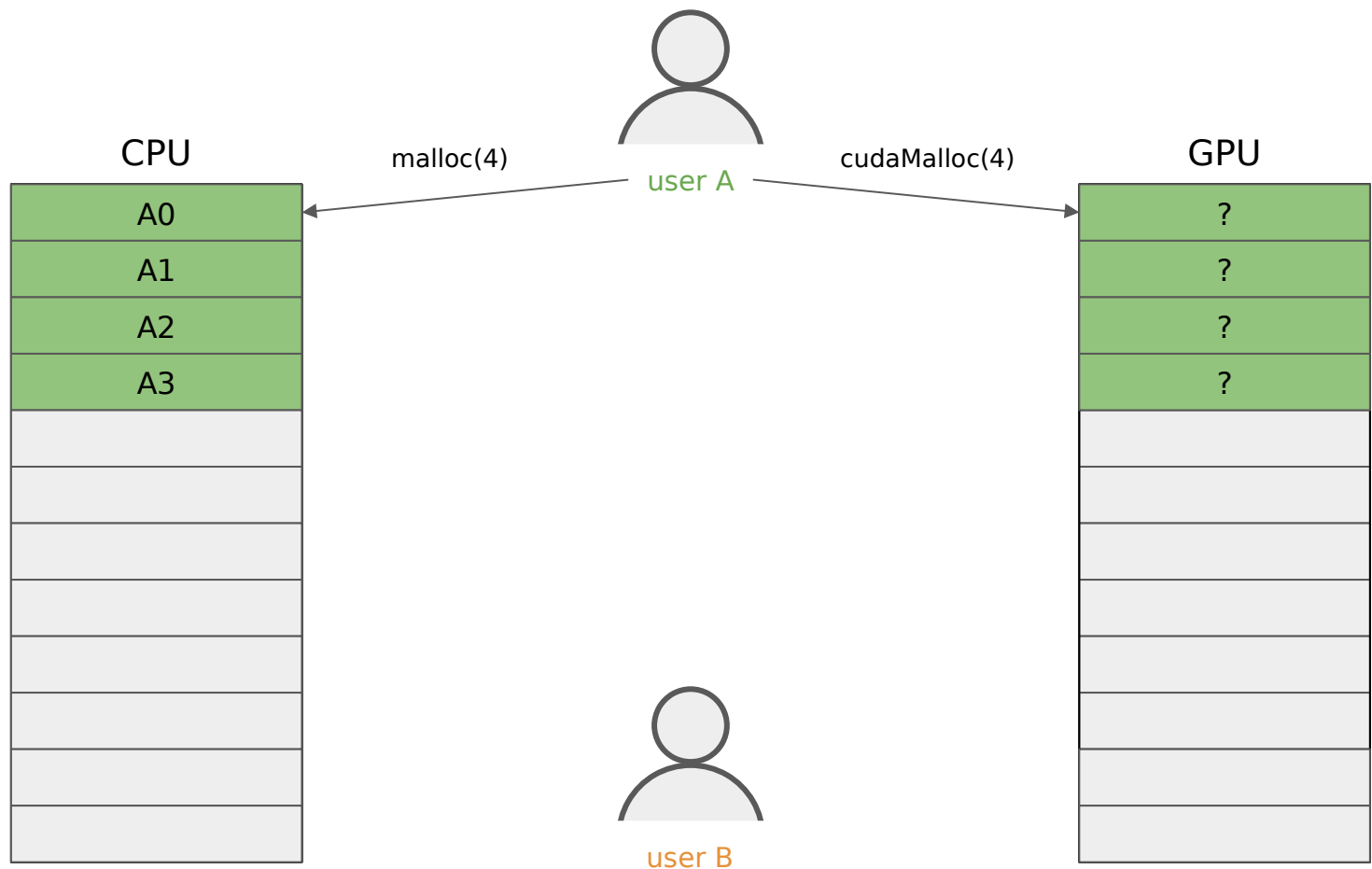
# Default CUDA Memory Model



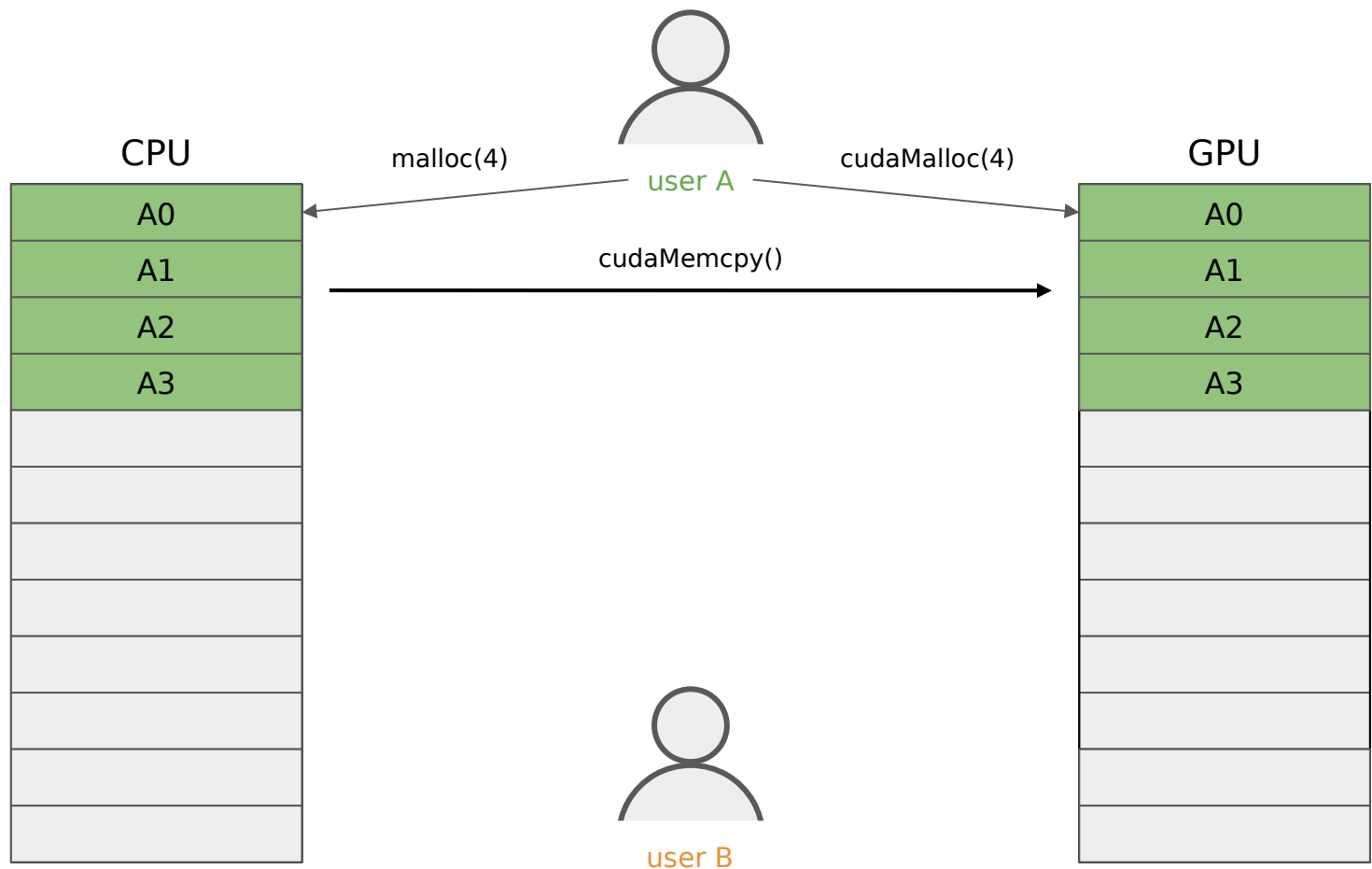
# Default CUDA Memory Model



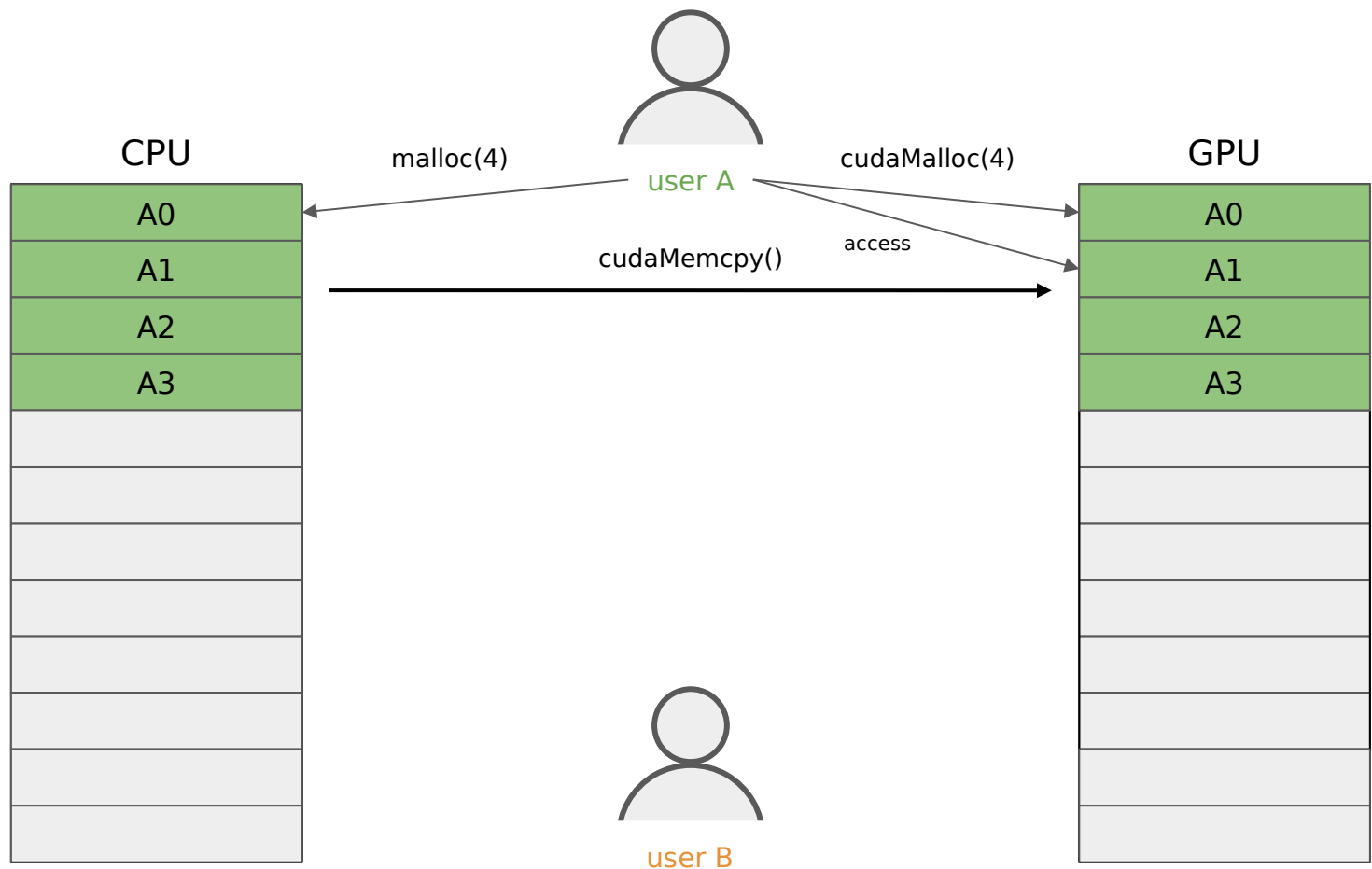
# Default CUDA Memory Model



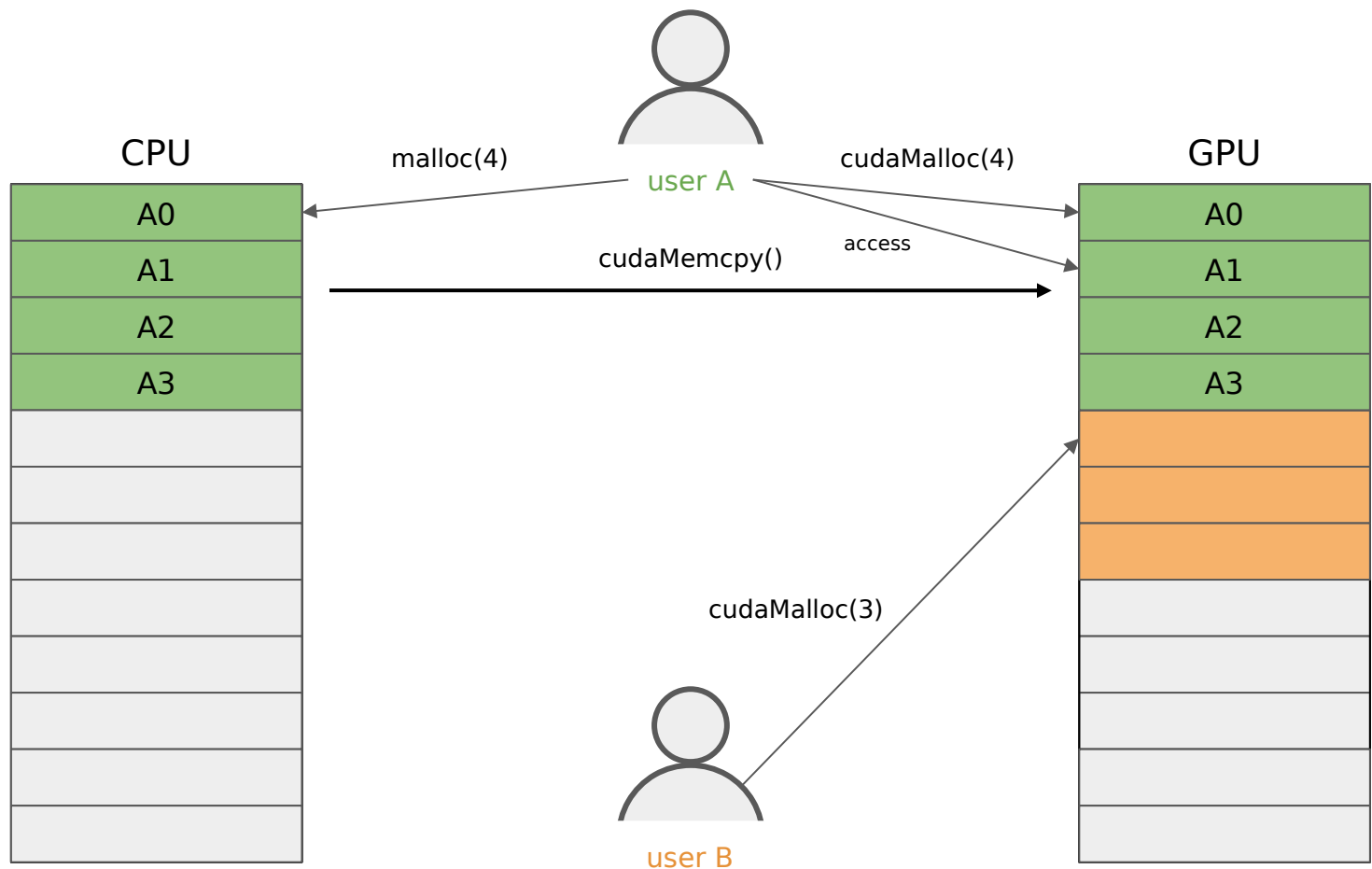
# Default CUDA Memory Model



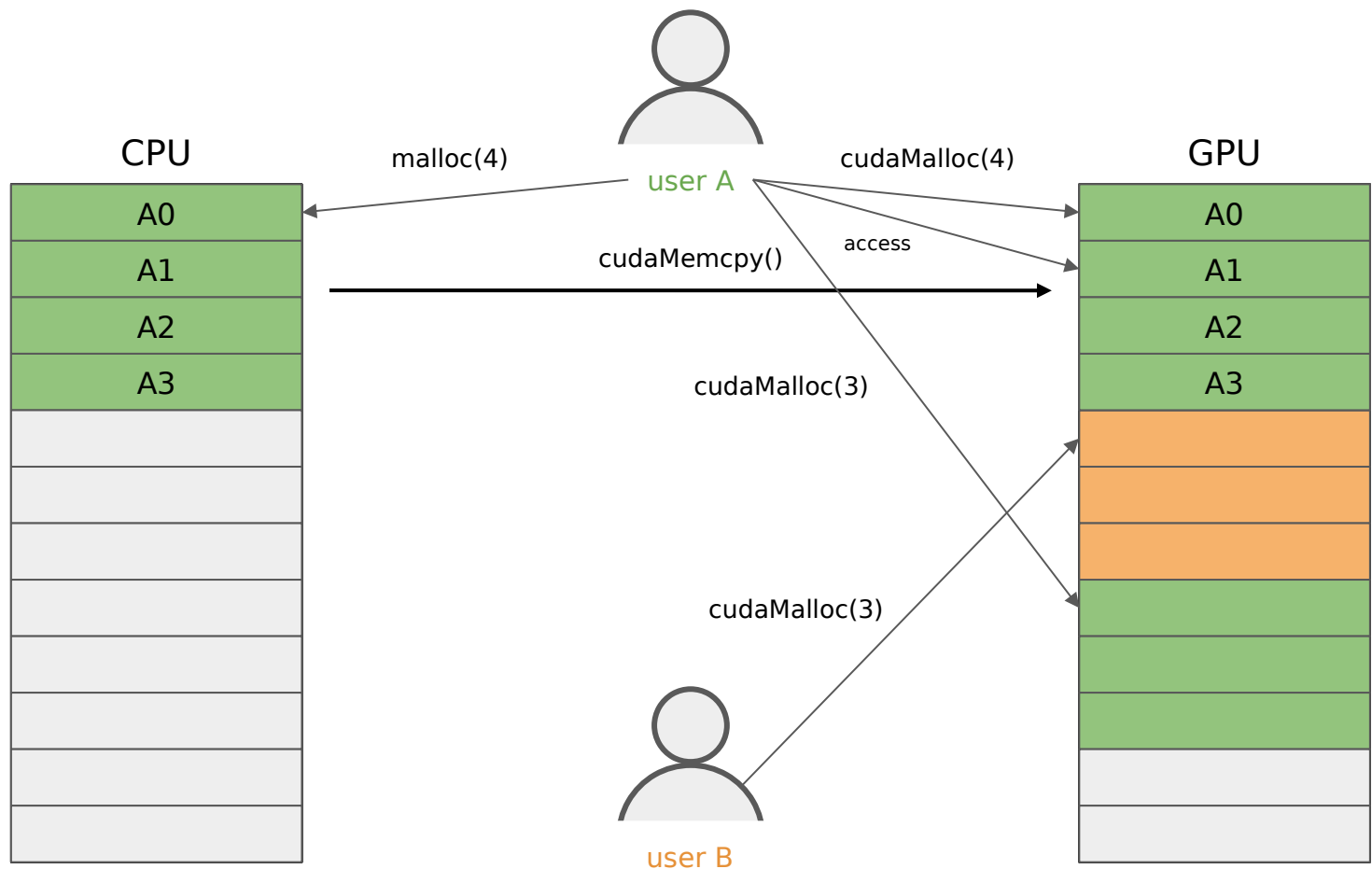
# Default CUDA Memory Model



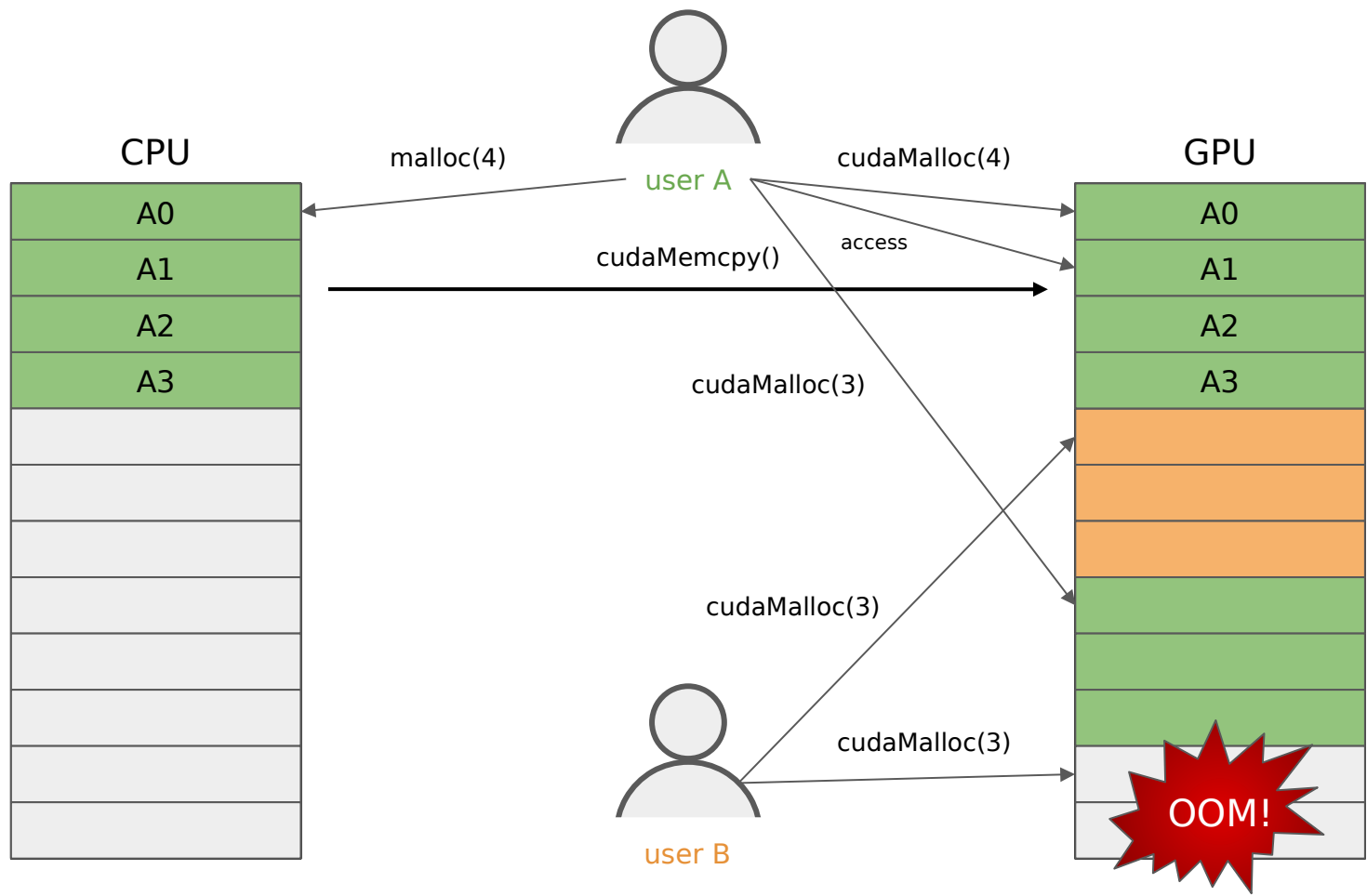
# Default CUDA Memory Model



# Default CUDA Memory Model



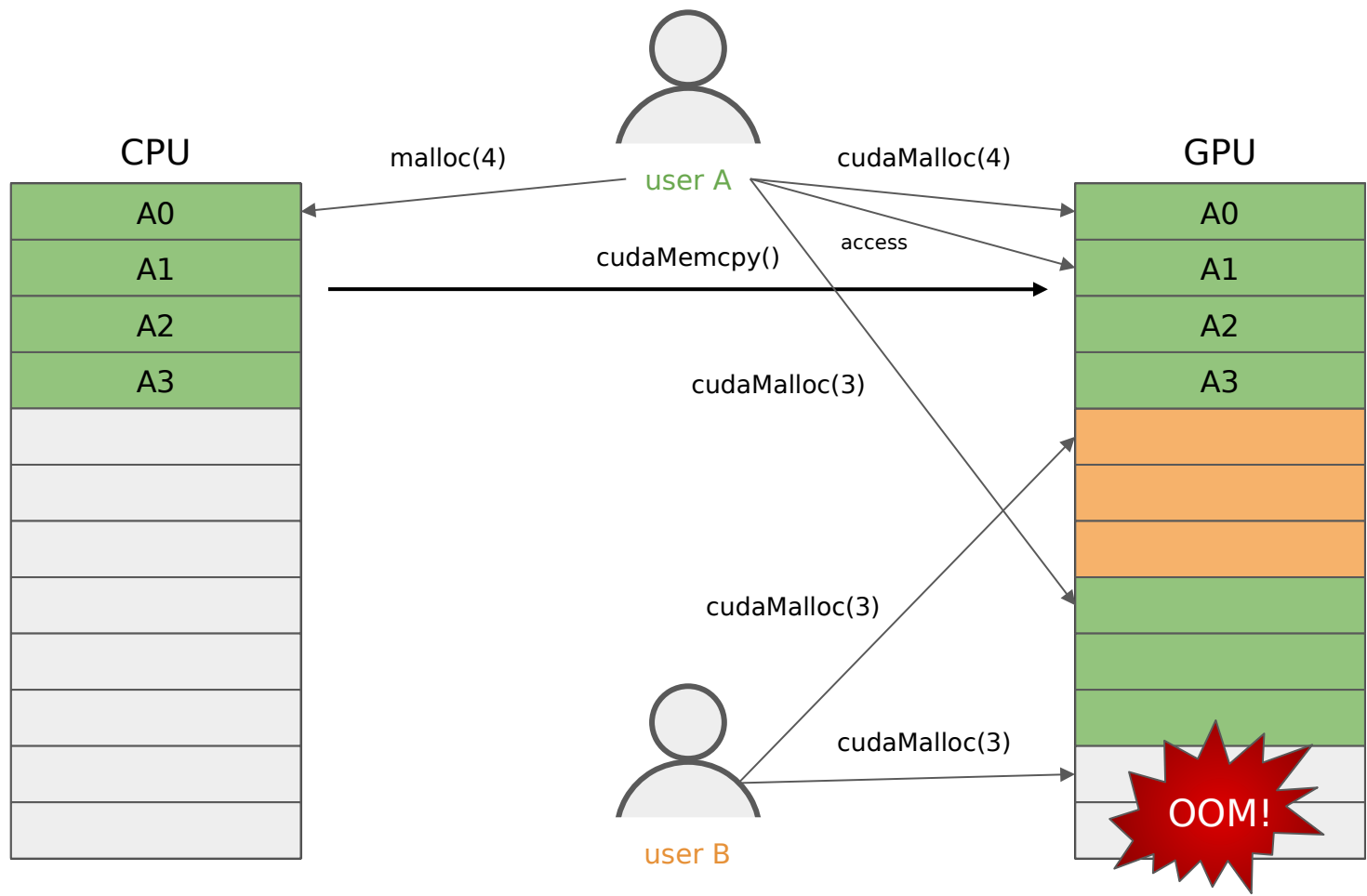
# Default CUDA Memory Model





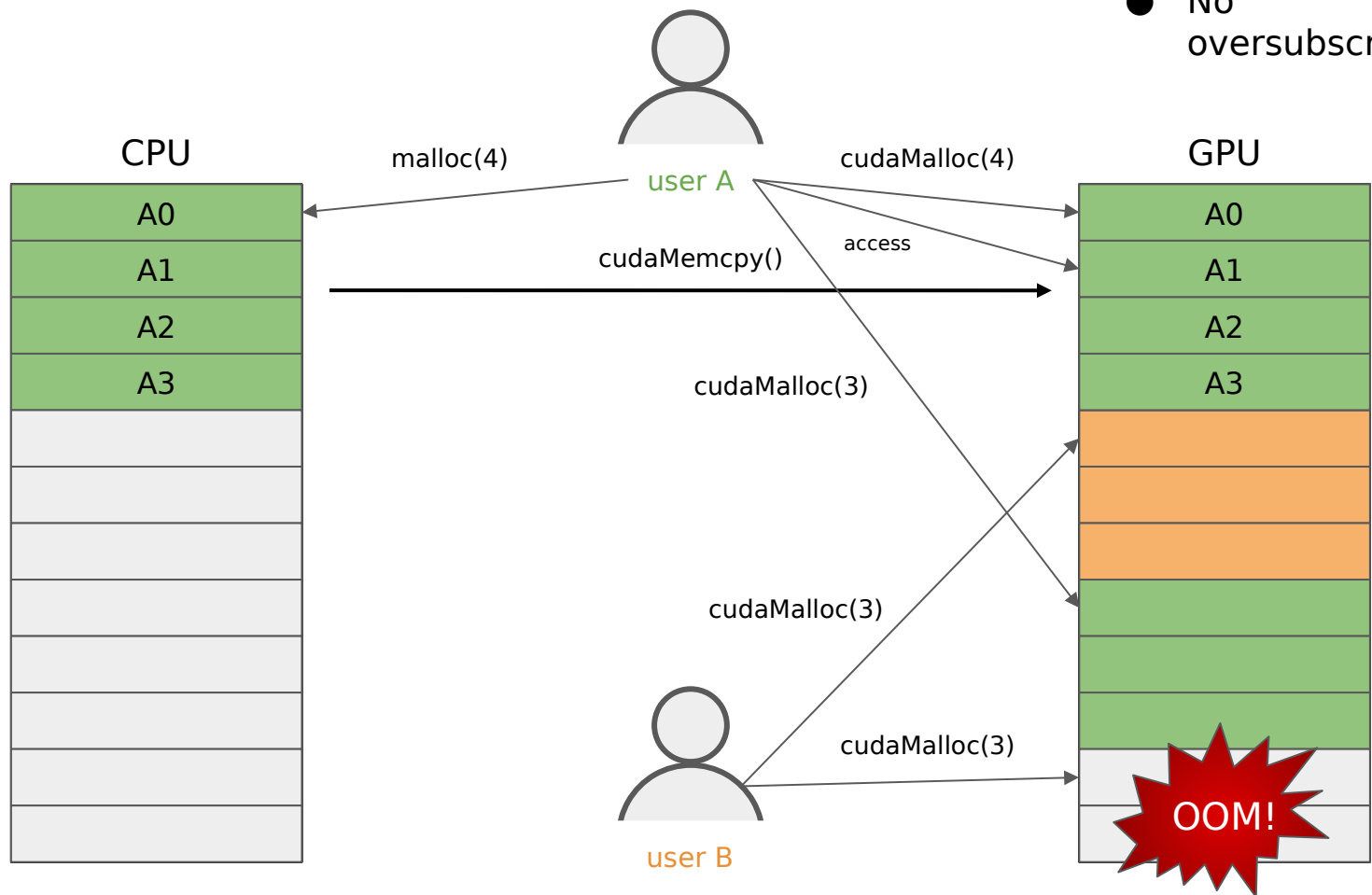
# Default CUDA Memory Model

● FCFS



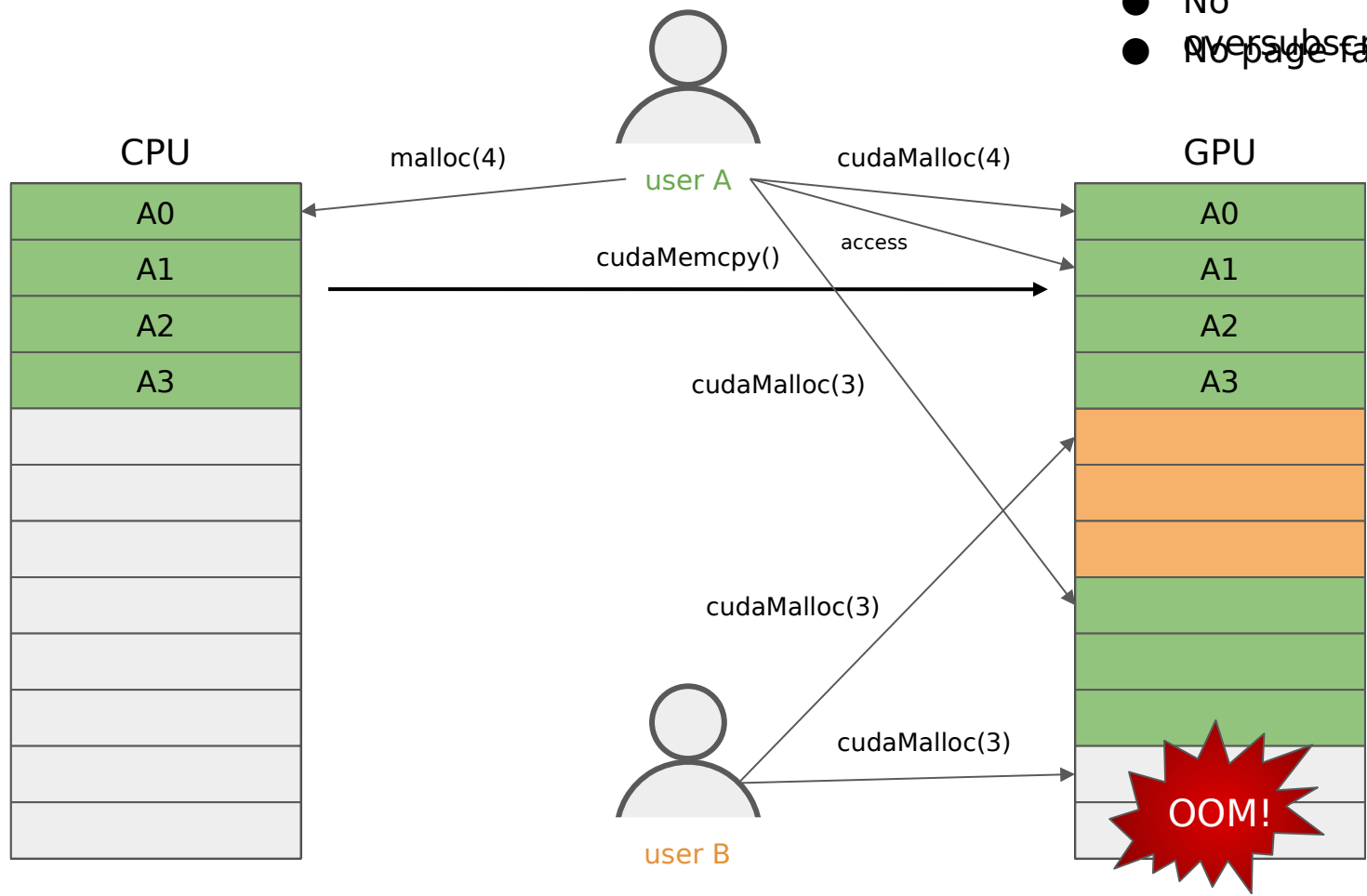
# Default CUDA Memory Model

- FCFS
- No oversubscription



# Default CUDA Memory Model

- FCFS
- No
- oversubscription
- No page faults

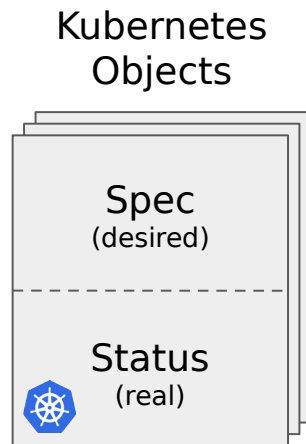


# Controller Pattern

Used ***everywhere*** in Kubernetes

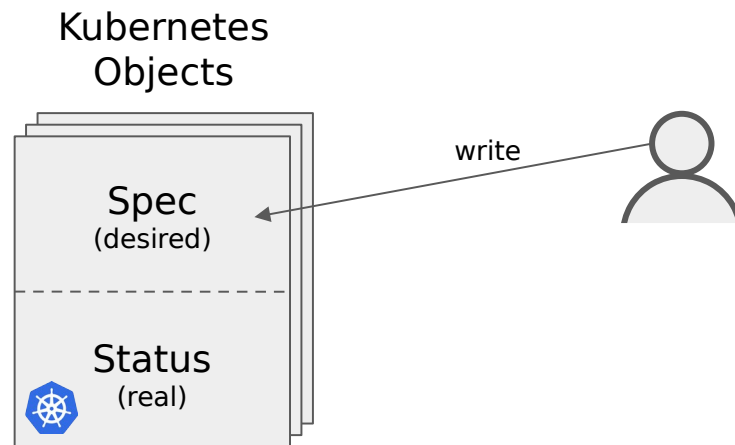
# Controller Pattern

Used ***everywhere*** in Kubernetes



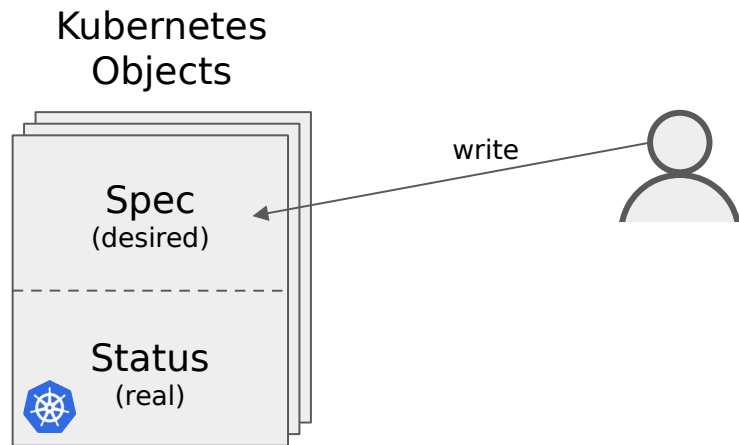
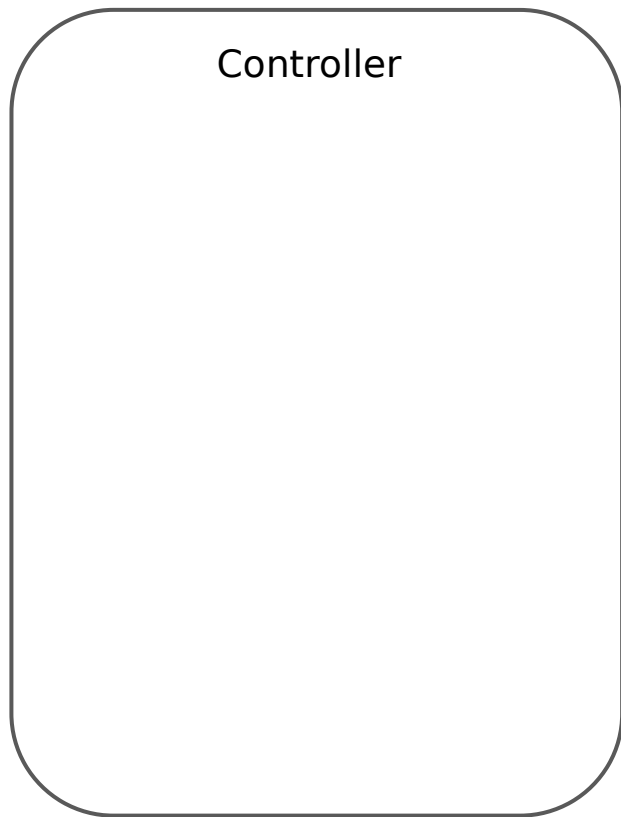
# Controller Pattern

Used ***everywhere*** in Kubernetes



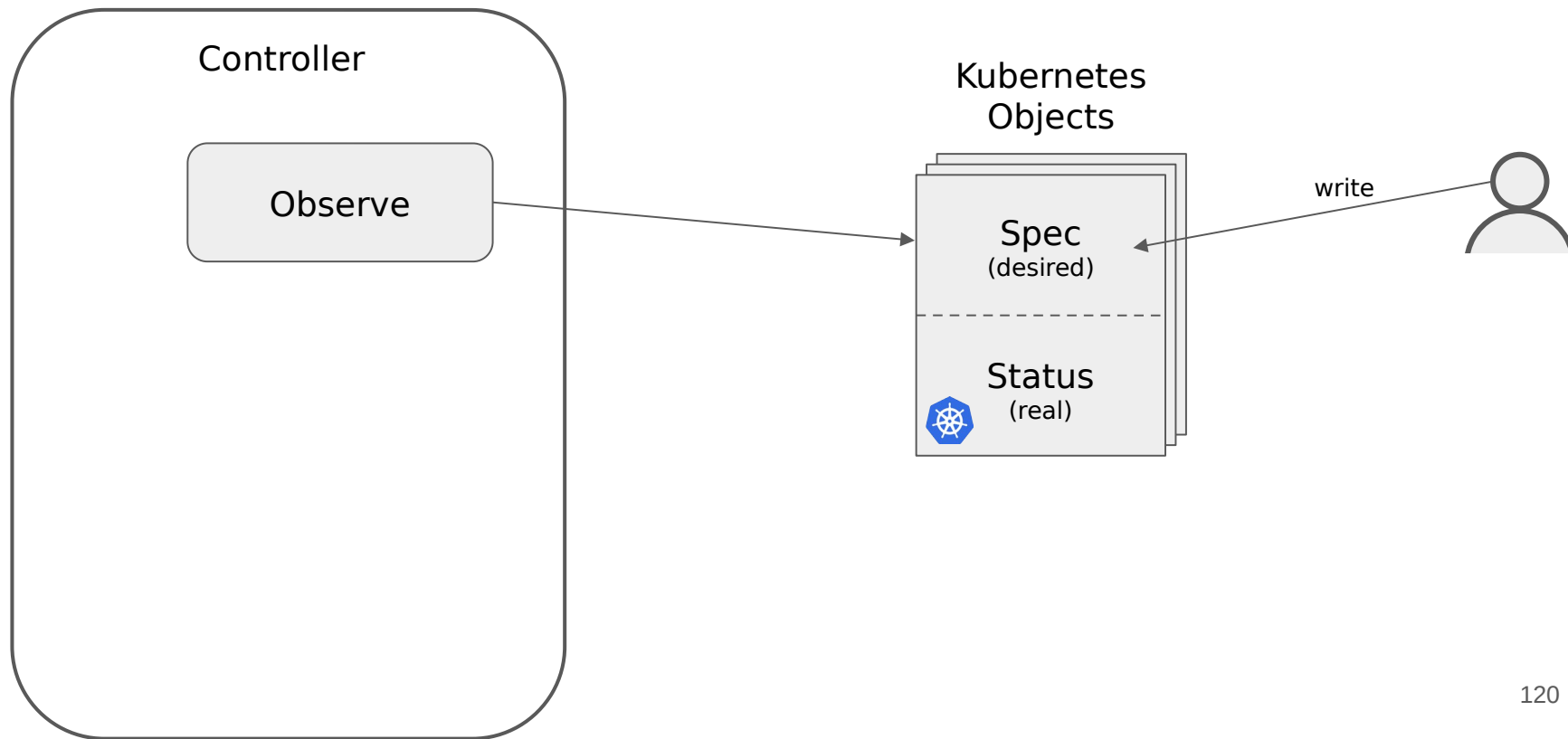
# Controller Pattern

Used ***everywhere*** in Kubernetes



# Controller Pattern

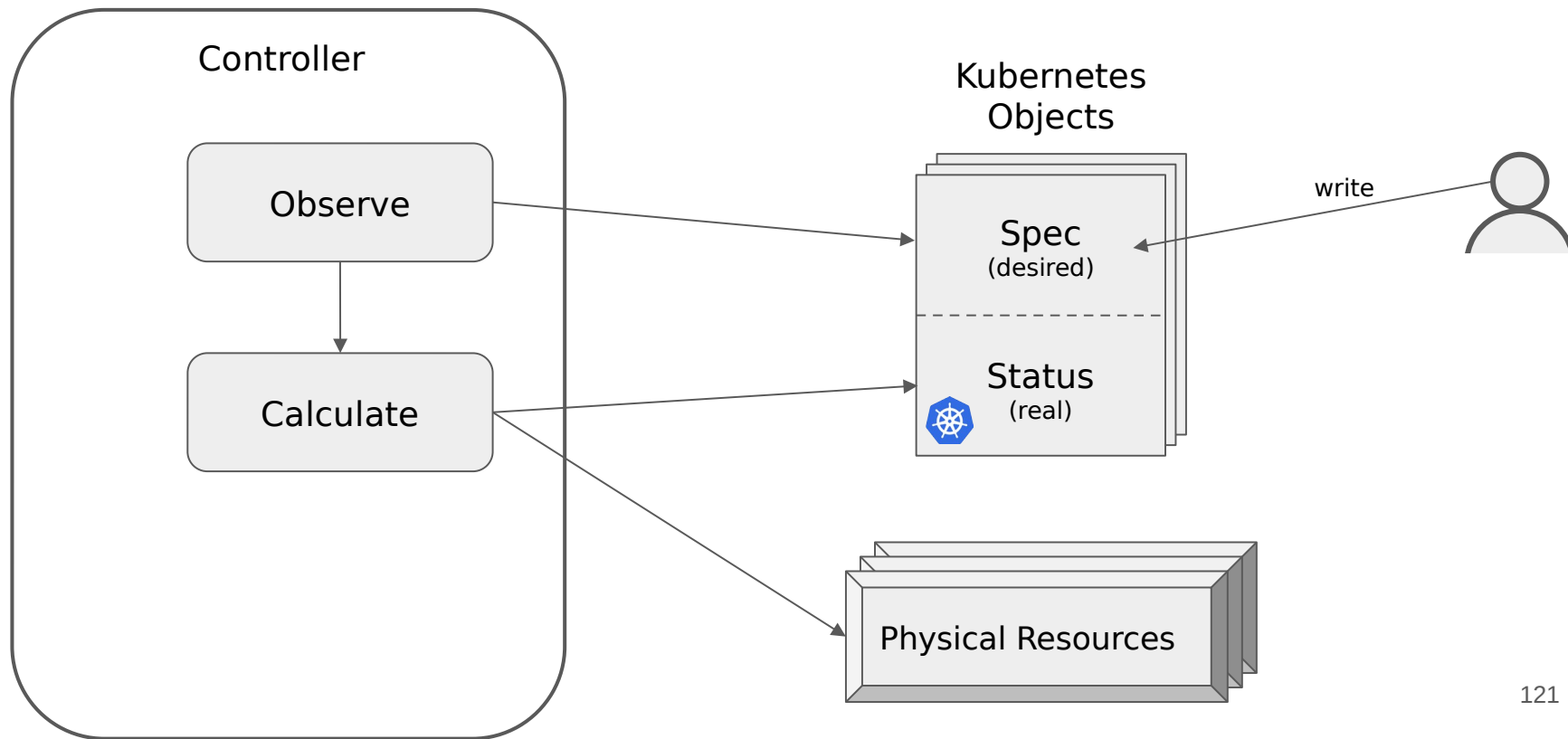
Used ***everywhere*** in Kubernetes





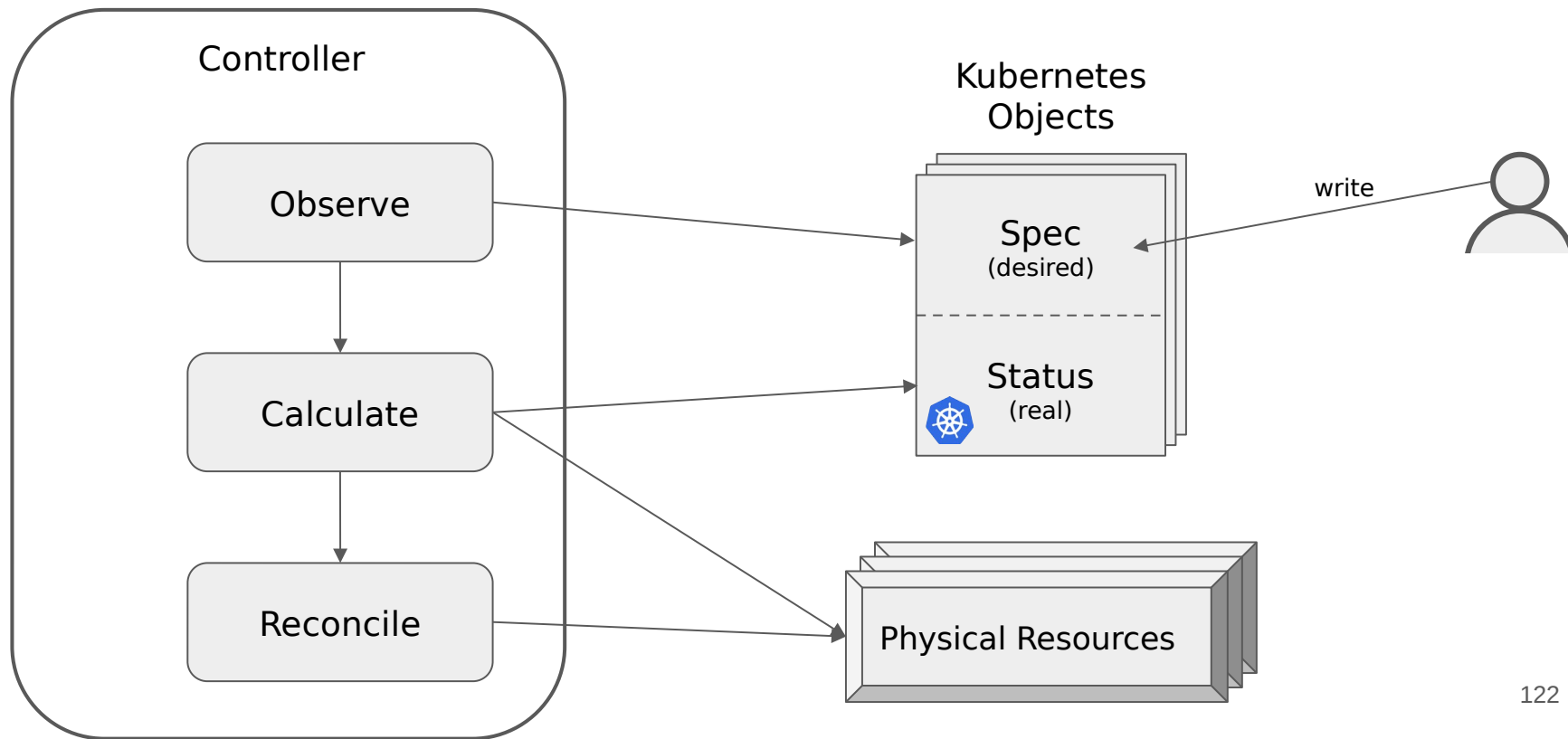
# Controller Pattern

Used ***everywhere*** in Kubernetes



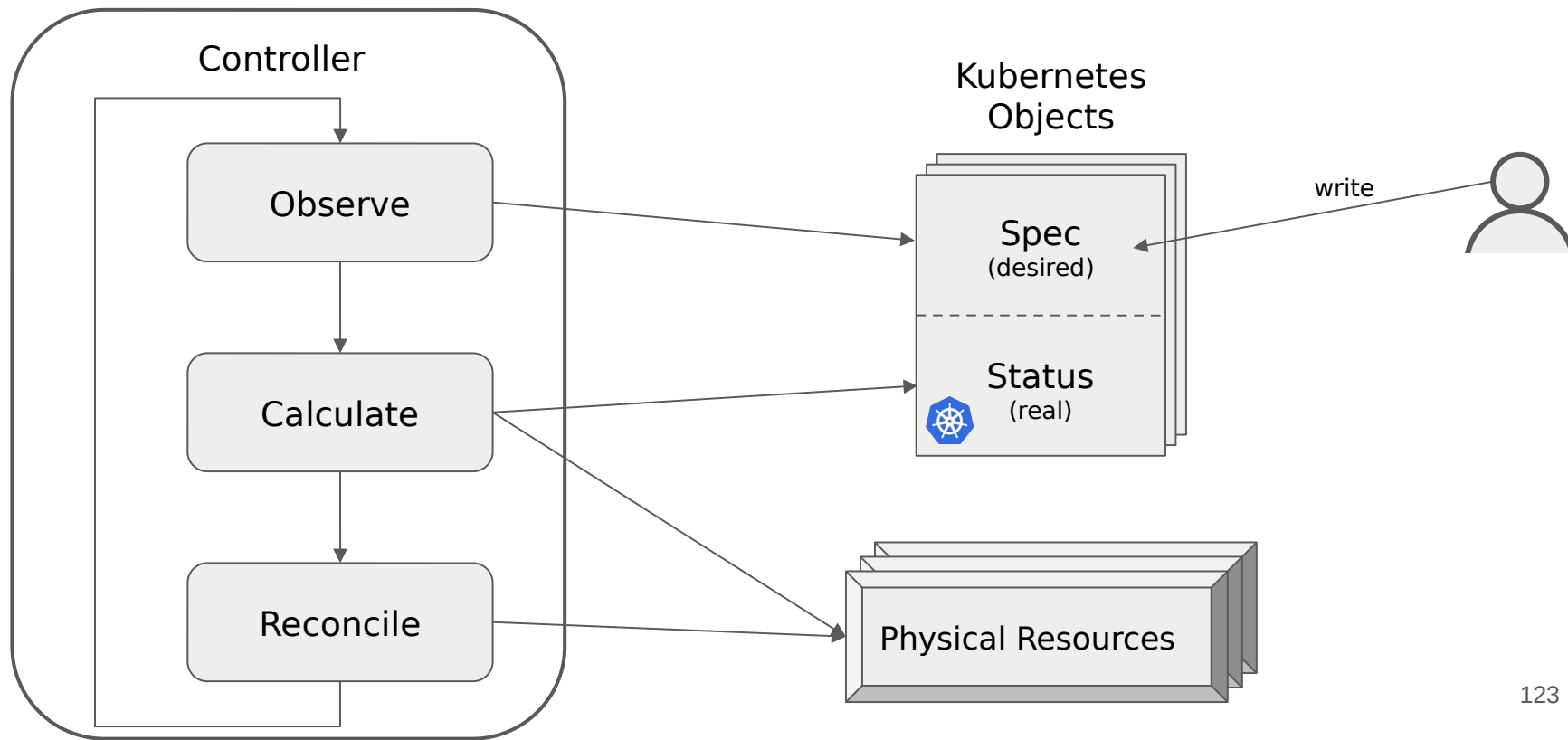
# Controller Pattern

Used ***everywhere*** in Kubernetes

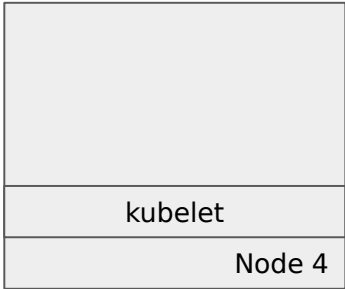
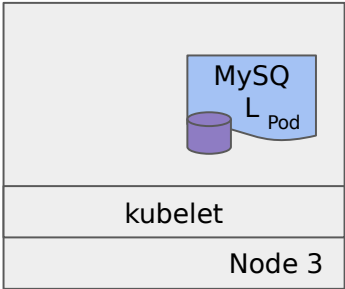
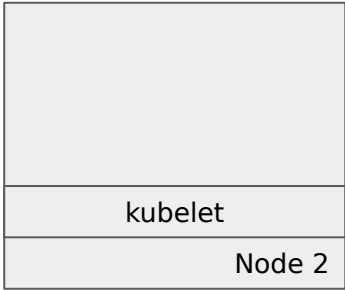
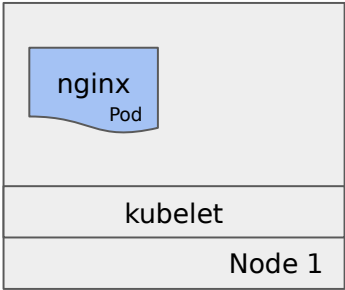
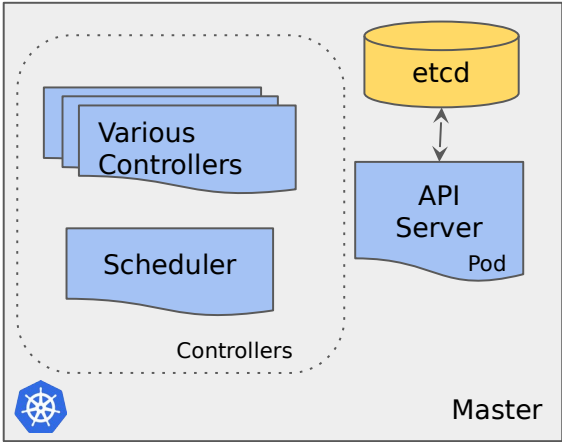


# Controller Pattern

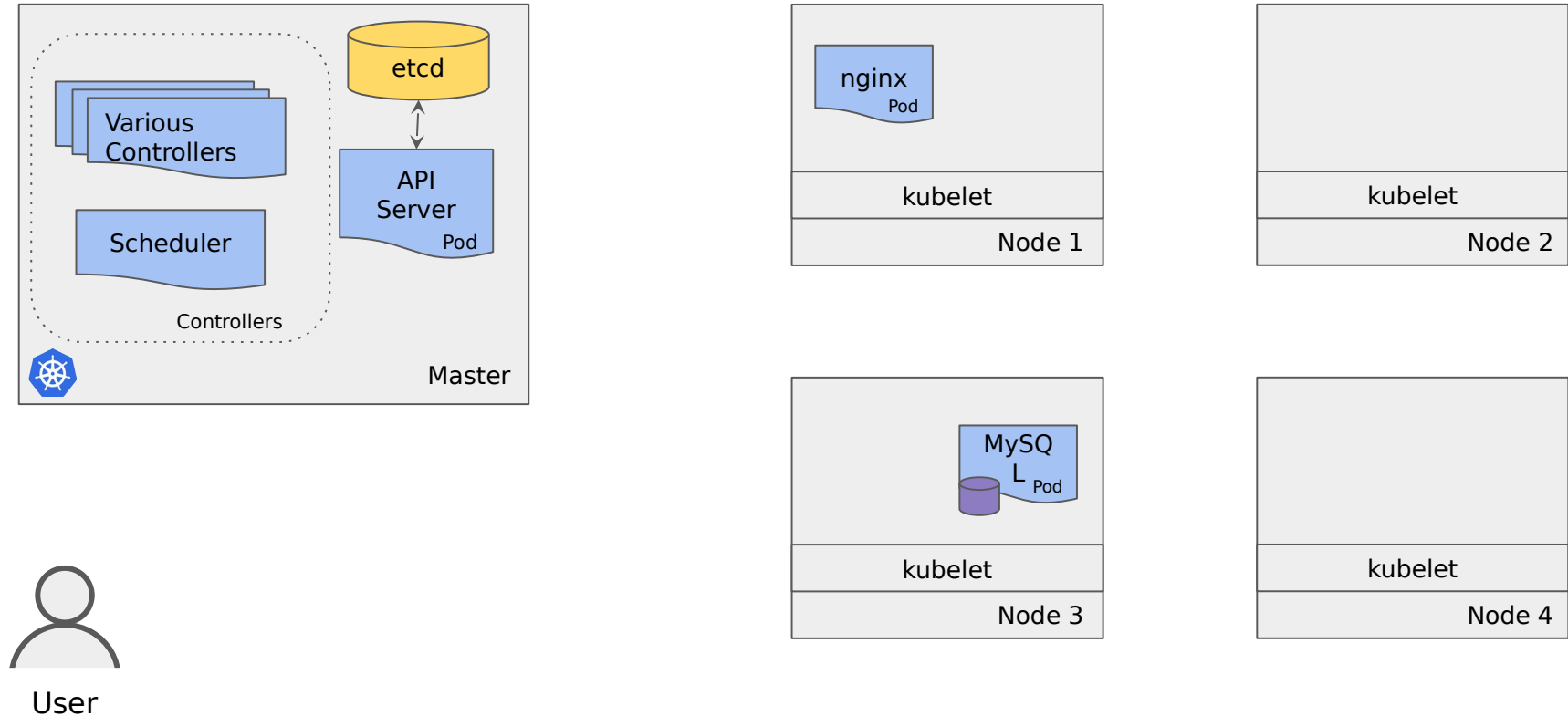
Used ***everywhere*** in Kubernetes



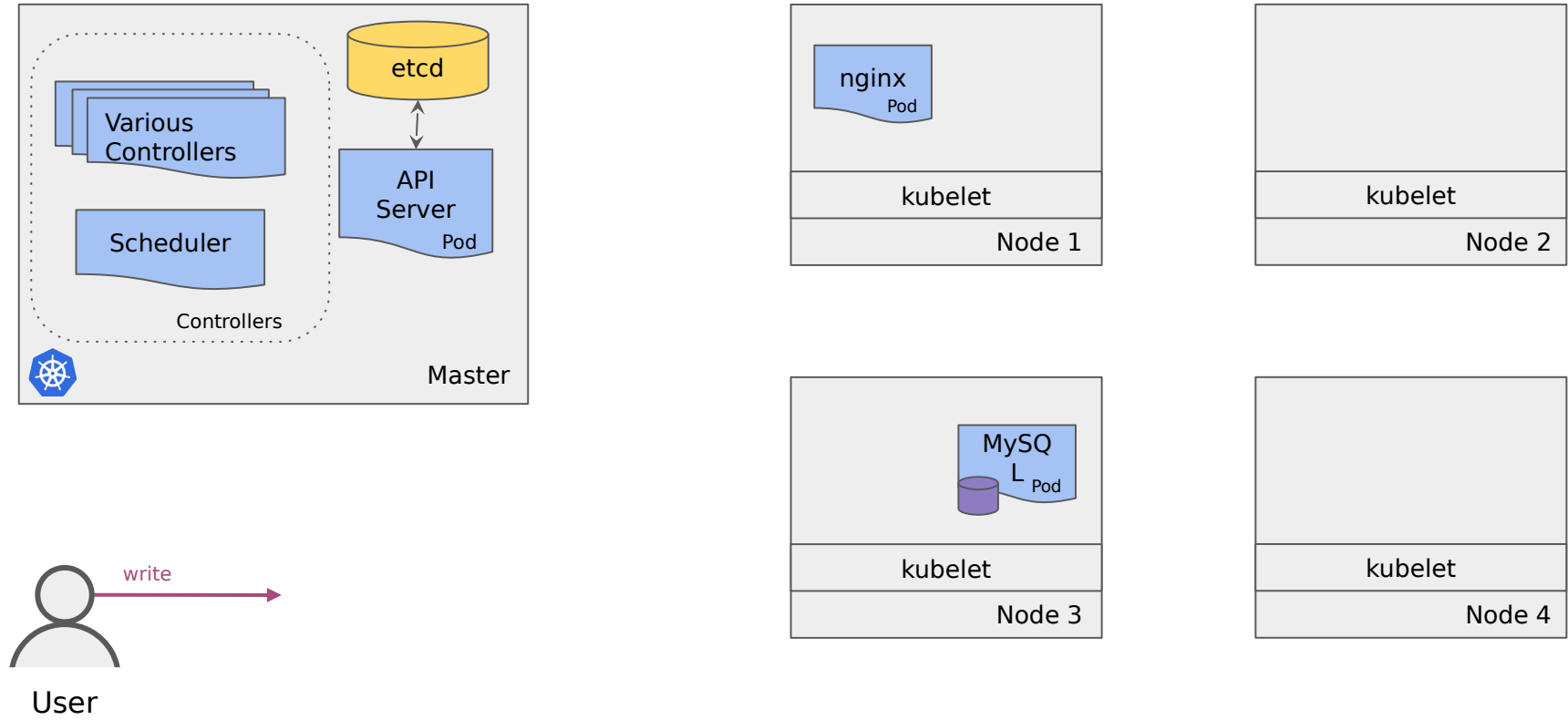
# Kubernetes Scheduling



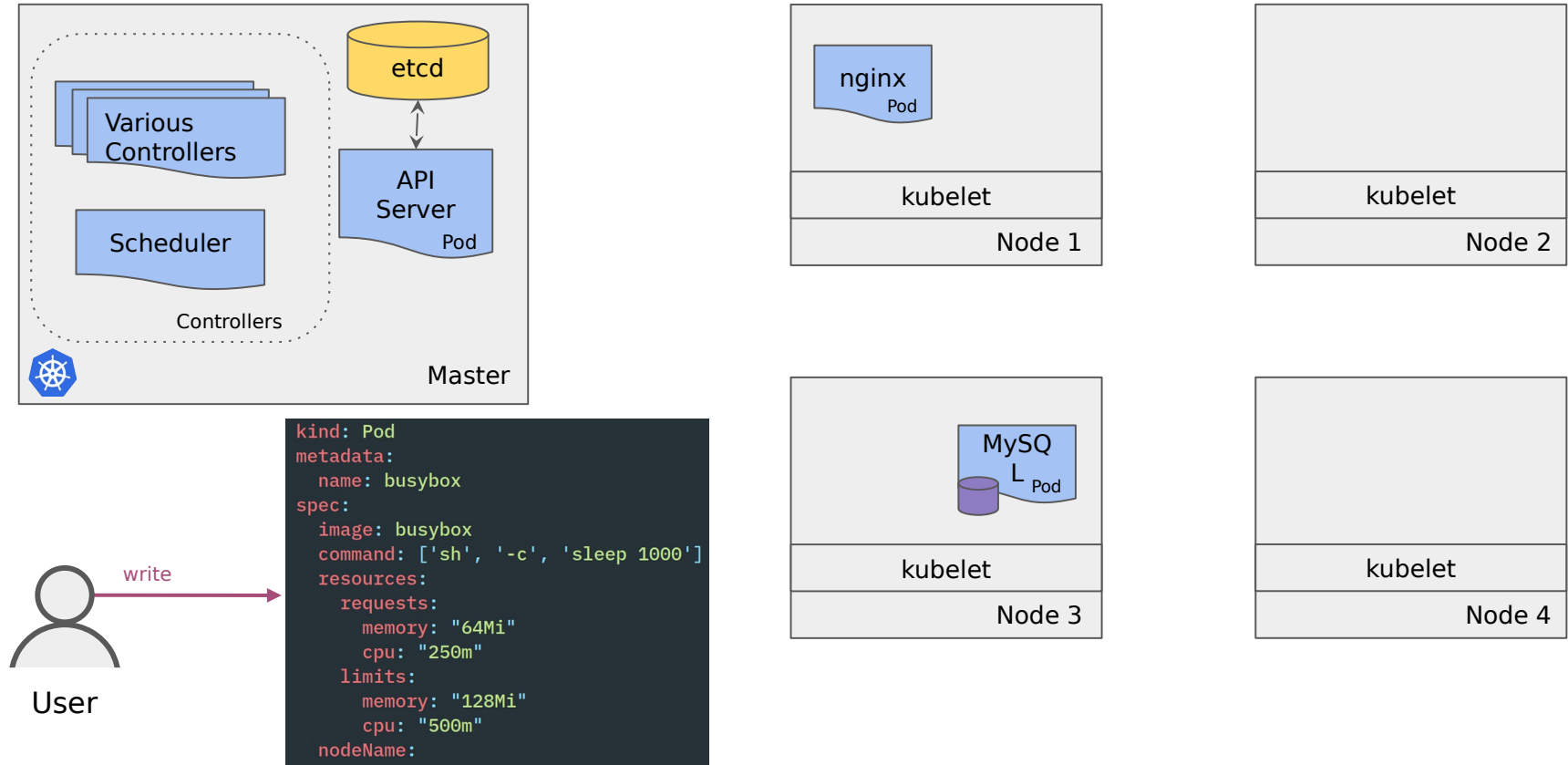
# Kubernetes Scheduling



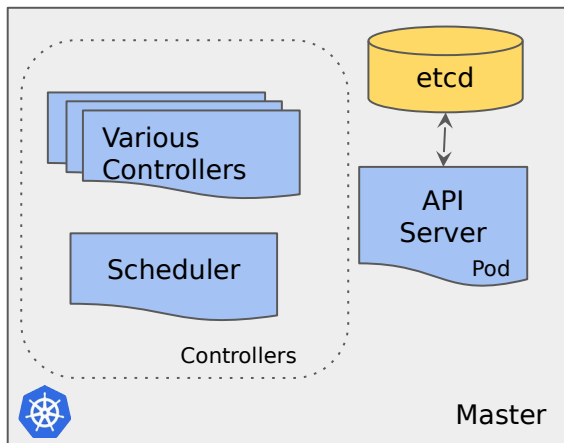
# Kubernetes Scheduling



# Kubernetes Scheduling



# Kubernetes Scheduling

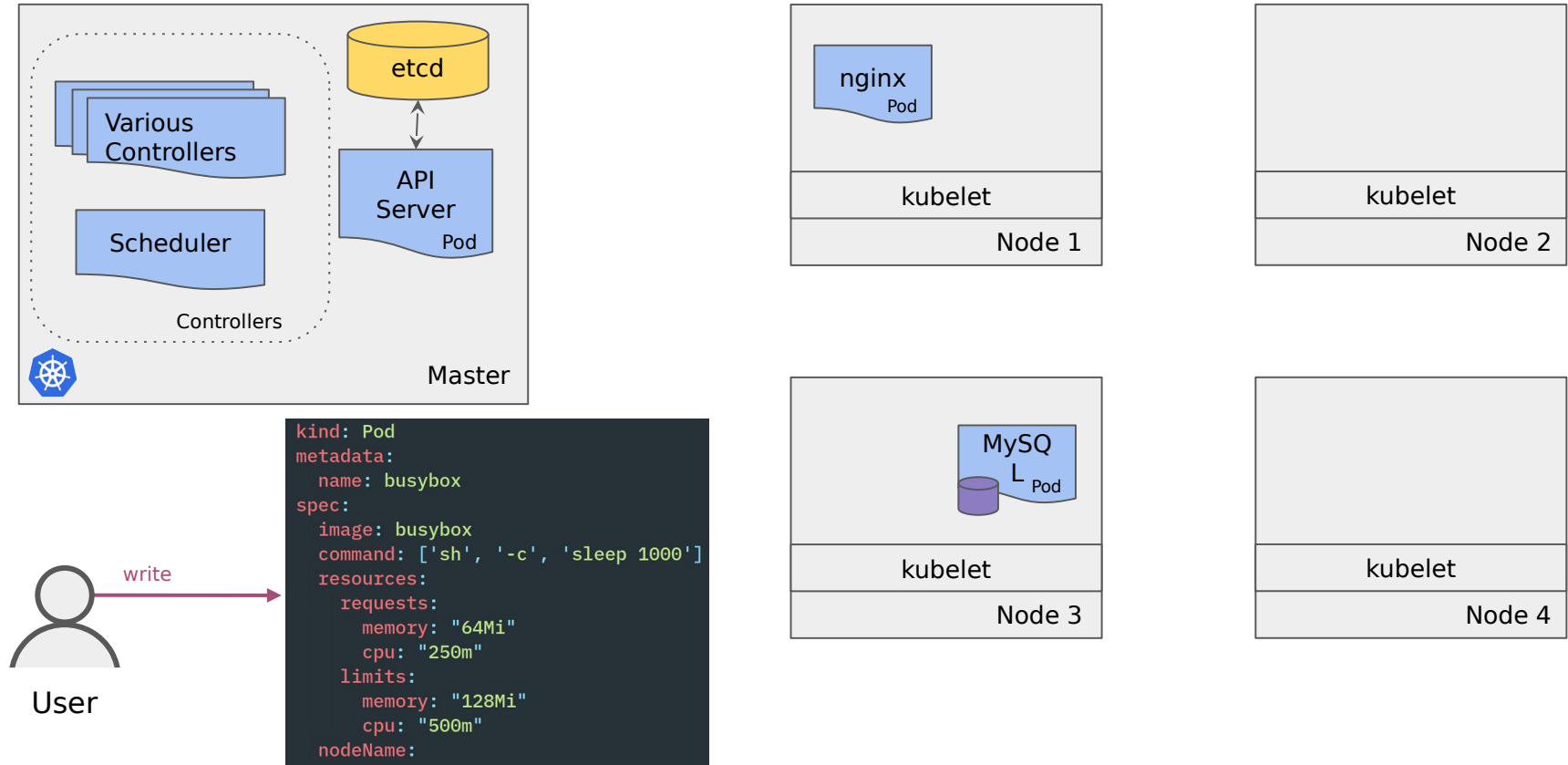


```
kind: Pod
metadata:
  name: busybox
spec:
  image: busybox
  command: ['sh', '-c', 'sleep 1000']
  resources:
    requests:
      memory: "64Mi"
      cpu: "250m"
    limits:
      memory: "128Mi"
      cpu: "500m"
  nodeName:
```

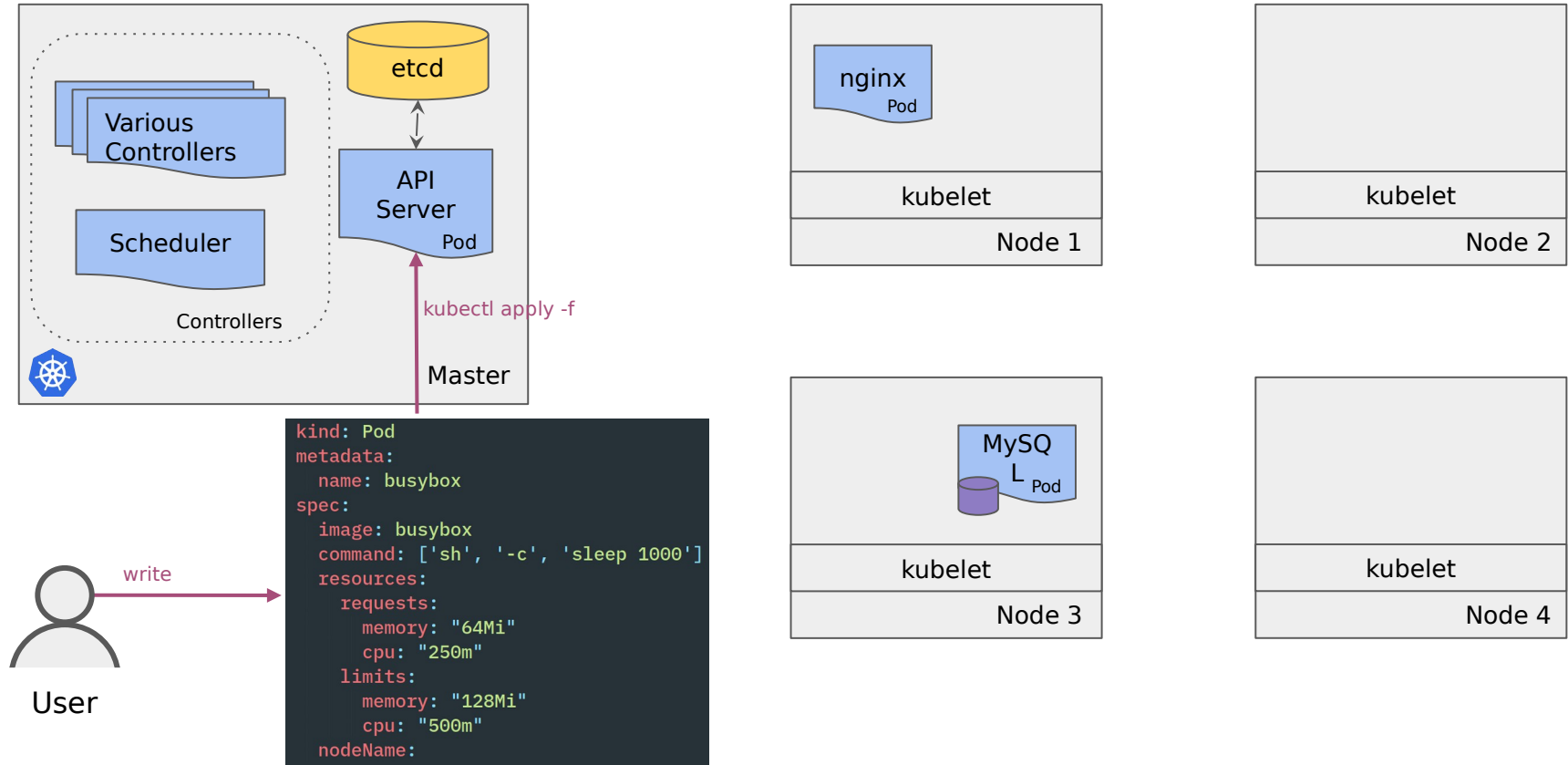
```
kind: Pod
metadata:
  name: busybox
spec:
  image: busybox
  command: ['sh', '-c', 'sleep 1000']
  resources:
    requests:
      memory: "64Mi"
      cpu: "250m"
    limits:
      memory: "128Mi"
      cpu: "500m"
  nodeName:
```



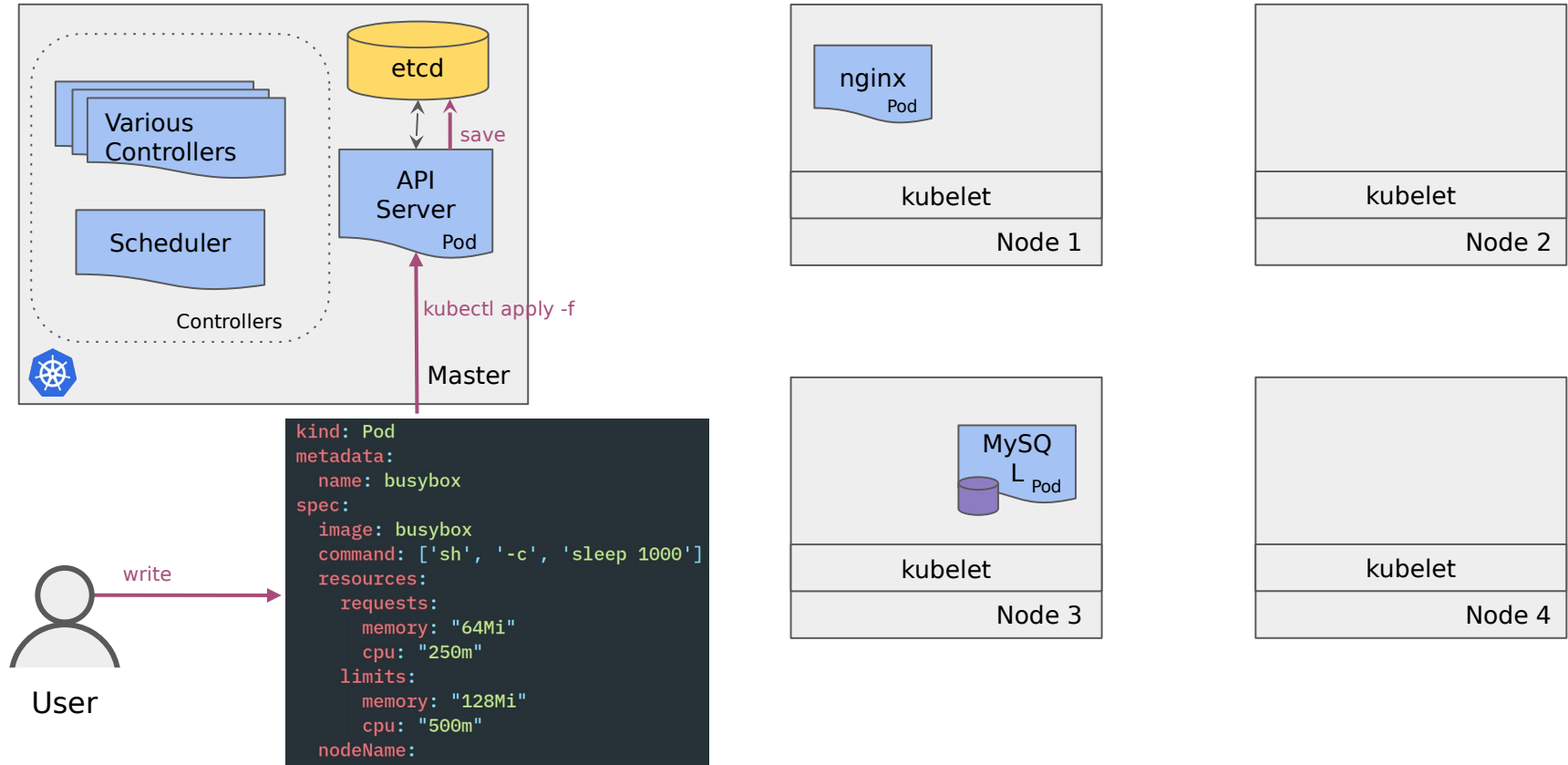
# Kubernetes Scheduling



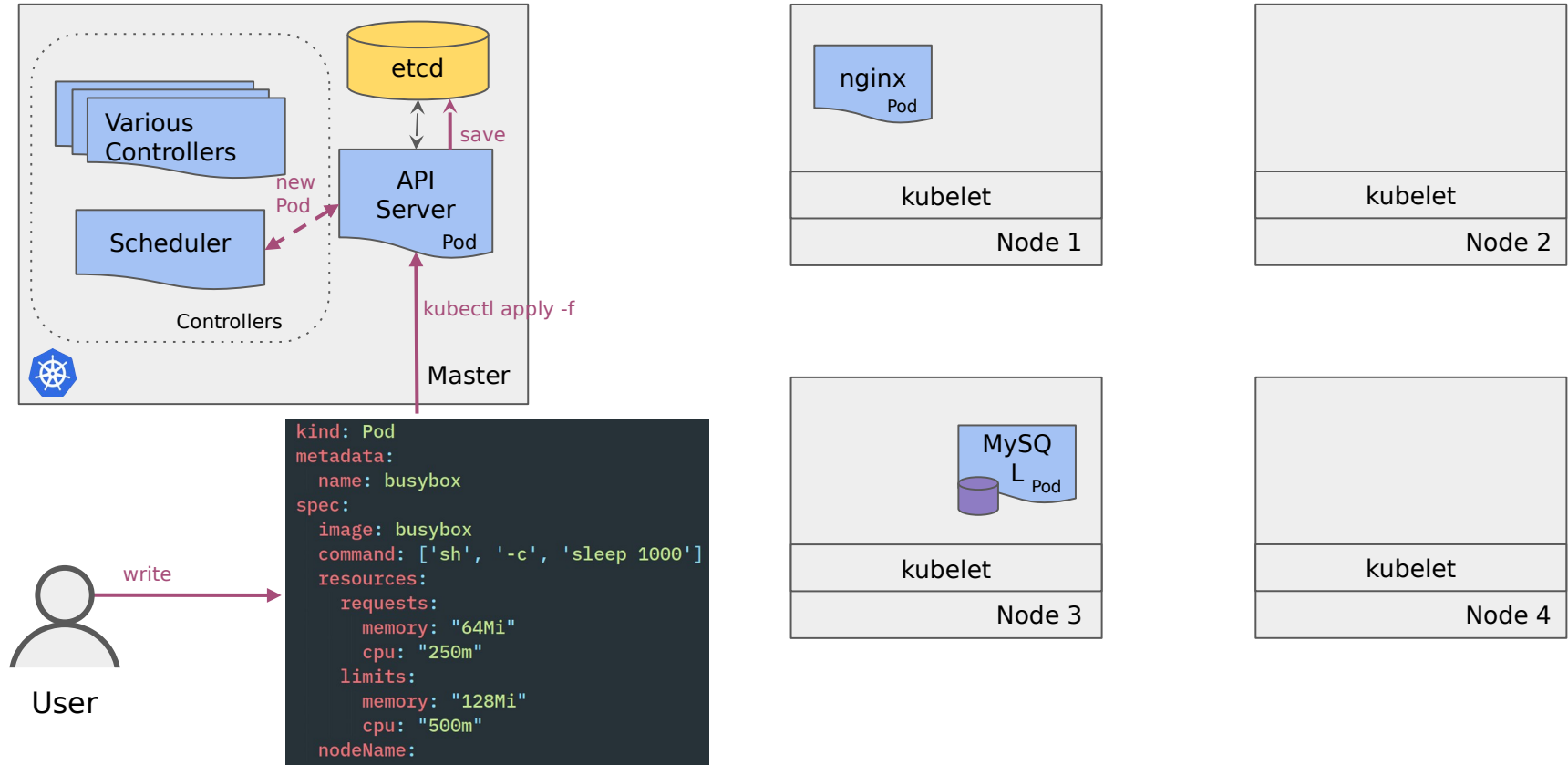
# Kubernetes Scheduling



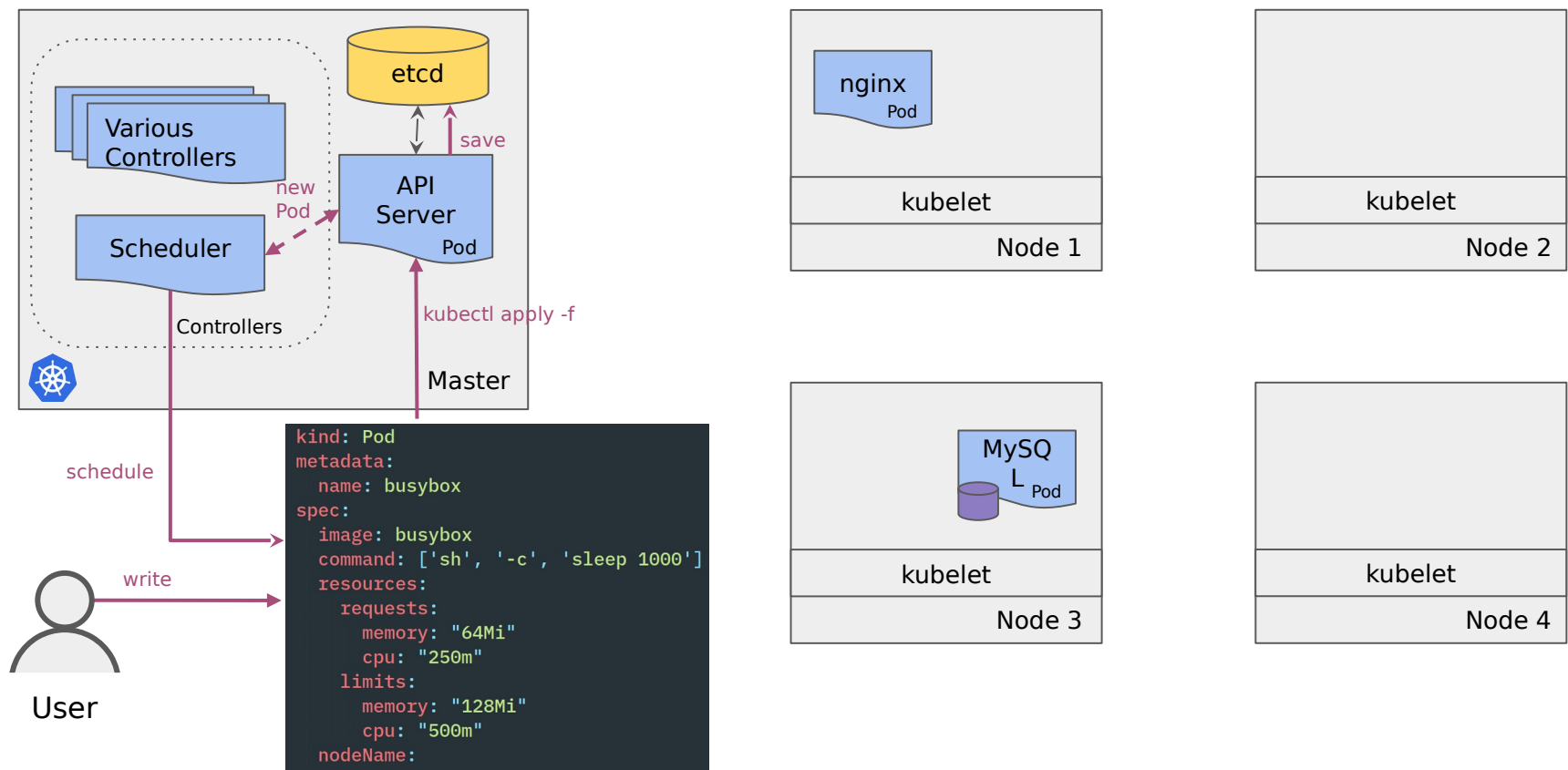
# Kubernetes Scheduling



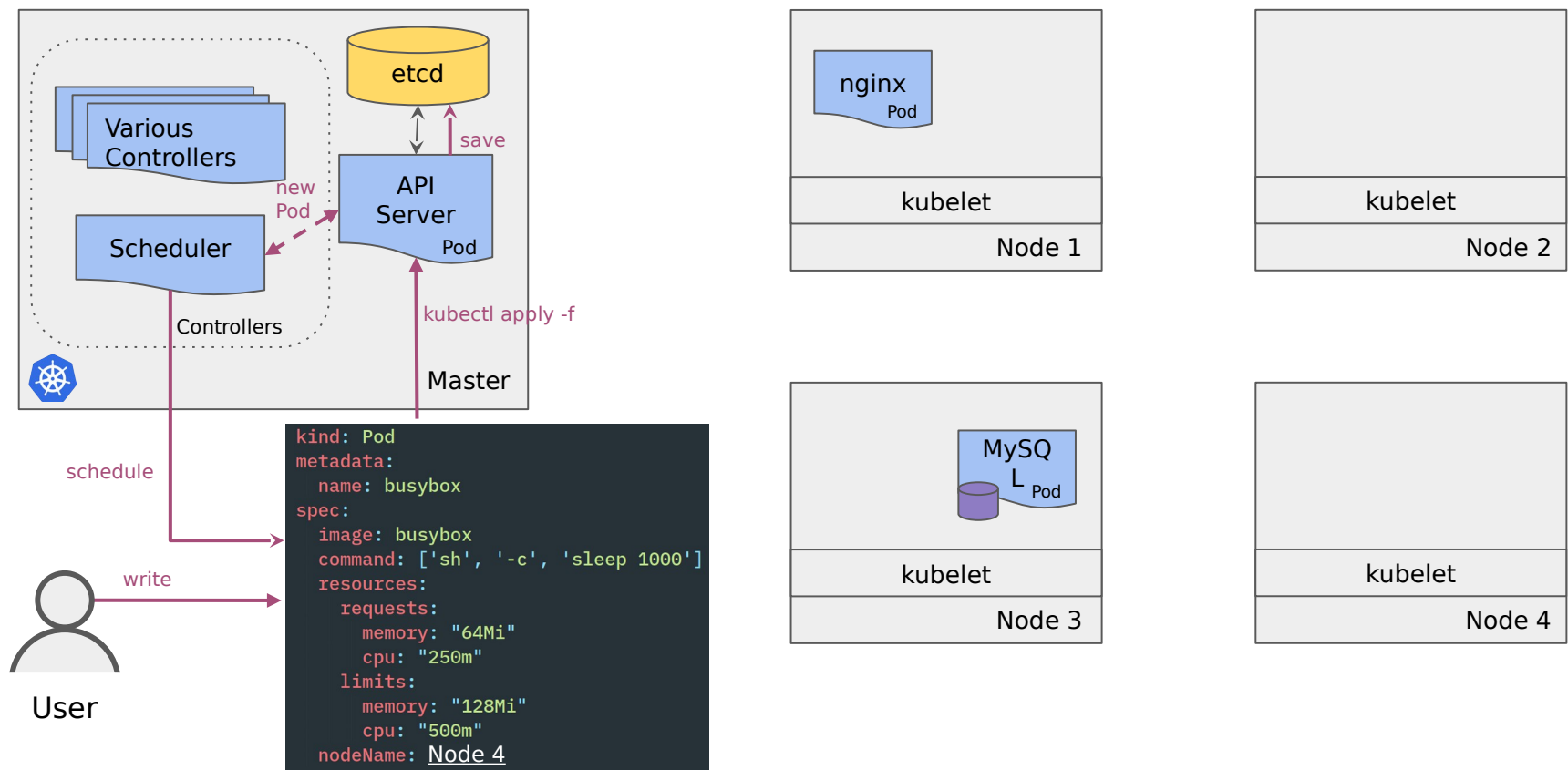
# Kubernetes Scheduling



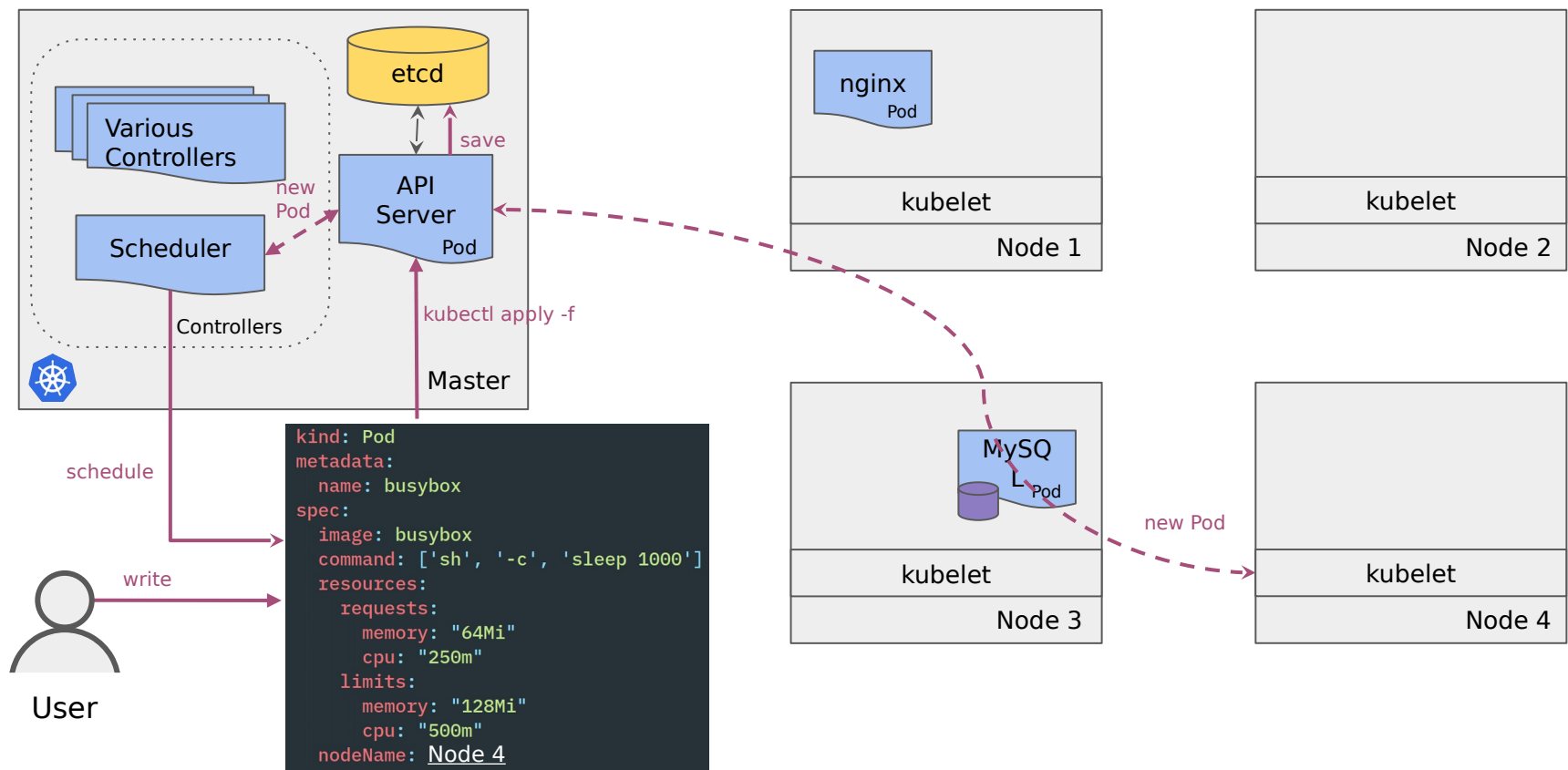
# Kubernetes Scheduling



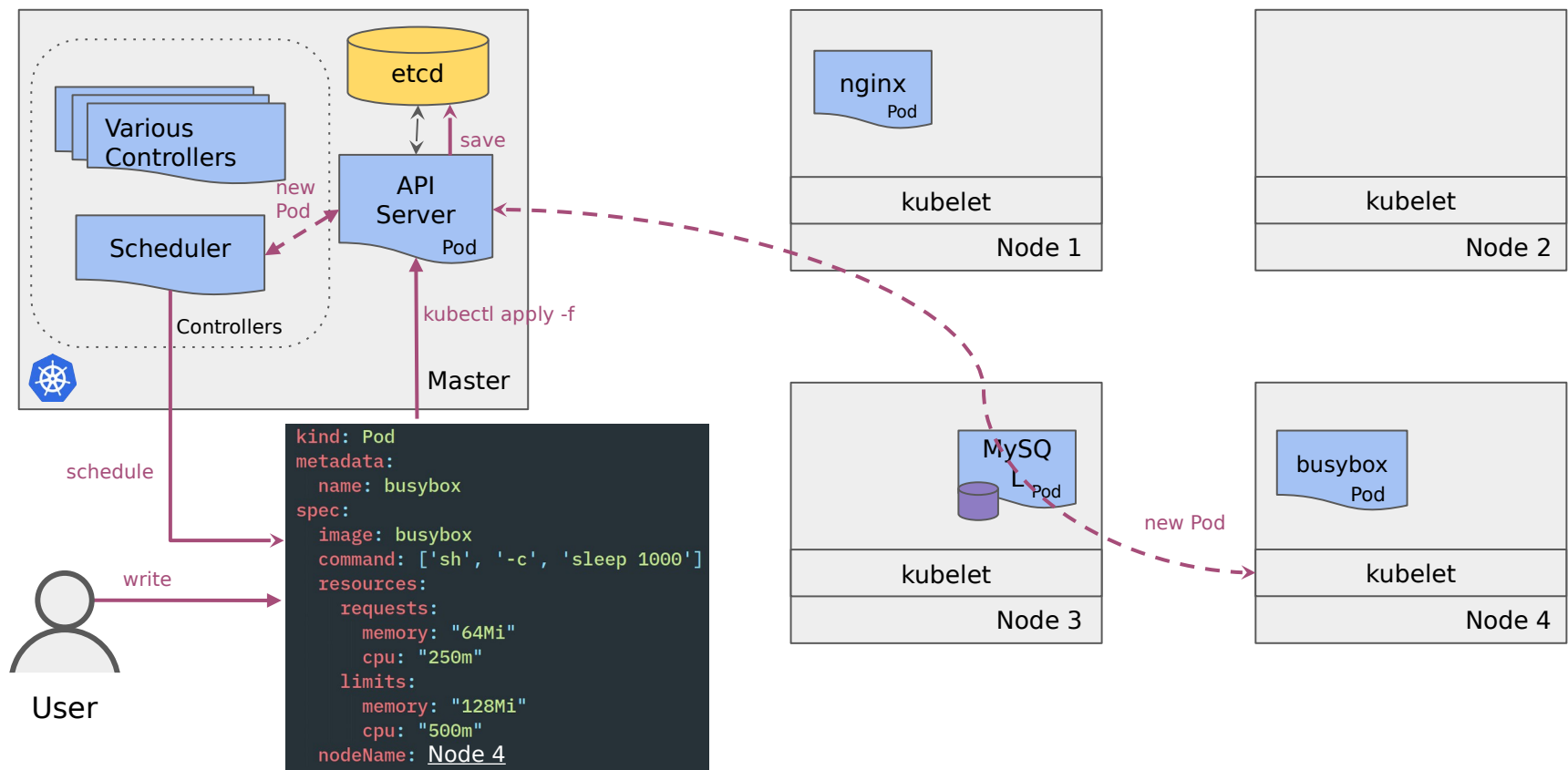
# Kubernetes Scheduling



# Kubernetes Scheduling

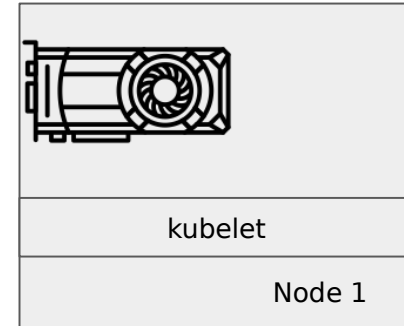
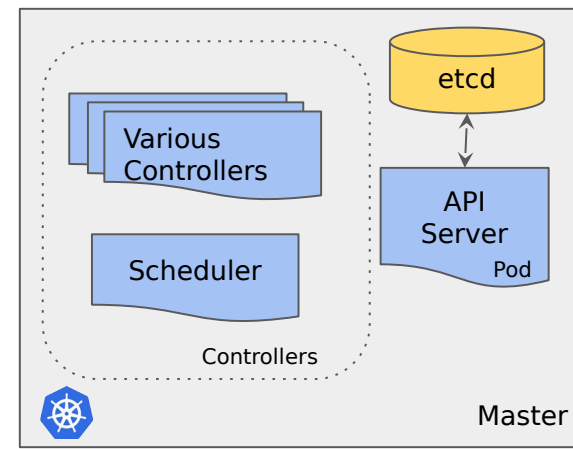


# Kubernetes Scheduling

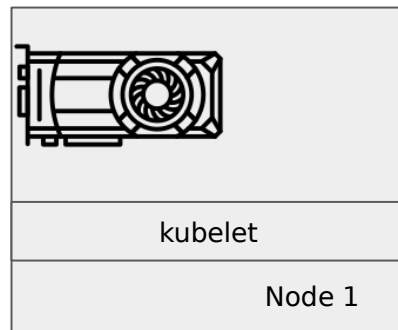
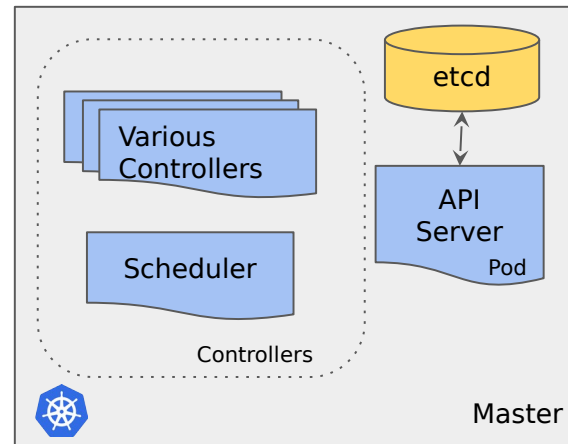




# The Problem (once more)



# The Problem (once more)



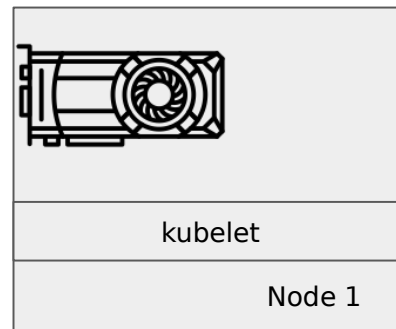
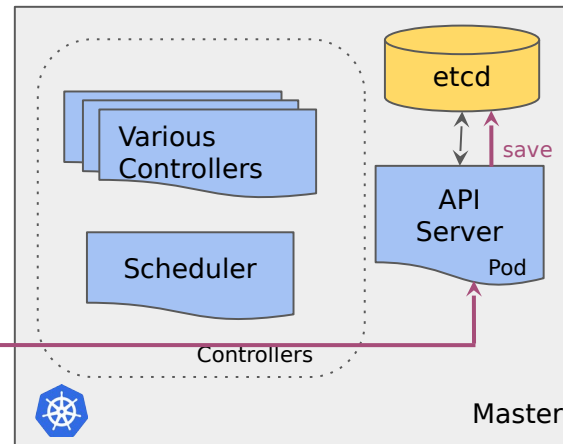
# The Problem (once more)



write

```
kind: Pod
metadata:
  name: jupyter-1
spec:
  image: jupyter:latest
  resources:
    limits:
      nvidia.com/gpu: 1
  nodeName:
```

kubectl apply -f



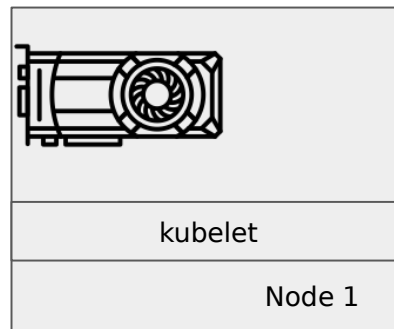
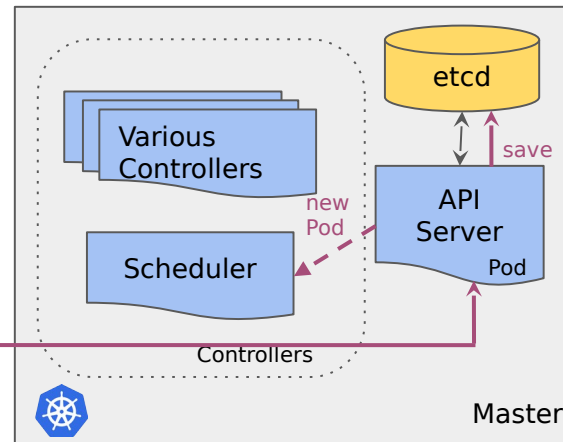
# The Problem (once more)



write

```
kind: Pod
metadata:
  name: jupyter-1
spec:
  image: jupyter:latest
  resources:
    limits:
      nvidia.com/gpu: 1
  nodeName:
```

kubectl apply -f



# The Problem (once more)

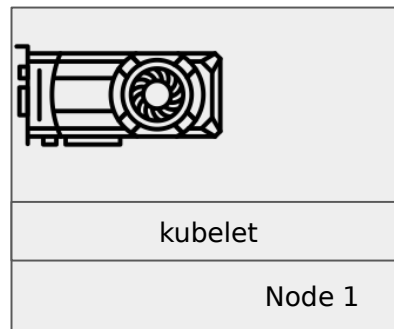
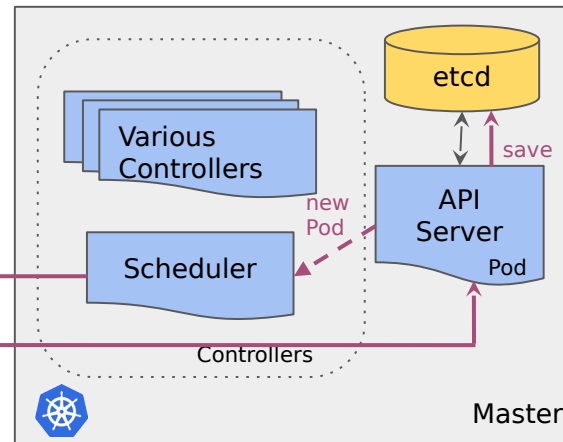


write

```
kind: Pod
metadata:
  name: jupyter-1
spec:
  image: jupyter:latest
  resources:
    limits:
      nvidia.com/gpu: 1
  nodeName:
```

schedule

kubectl apply -f



# The Problem (once more)

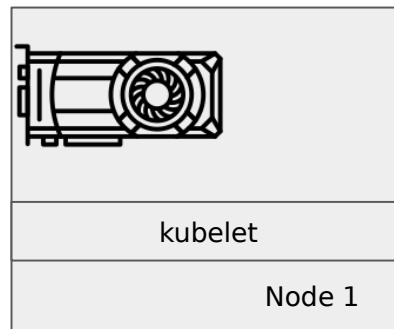
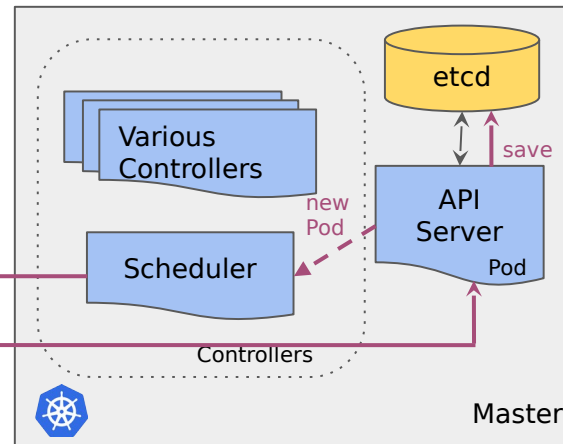


write

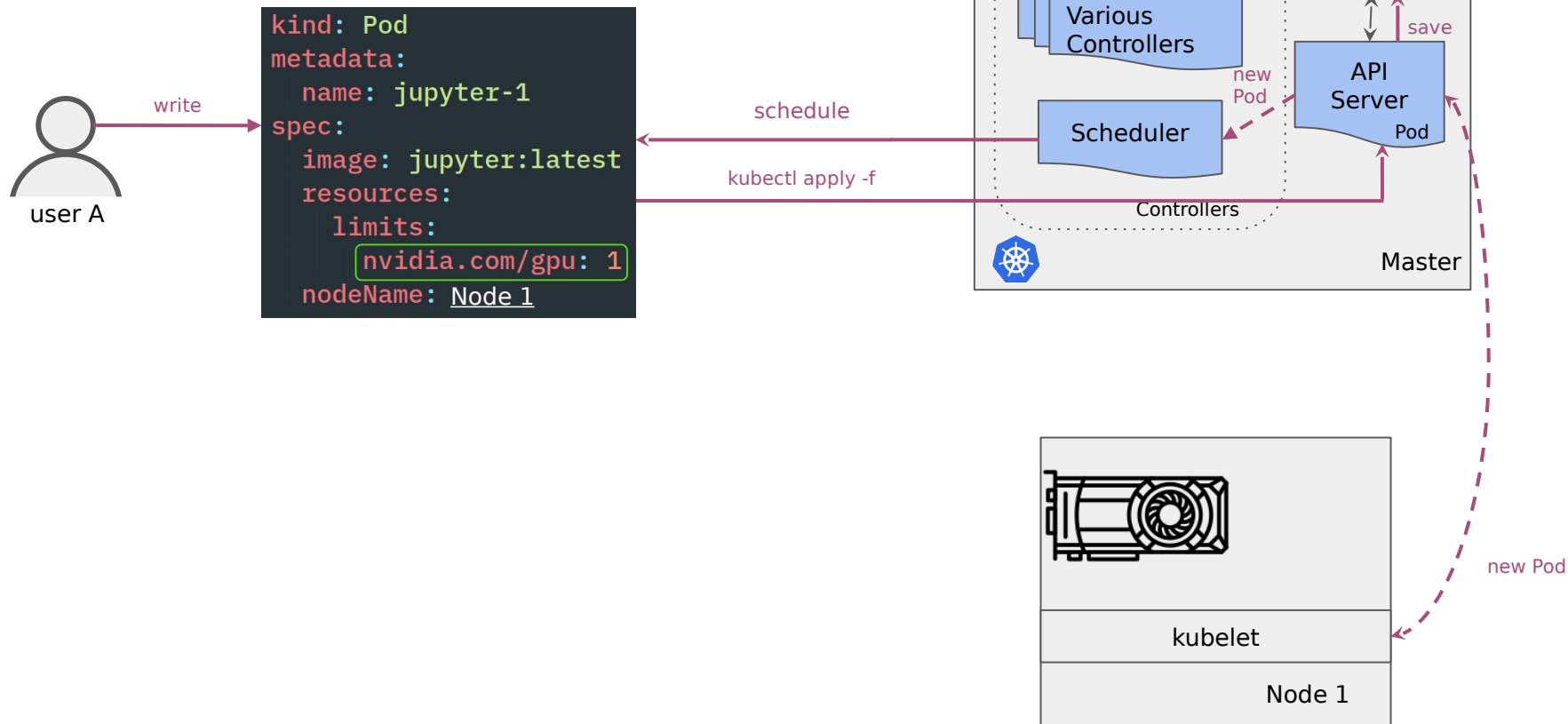
```
kind: Pod
metadata:
  name: jupyter-1
spec:
  image: jupyter:latest
  resources:
    limits:
      nvidia.com/gpu: 1
  nodeName: Node 1
```

schedule

kubectl apply -f



# The Problem (once more)



# The Problem (once more)

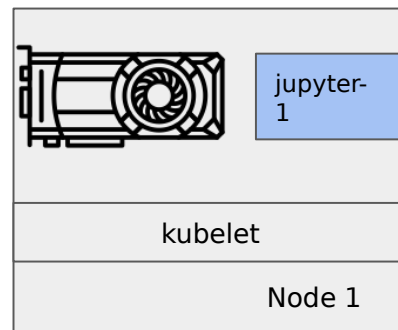
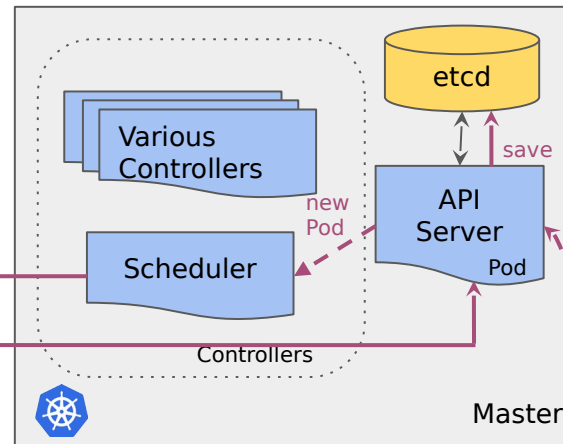


write

```
kind: Pod
metadata:
  name: jupyter-1
spec:
  image: jupyter:latest
  resources:
    limits:
      nvidia.com/gpu: 1
  nodeName: Node 1
```

schedule

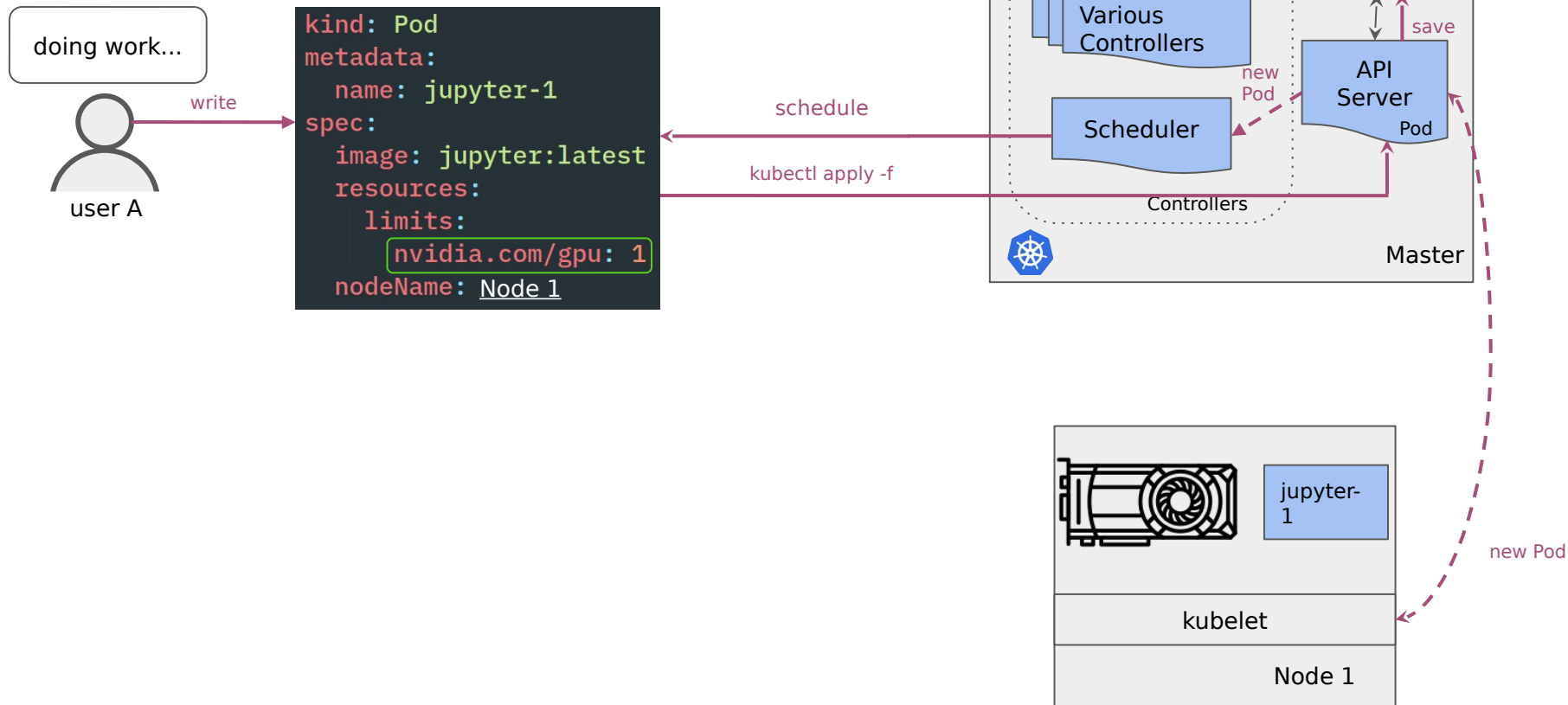
kubectl apply -f



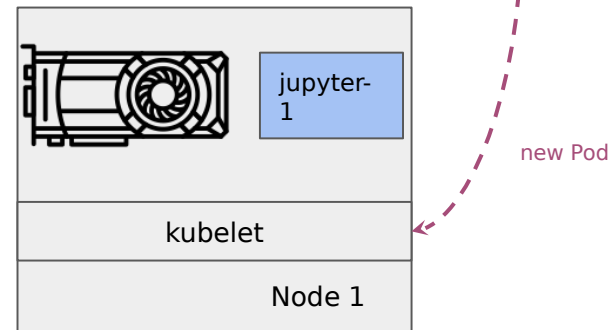
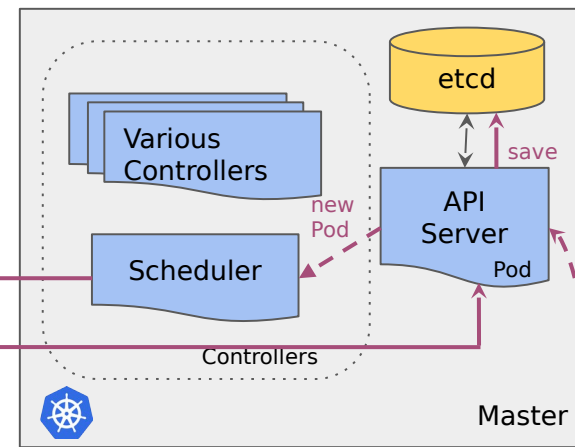
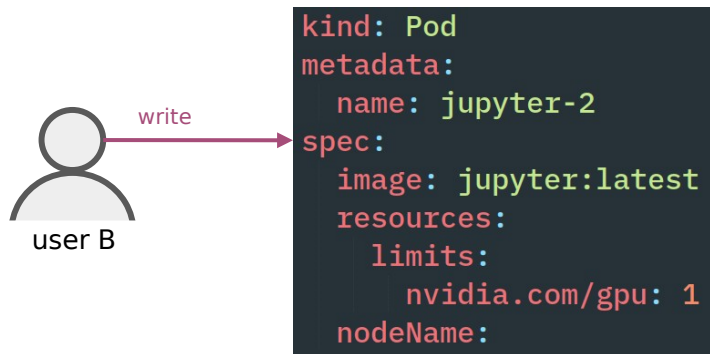
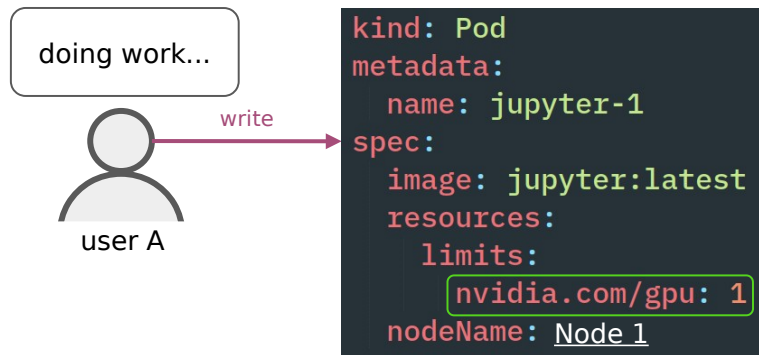
new Pod



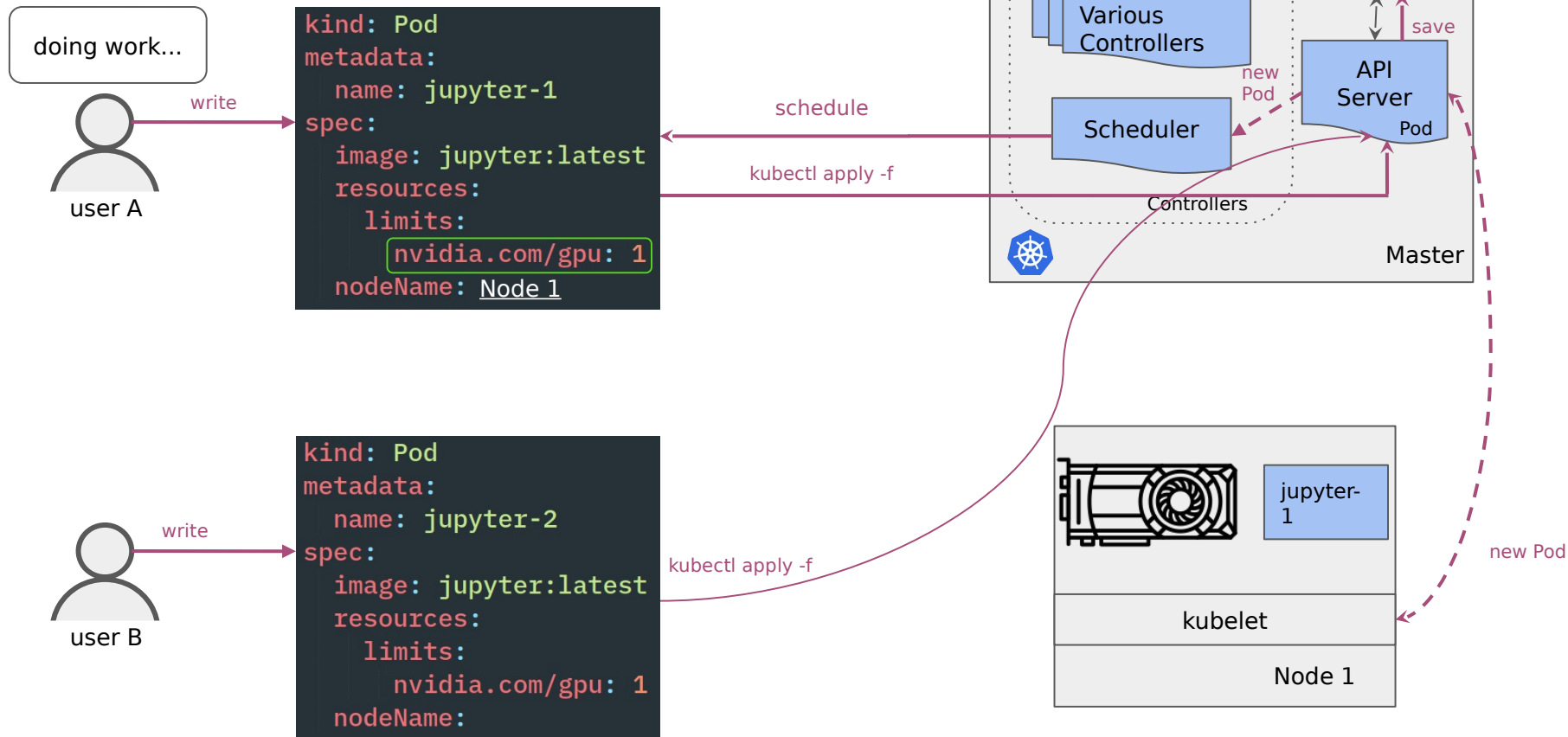
# The Problem (once more)



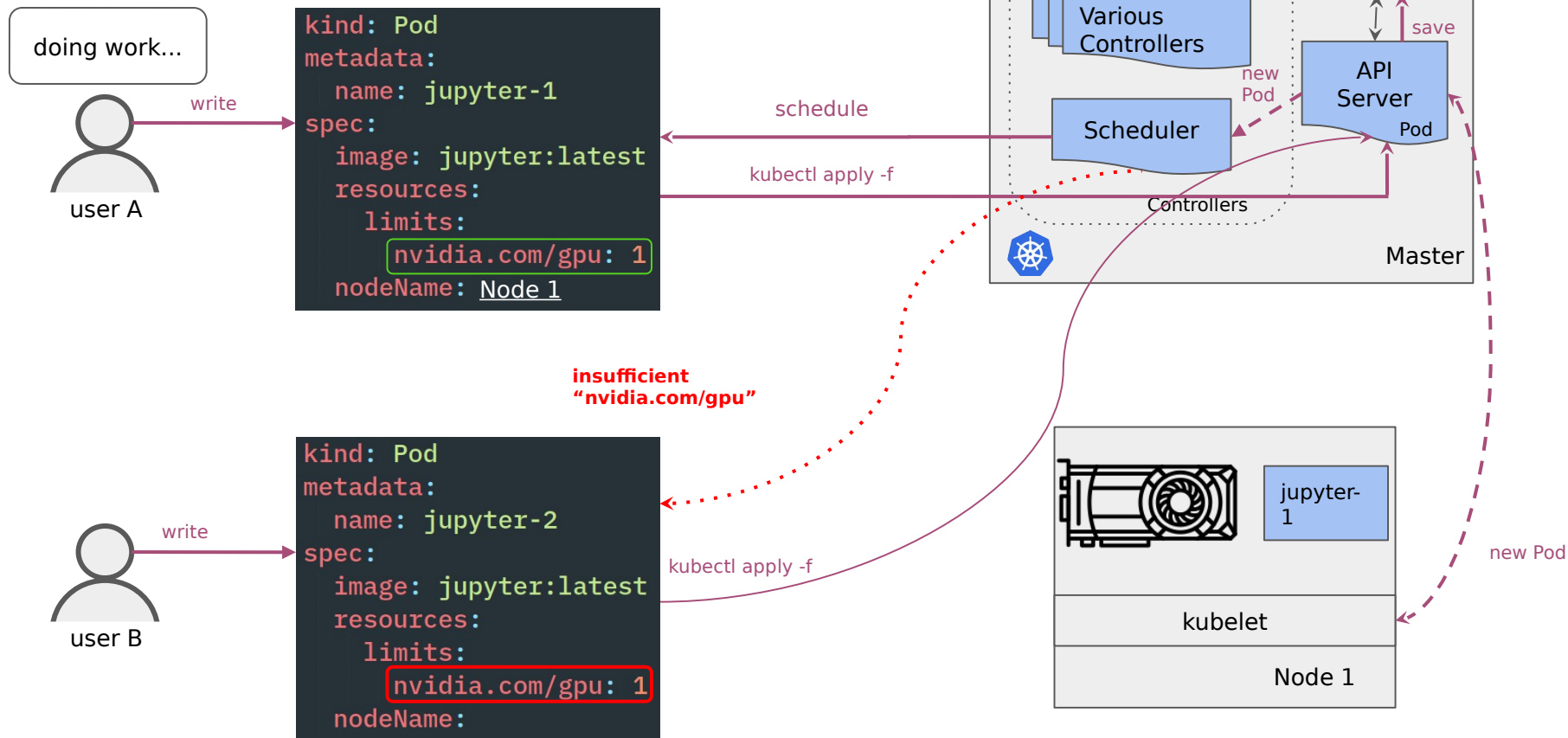
# The Problem (once more)



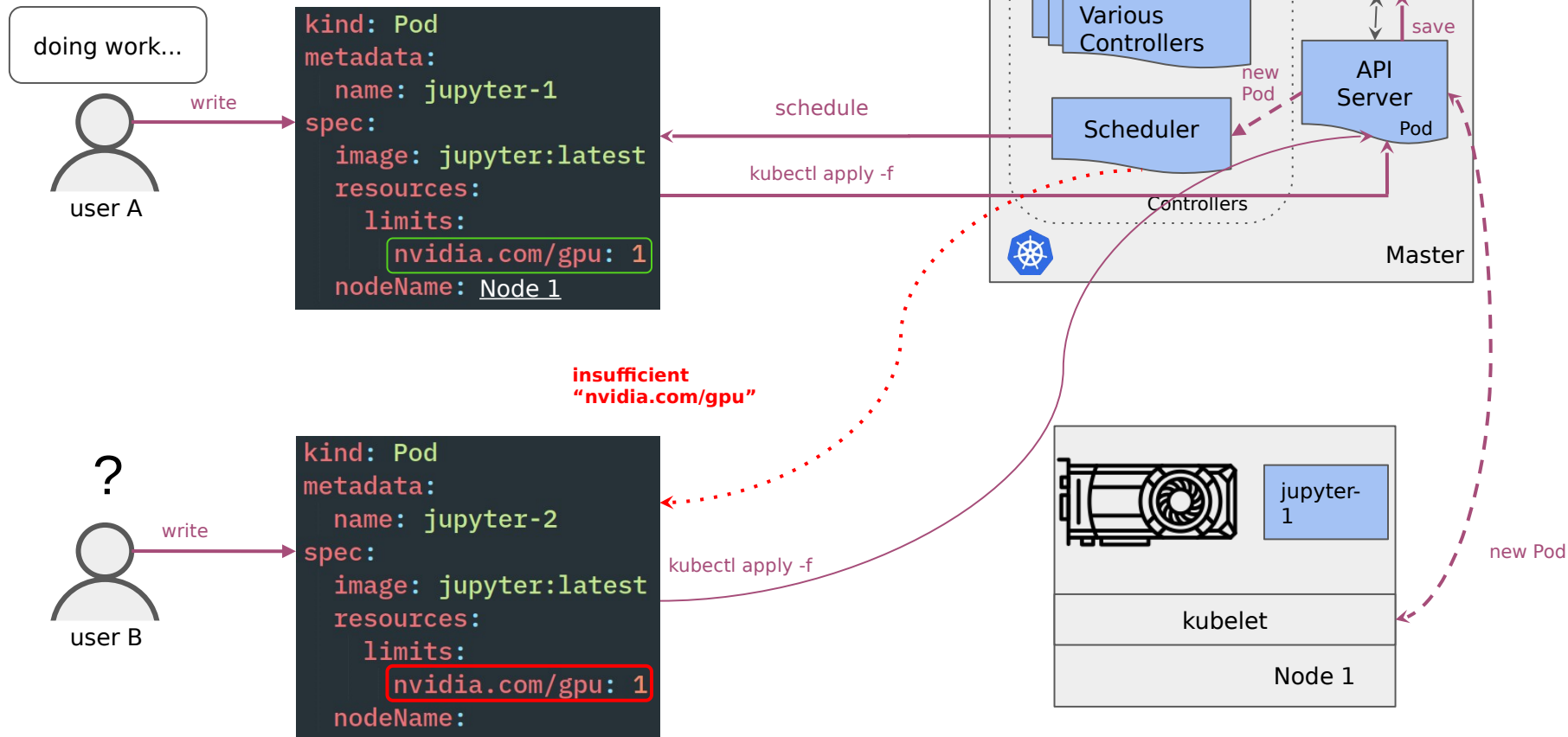
# The Problem (once more)



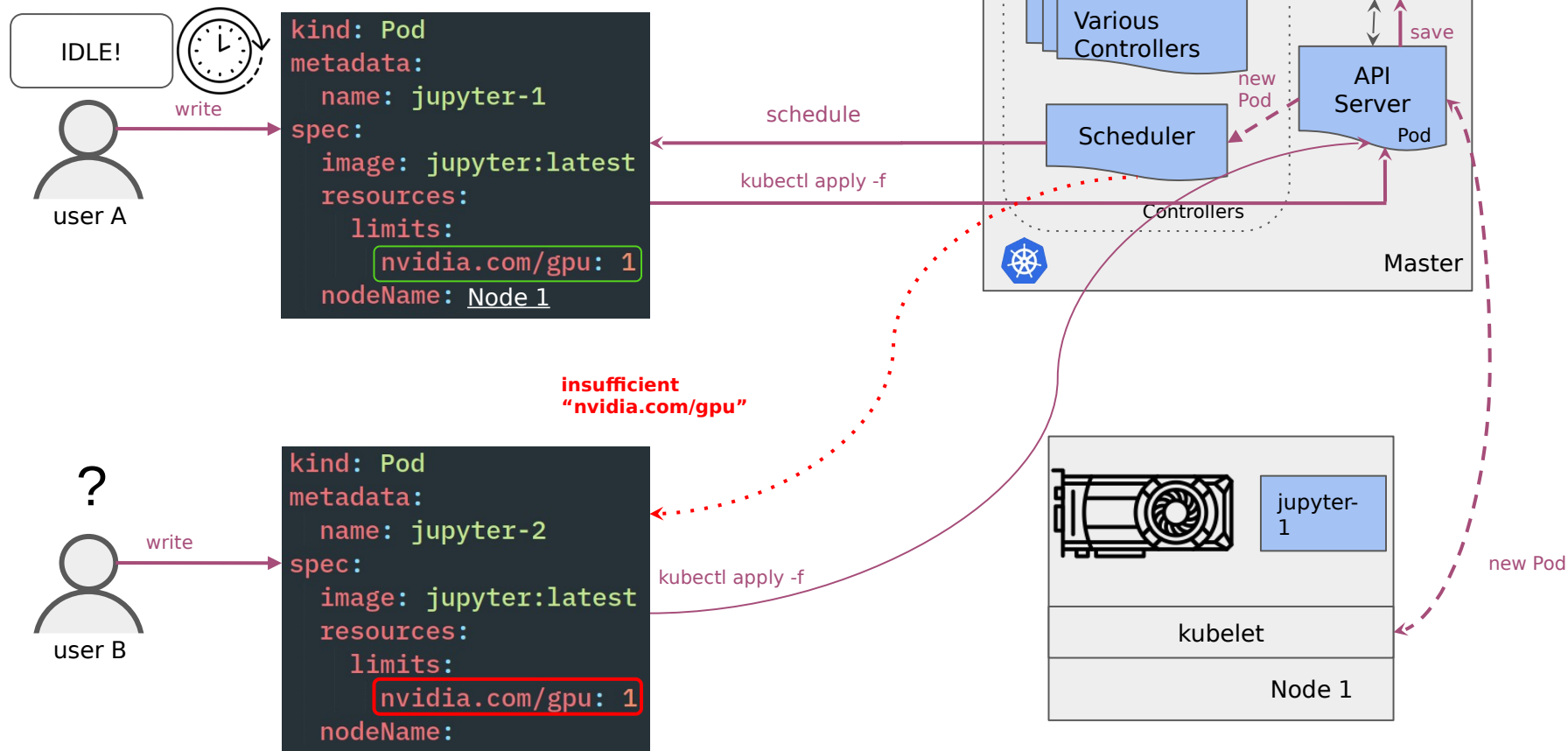
# The Problem (once more)



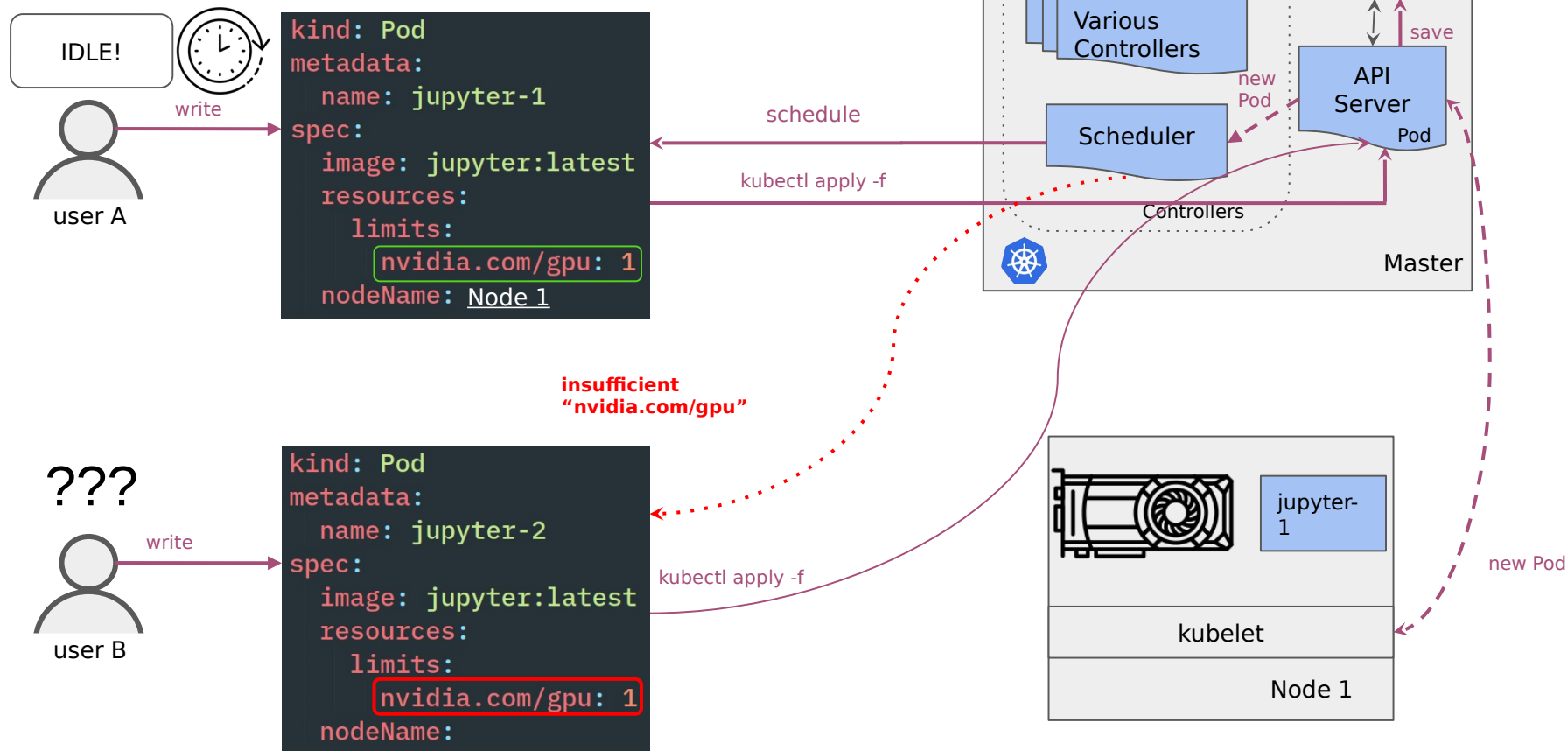
# The Problem (once more)



# The Problem (once more)



# The Problem (once more)



# Existing Approaches

Scheduler

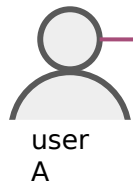
8 GB VRAM



Node 1



# Existing Approaches



write

```
kind: Pod
metadata:
  name: jupyter-1
spec:
  image: jupyter:latest
  resources:
    limits:
      aliyun.com/gpu-mem: 4
  nodeName:
```

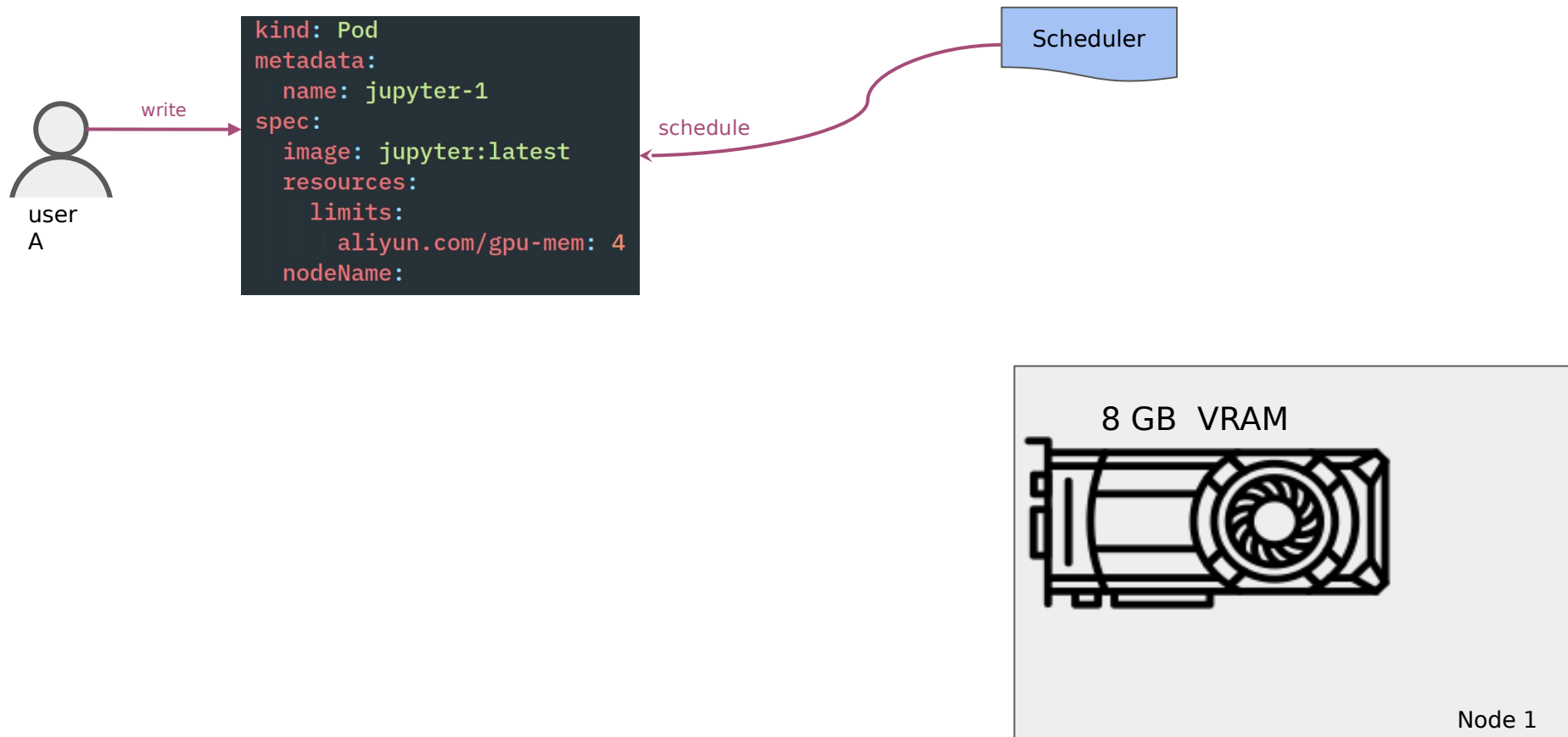
Scheduler

8 GB VRAM

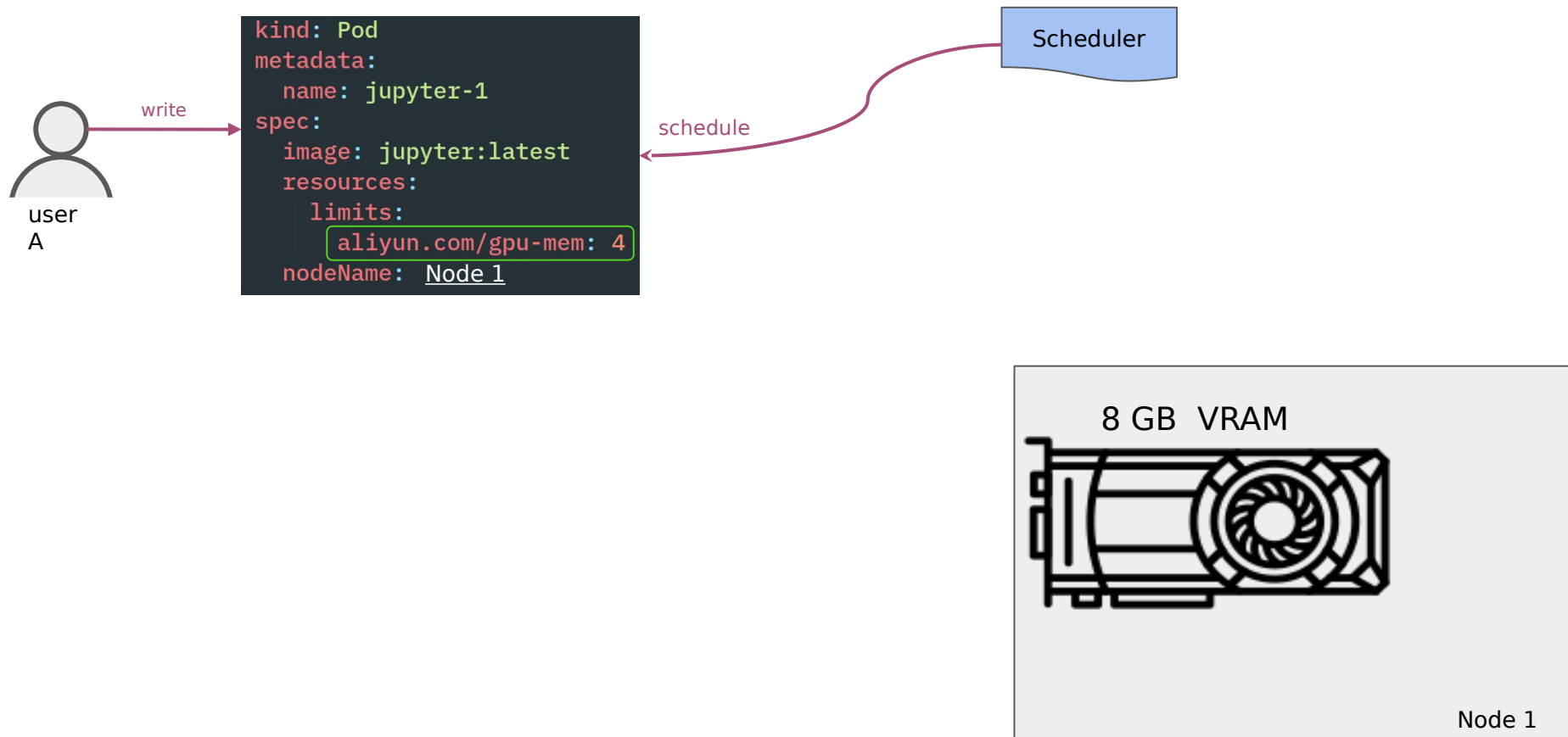


Node 1

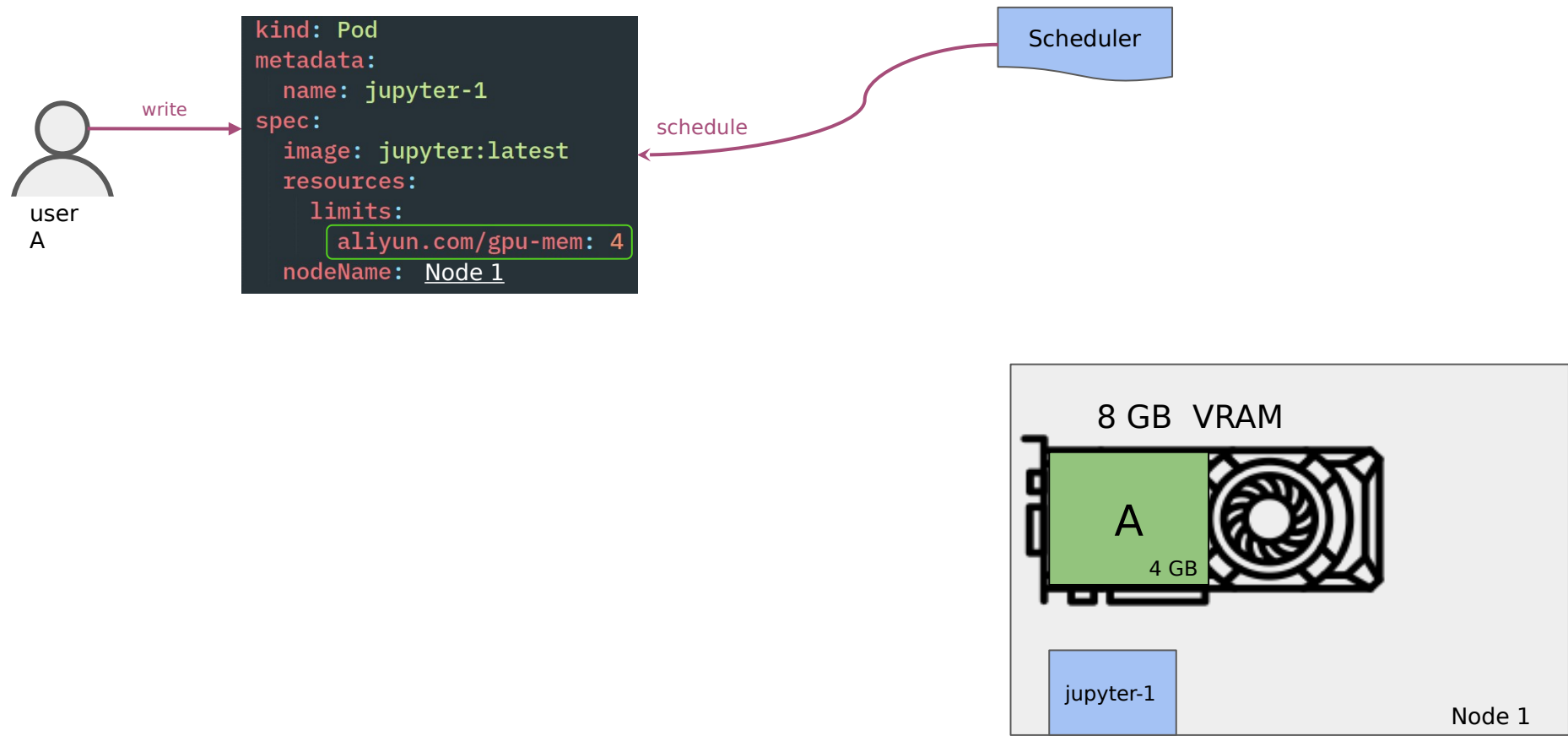
# Existing Approaches



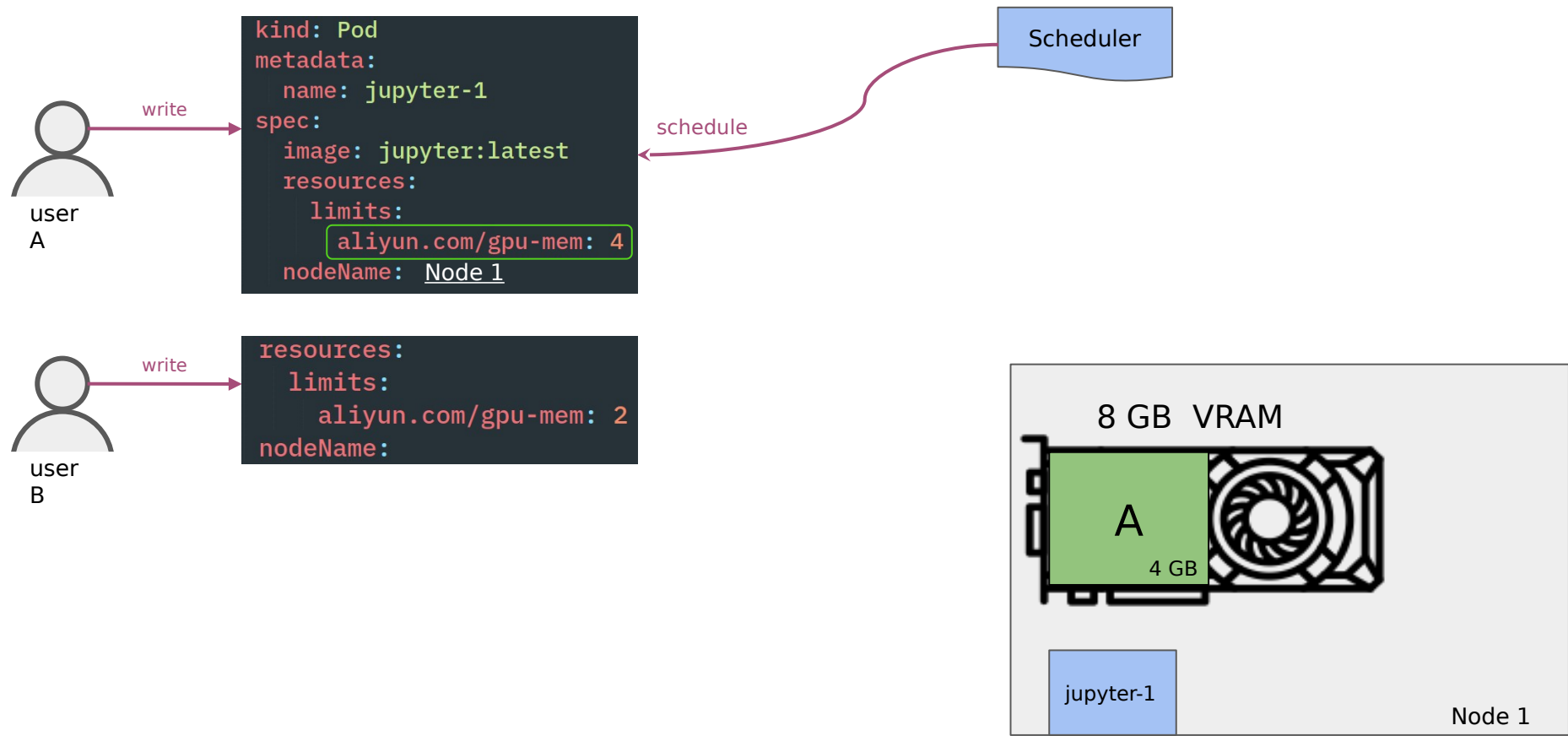
# Existing Approaches



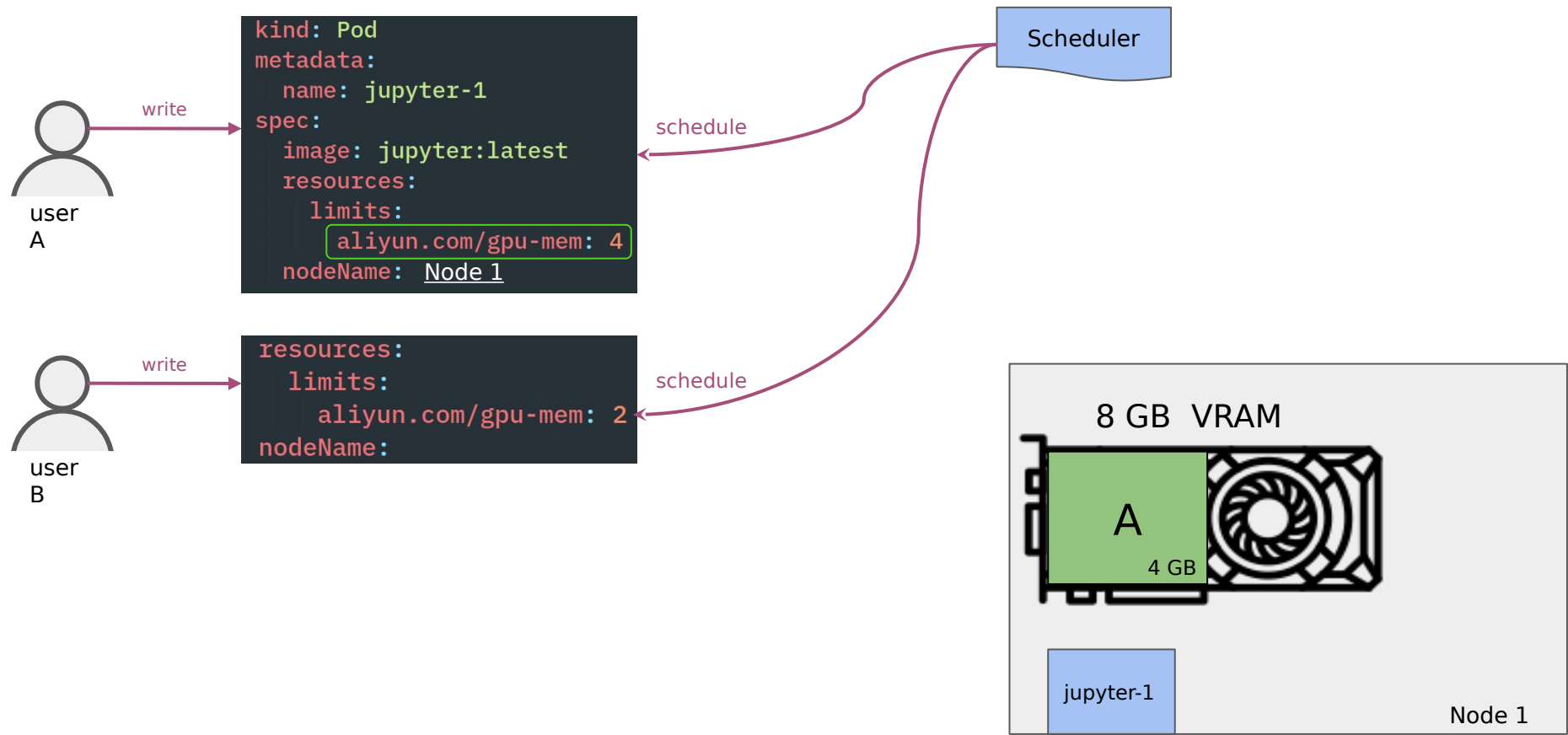
# Existing Approaches



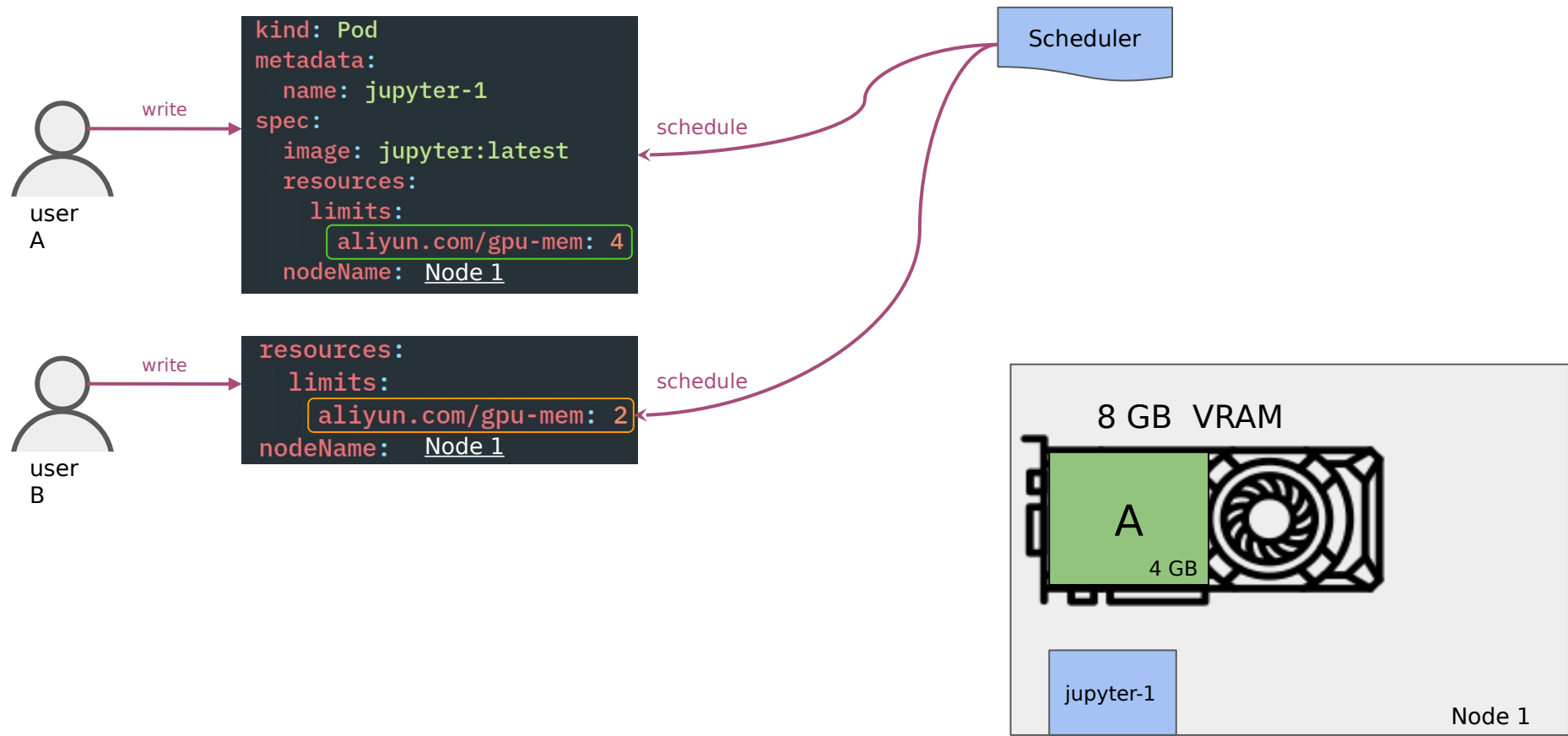
# Existing Approaches



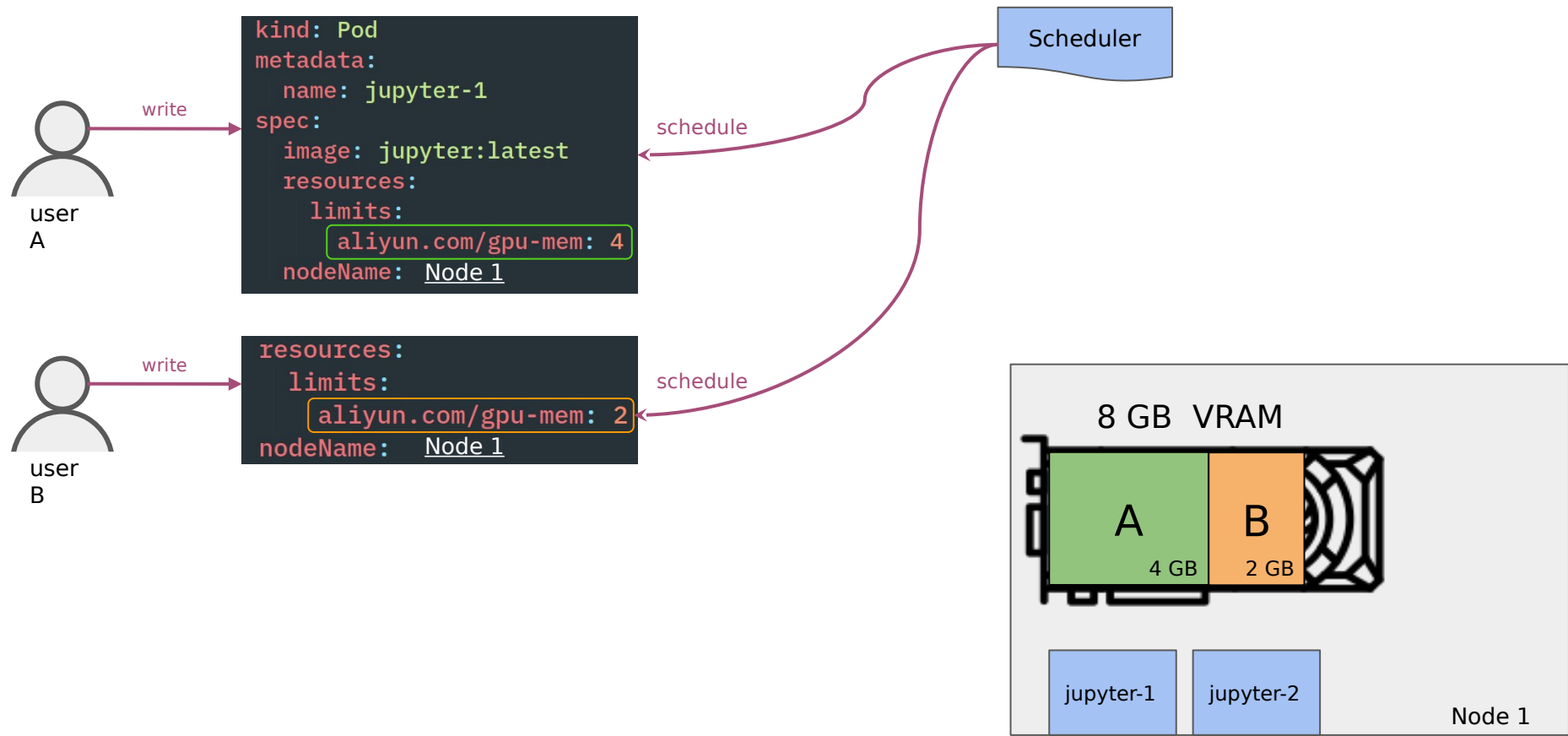
# Existing Approaches



# Existing Approaches

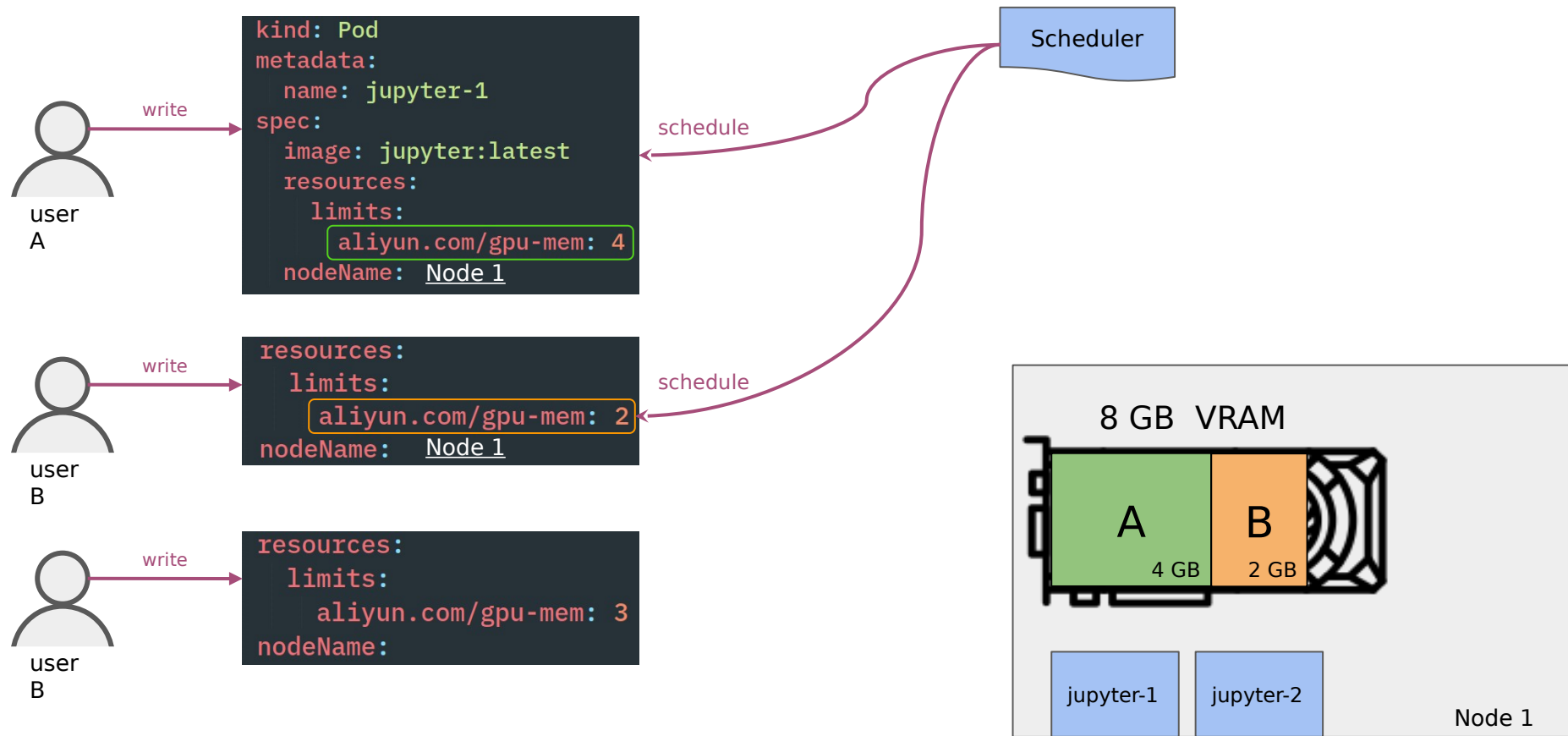


# Existing Approaches

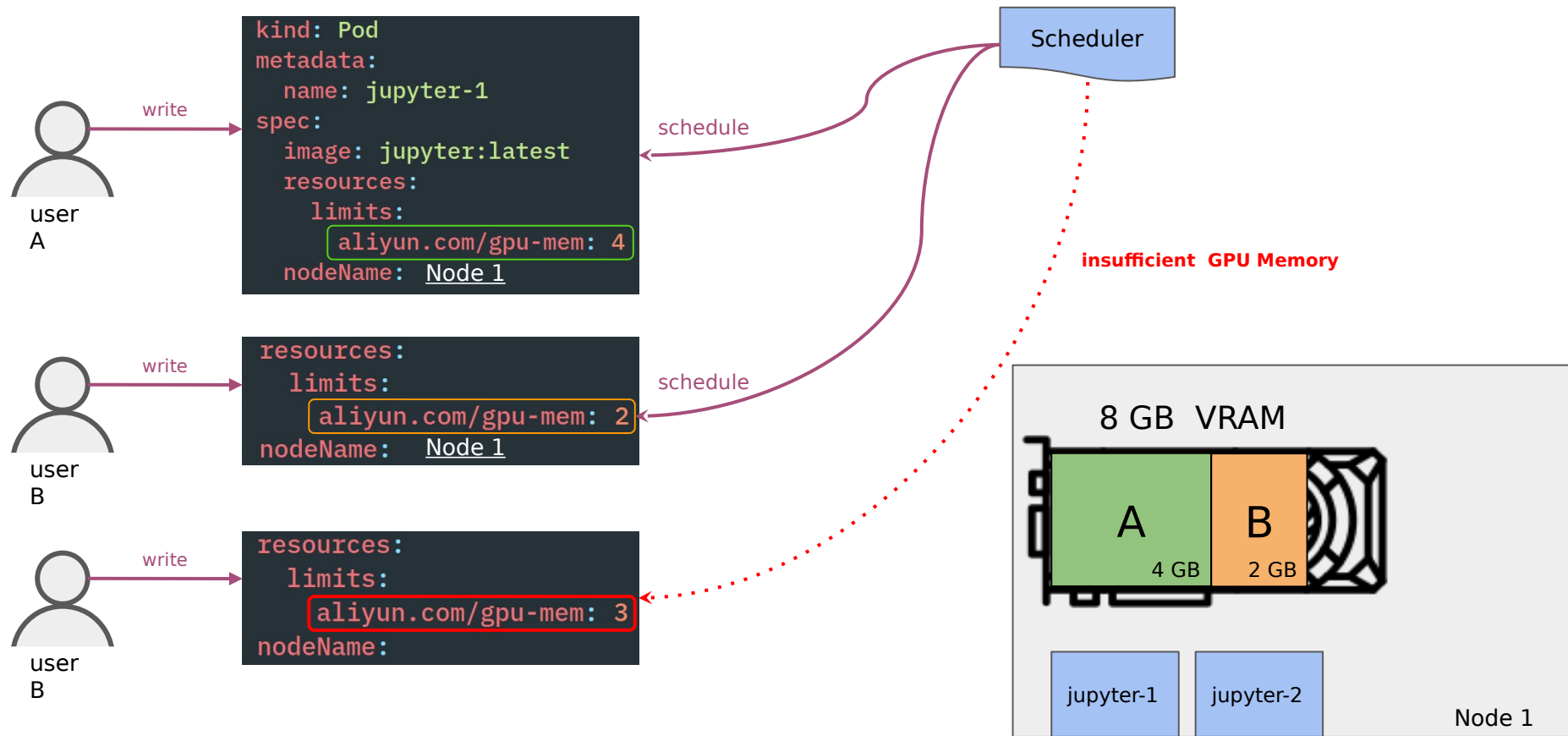




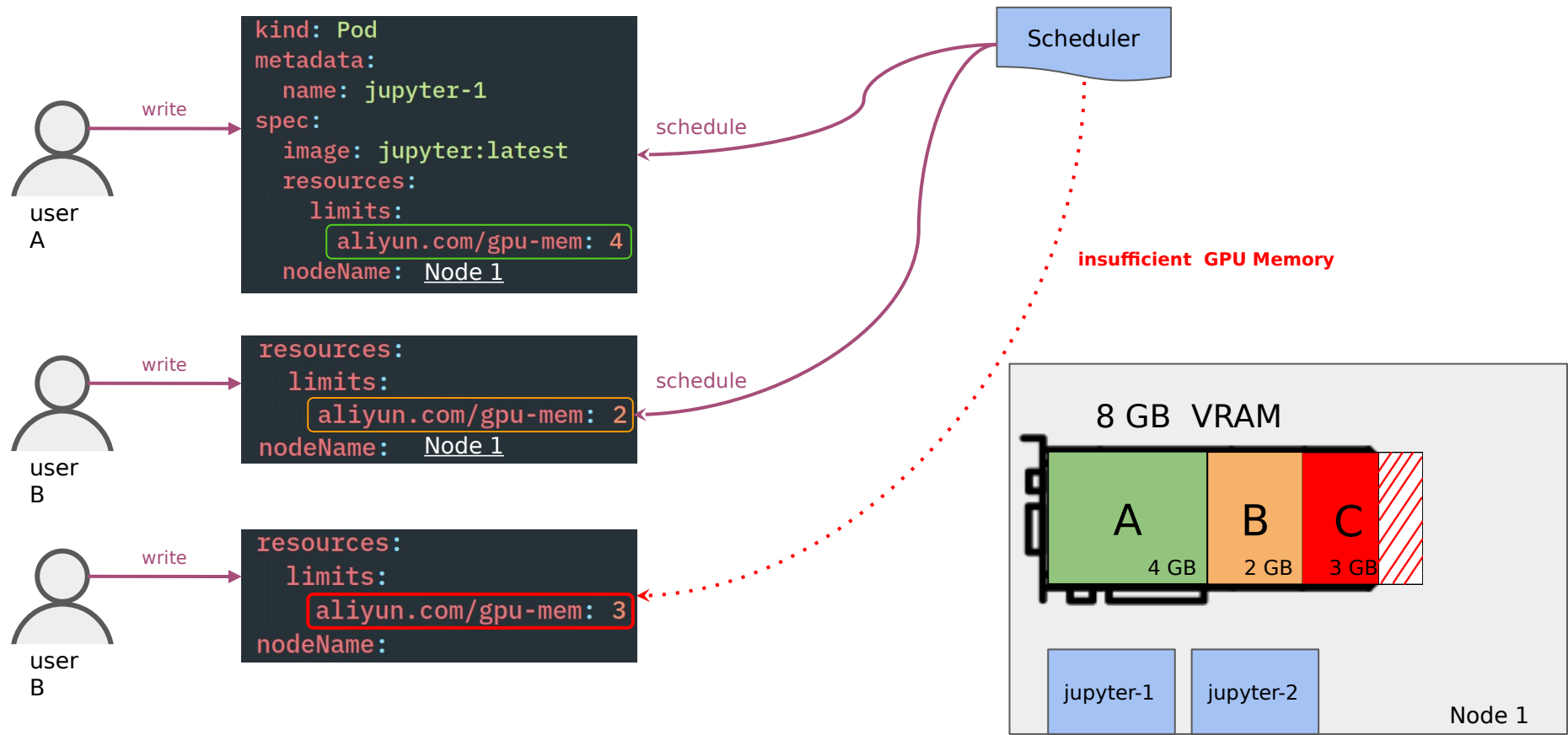
# Existing Approaches



# Existing Approaches



# Existing Approaches



# Existing Approaches



## ● Aliyun Scheduler Extender

- GPUs no longer exclusively assigned to Pods
- users specify GPU memory requests
- uses GPU memory to bin pack Pods in Kubernetes
- does not enforce memory limits (only used for scheduling)

## ● Kubeshare (HPDC '20)

- similar to Aliyun w.r.t. the scheduling part
- enforces GPU memory limits (assigns a memory slice to each process)
- wraps CUDA API to keep track of how much memory each process allocates



# Memory

CPU

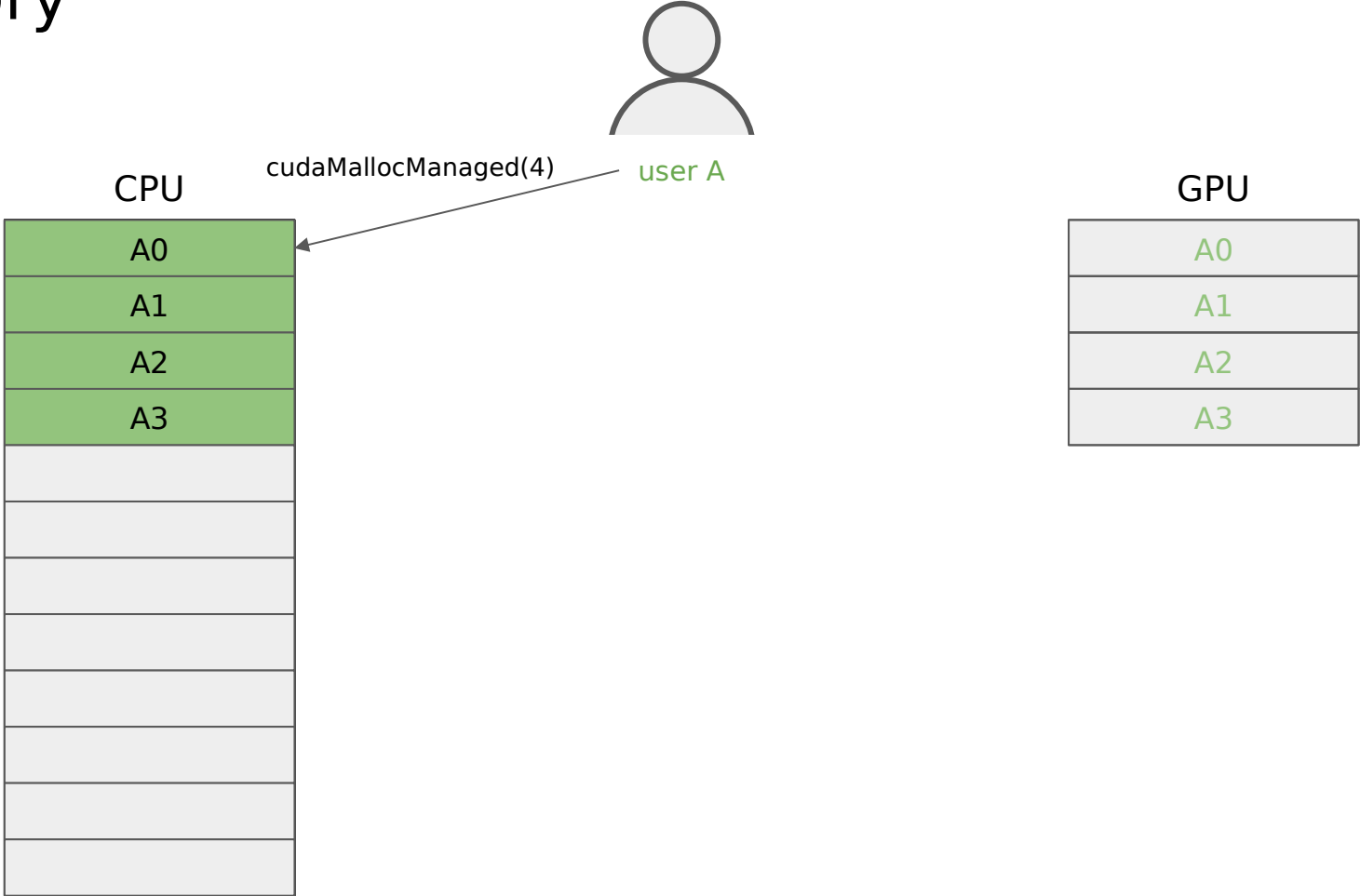


user A

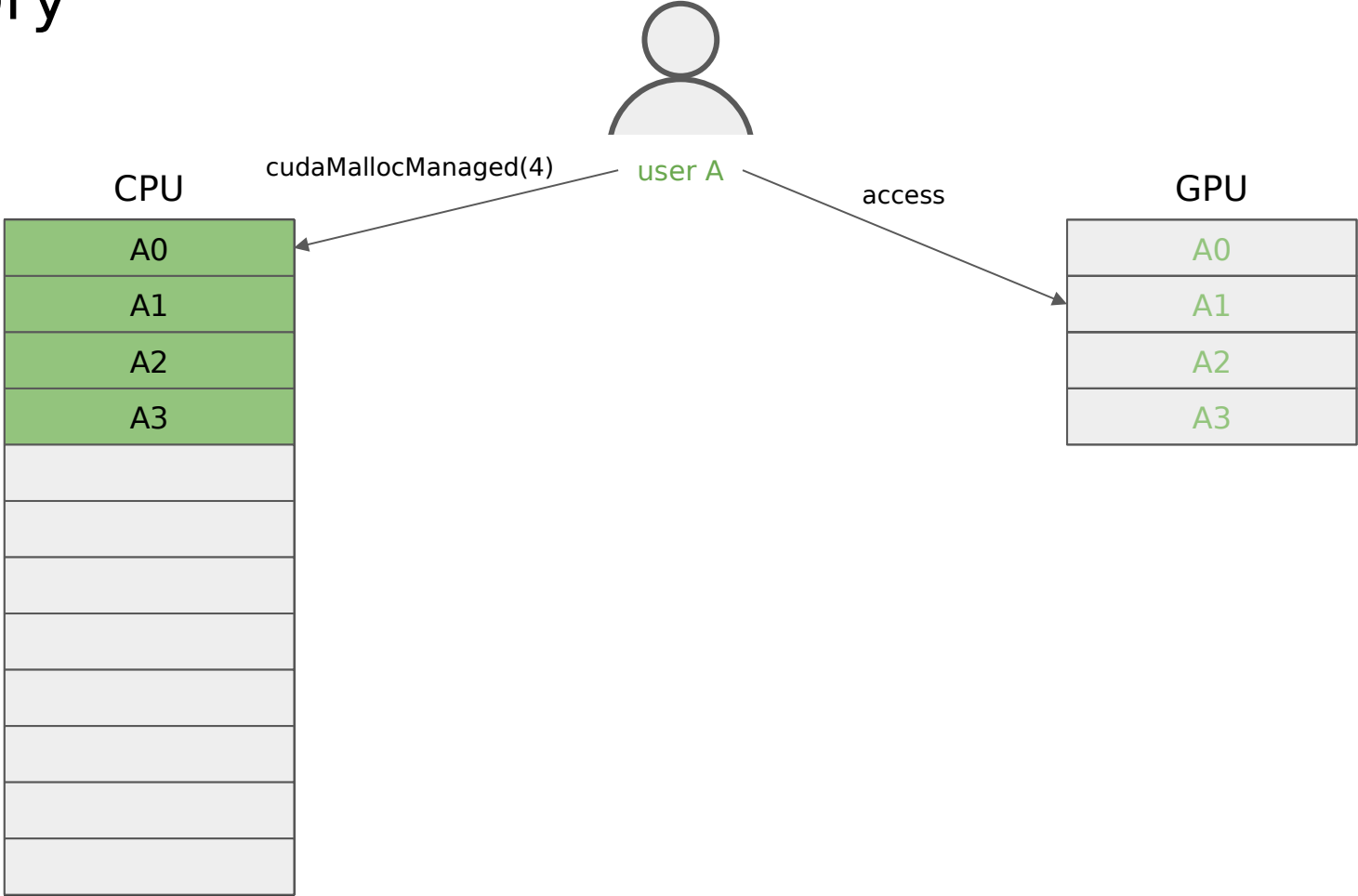
GPU



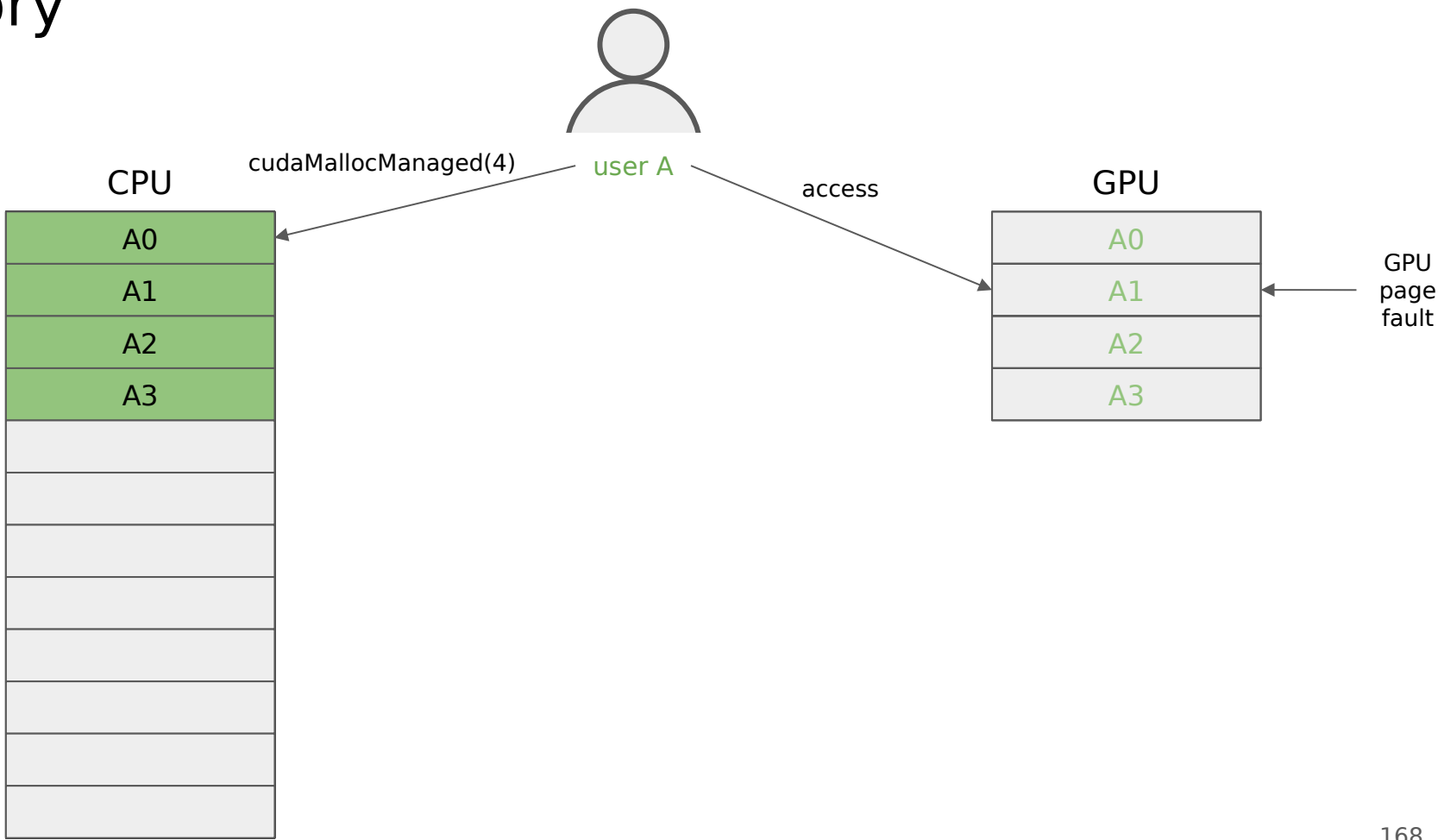
# Memory



# Memory

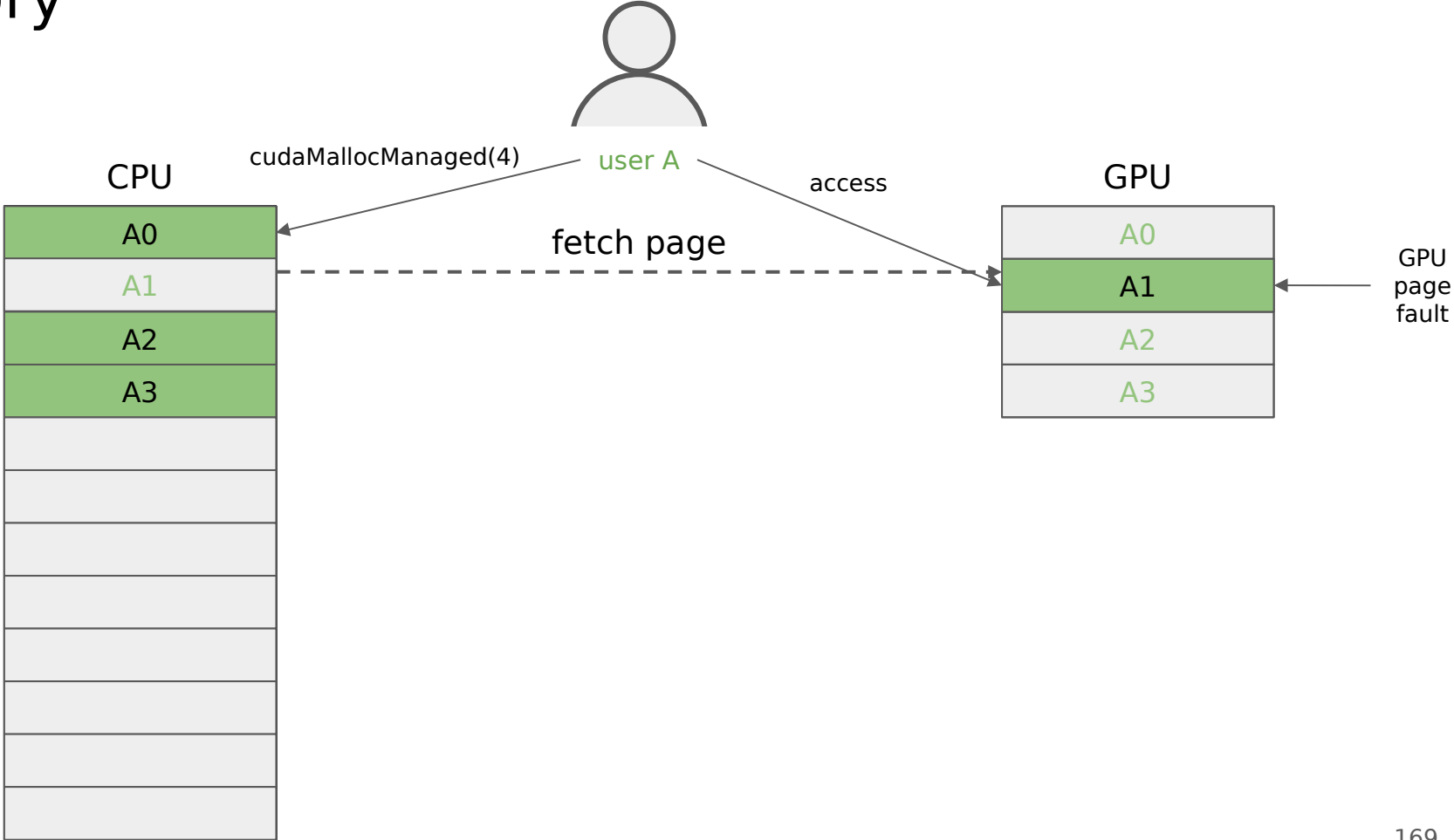


# Memory

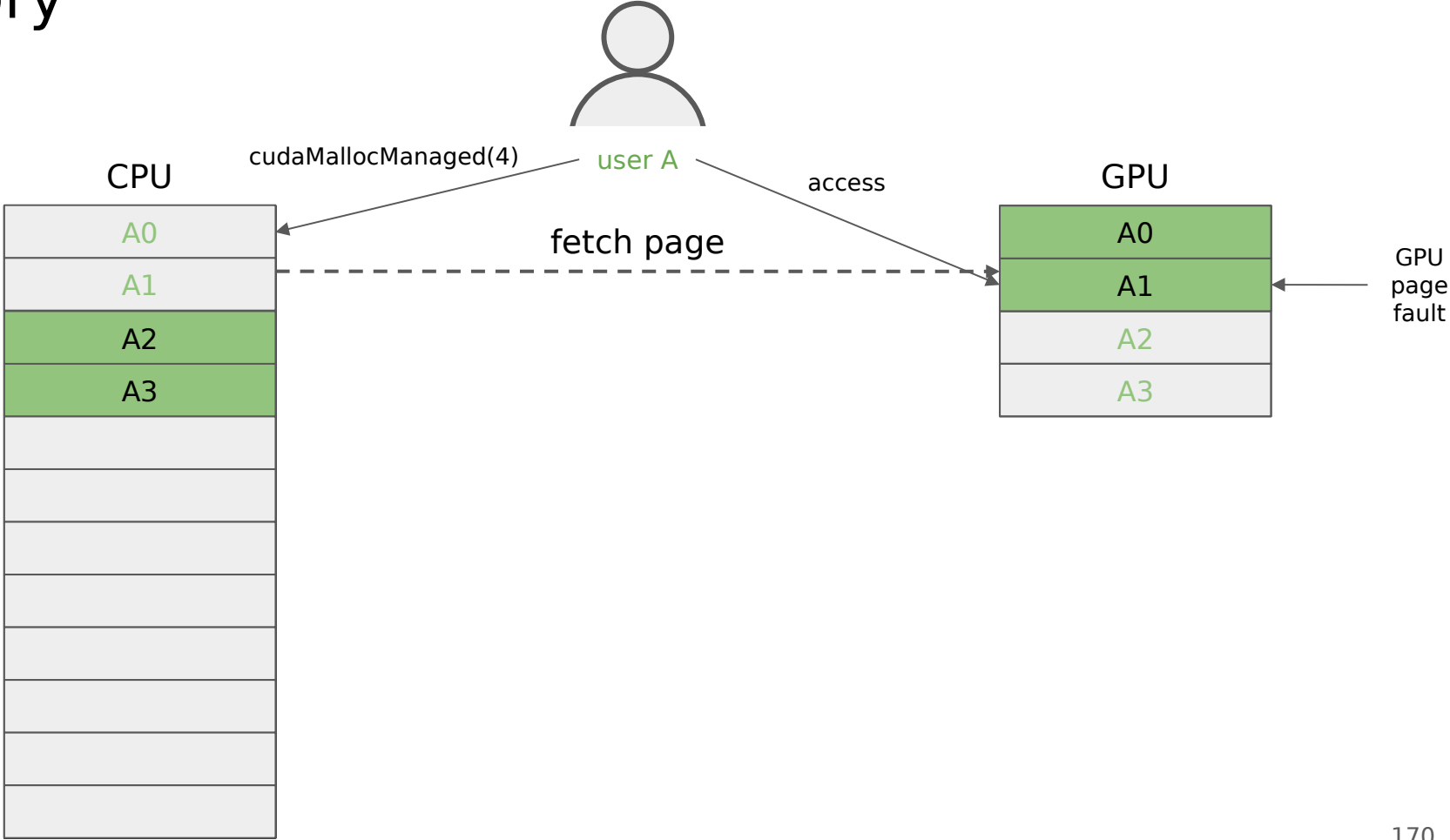




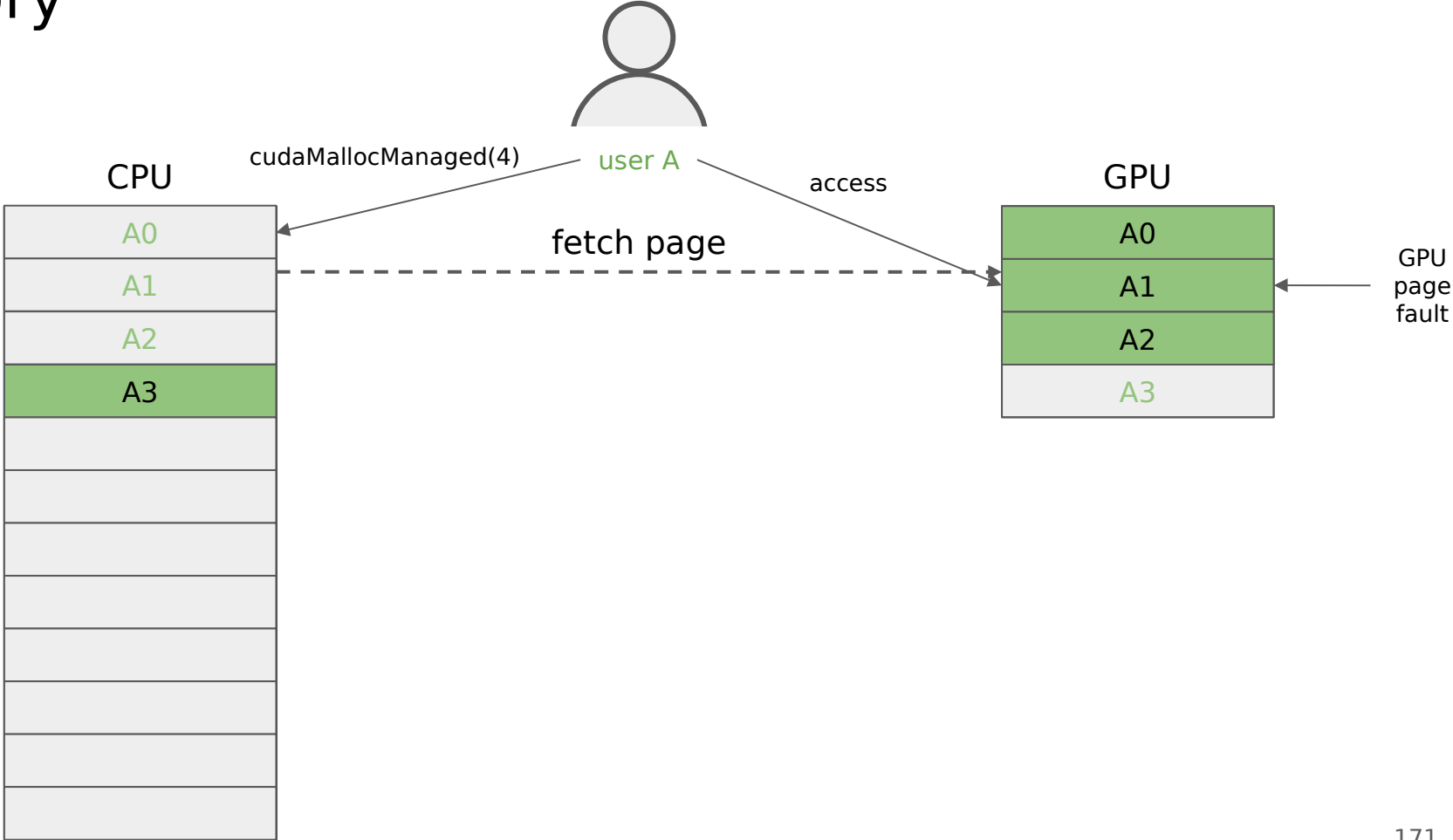
# Memory



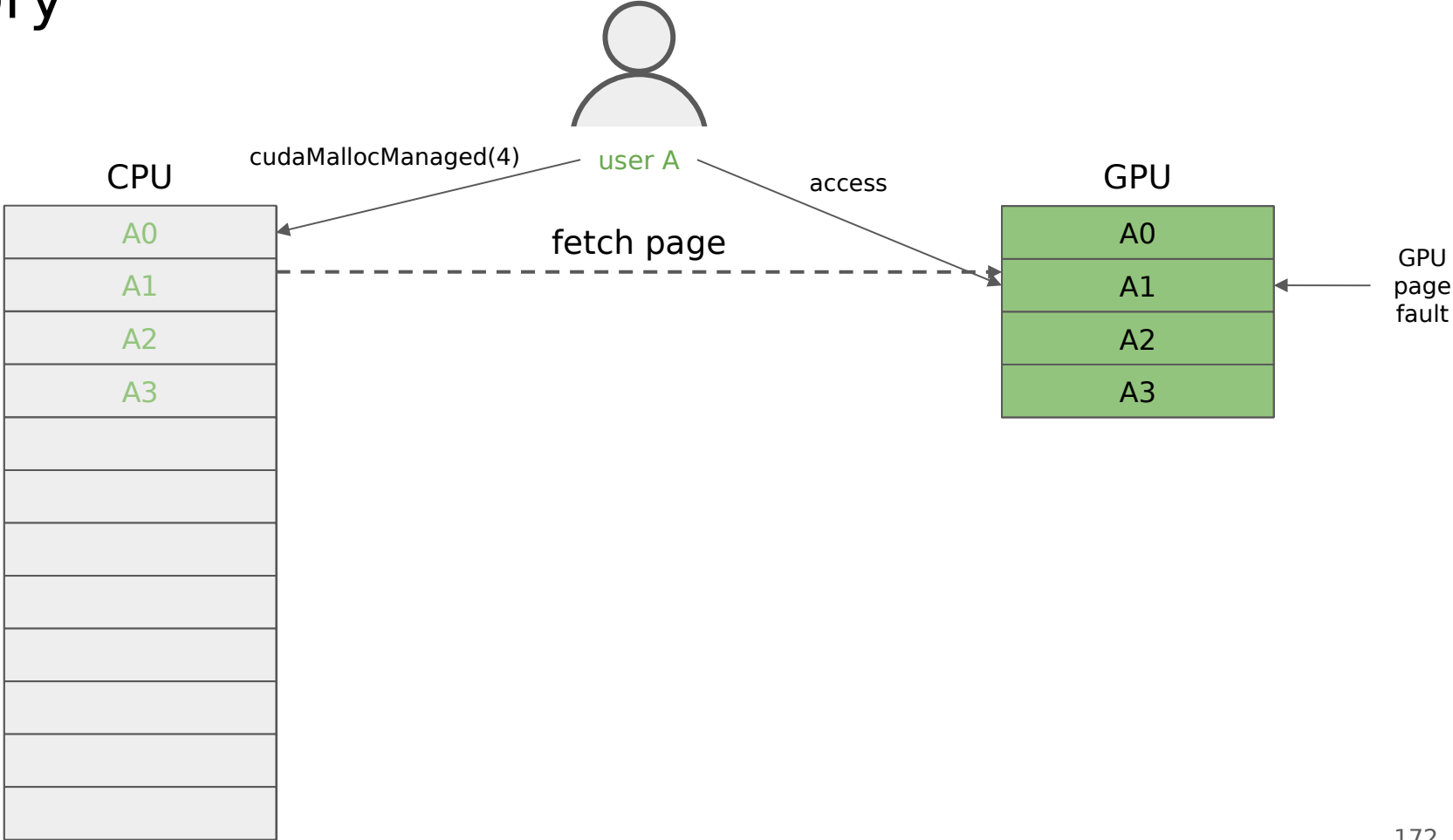
# Memory



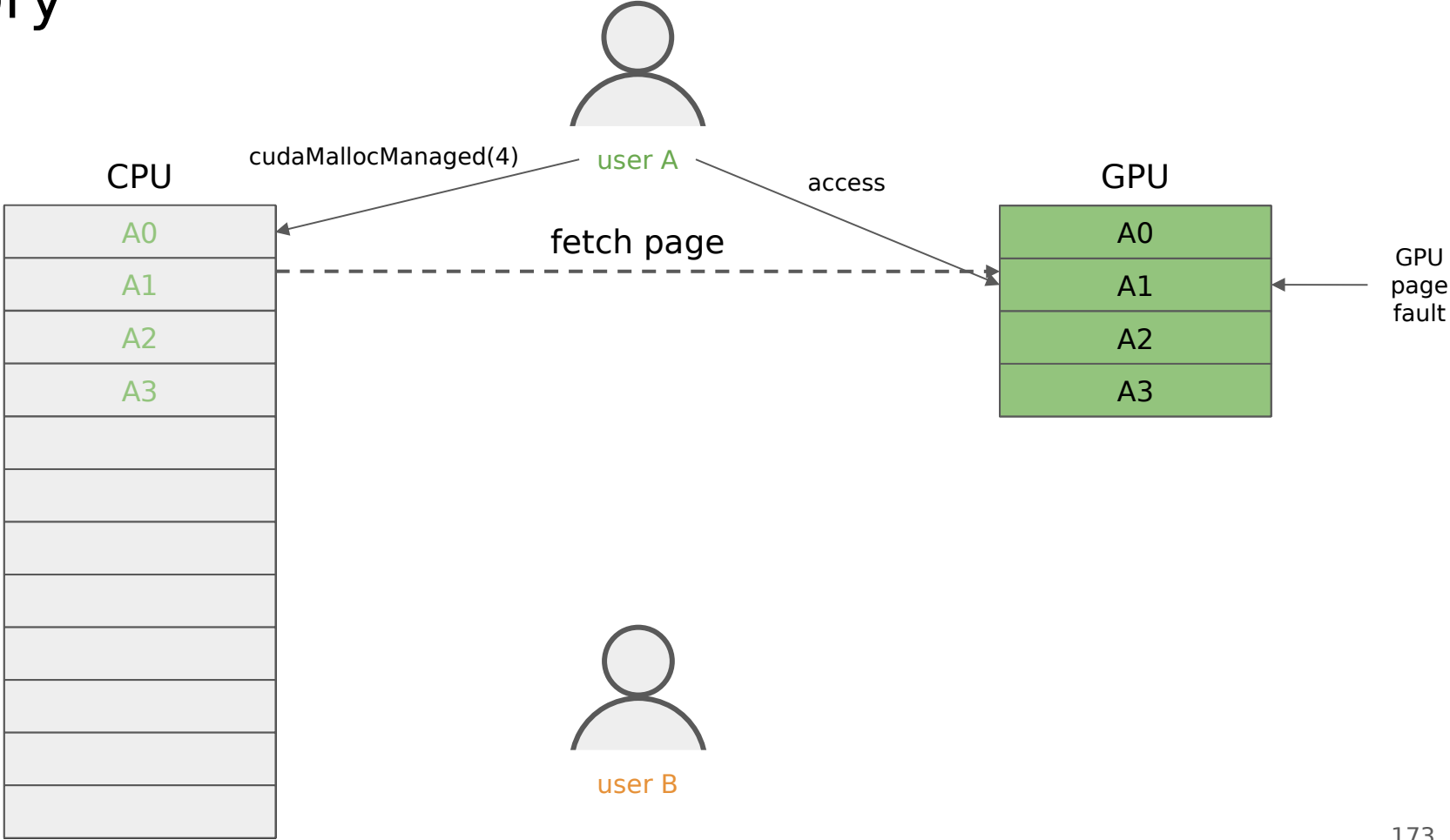
# Memory



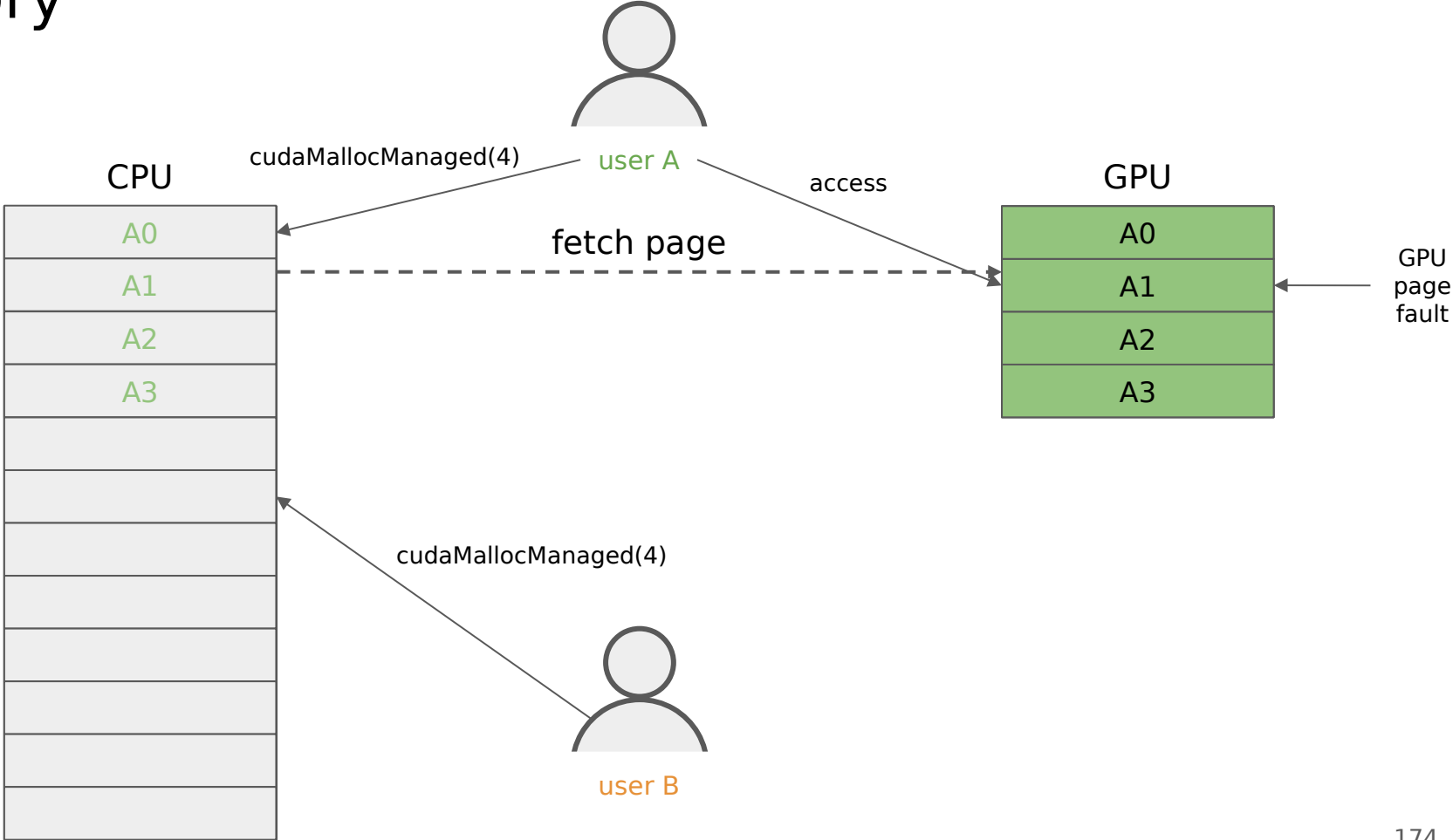
# Memory



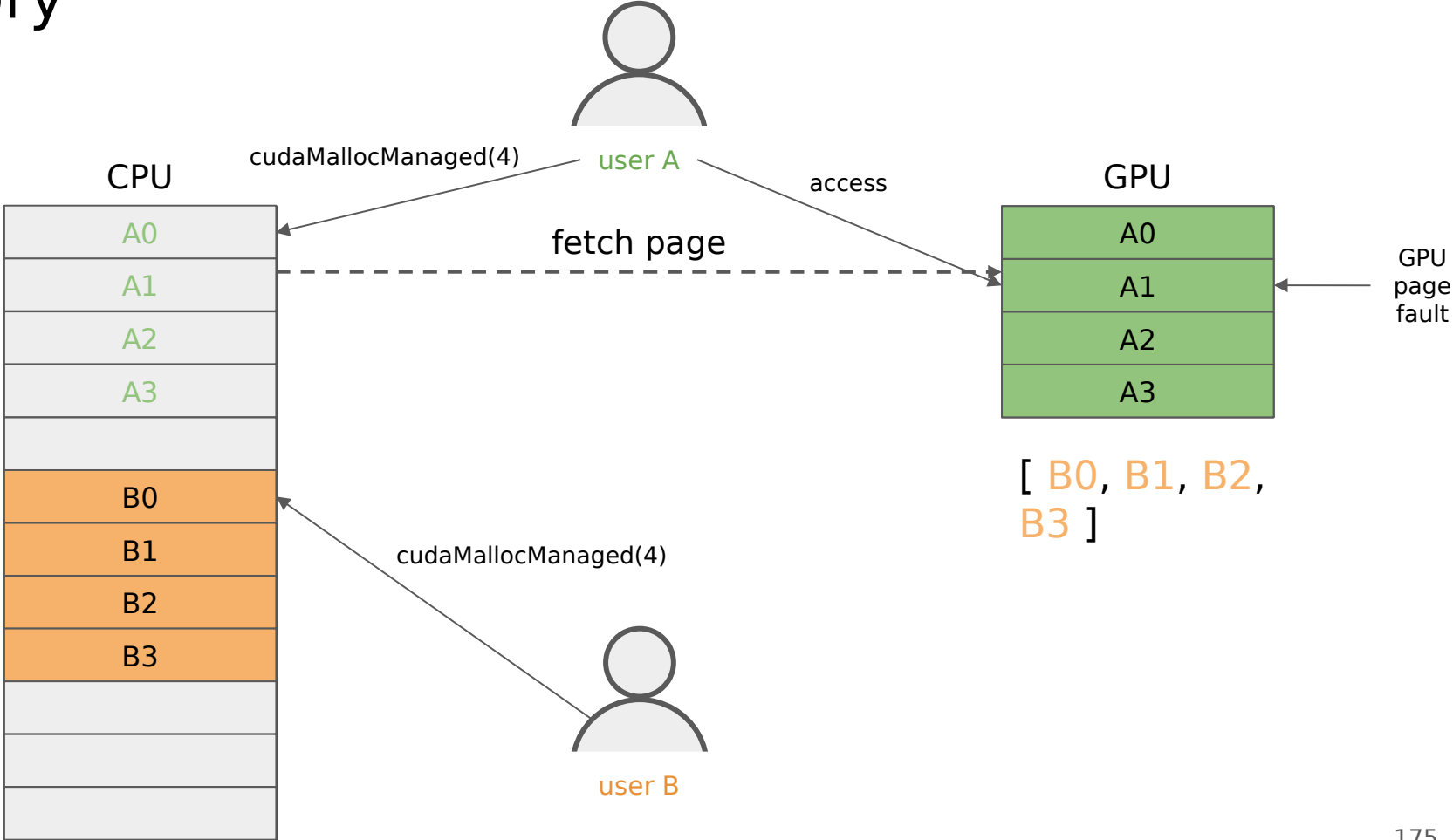
# Memory



# Memory

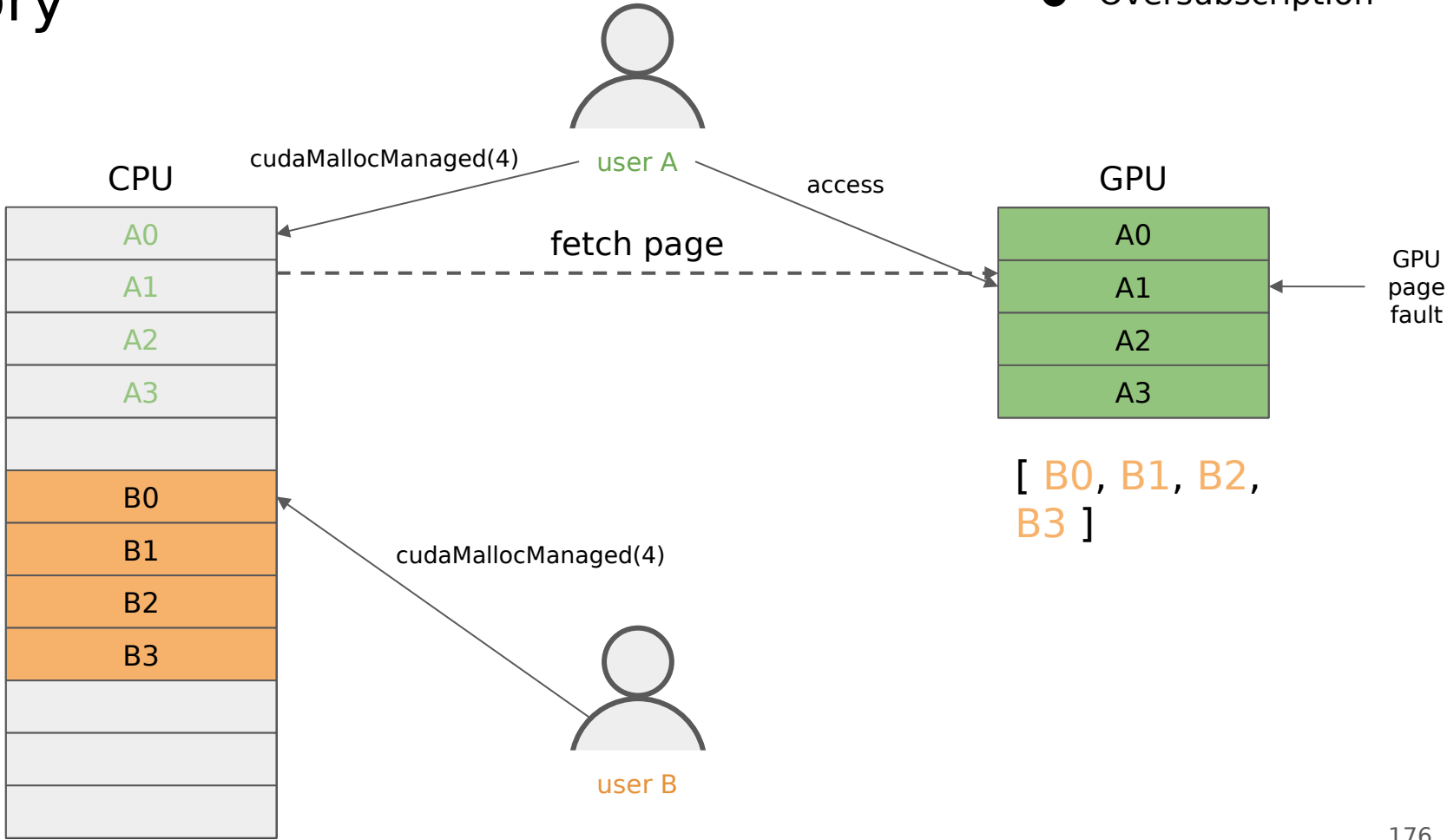


# Memory



# Memory

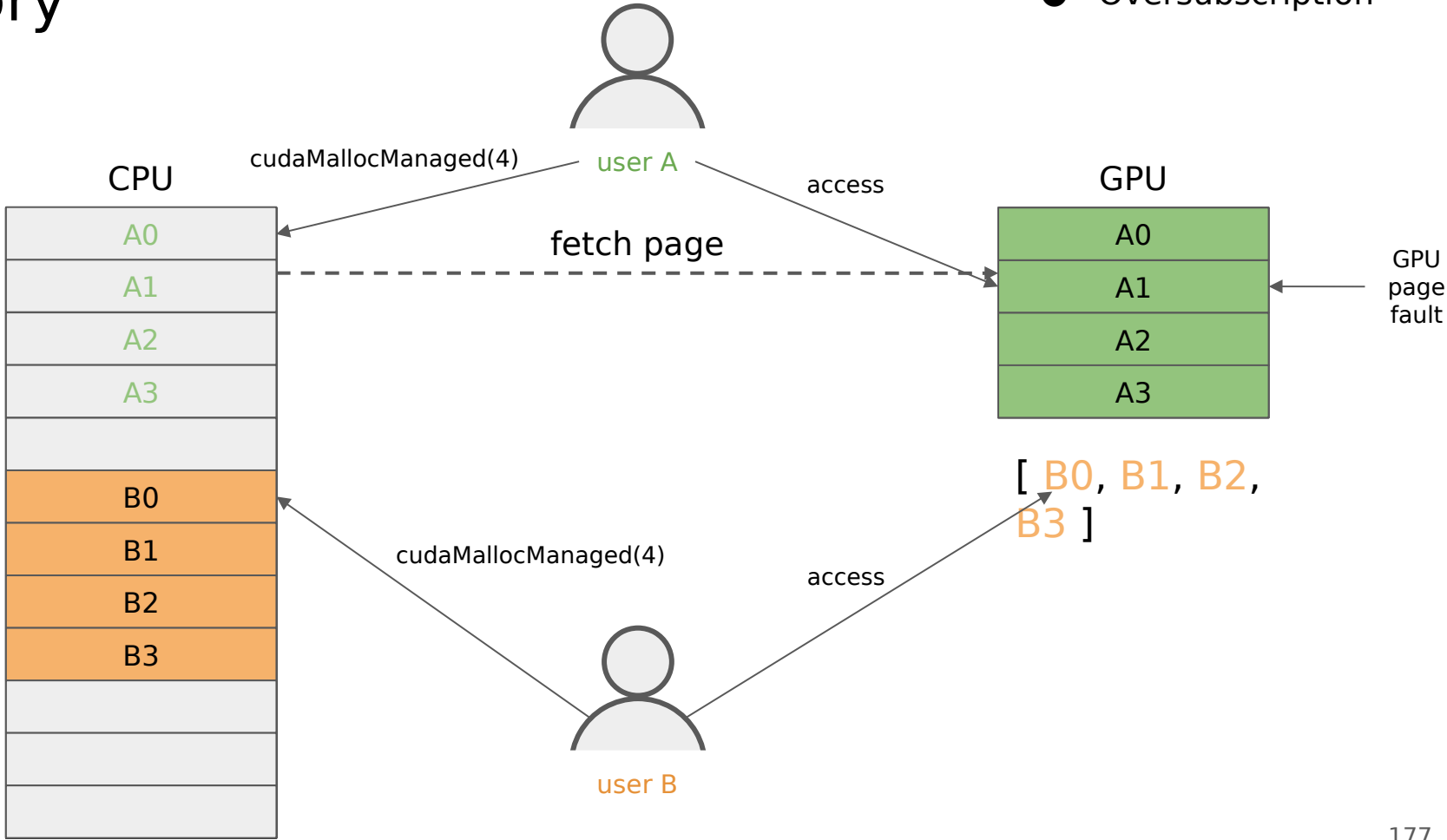
● Oversubscription





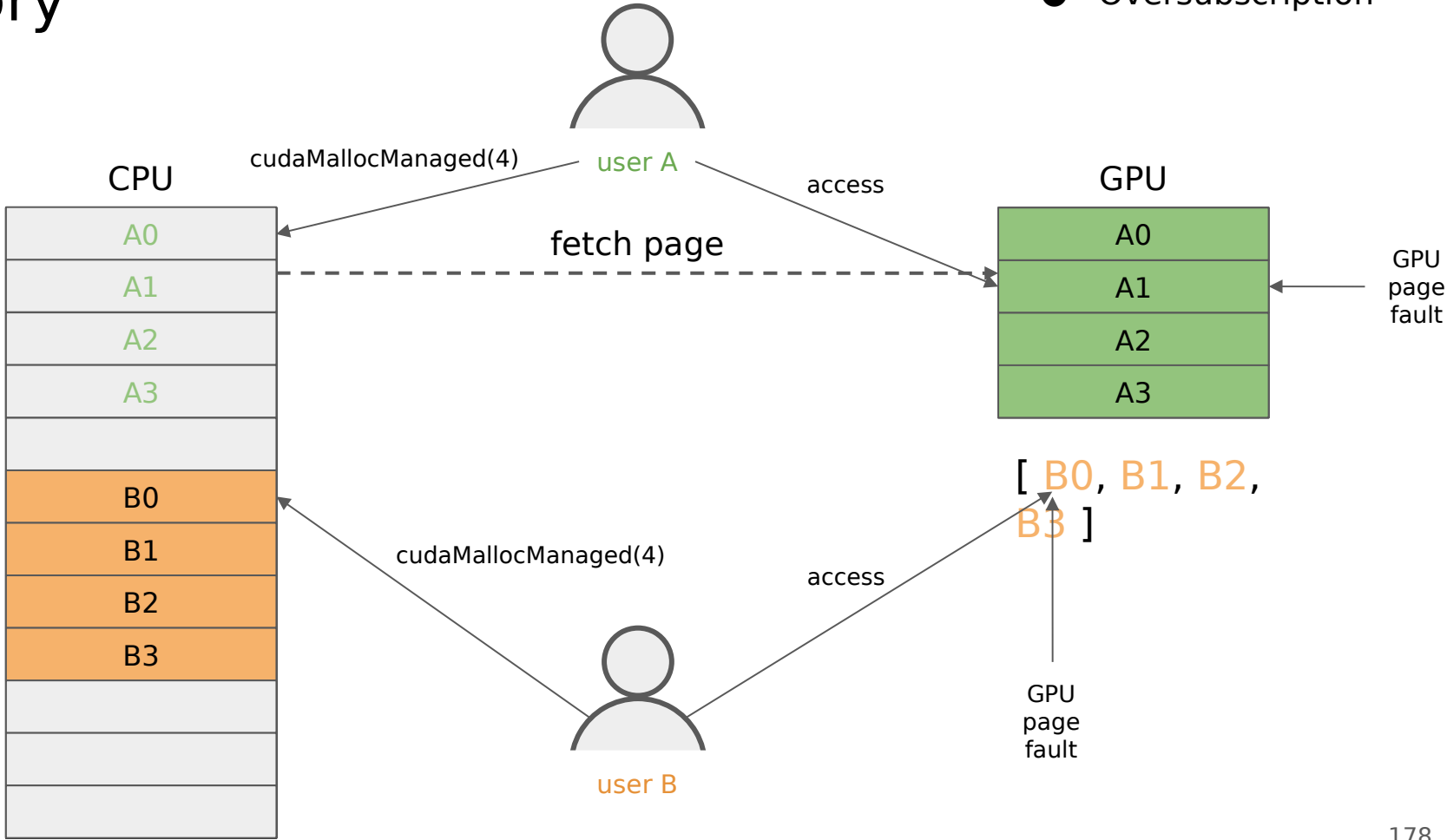
# Memory

● Oversubscription



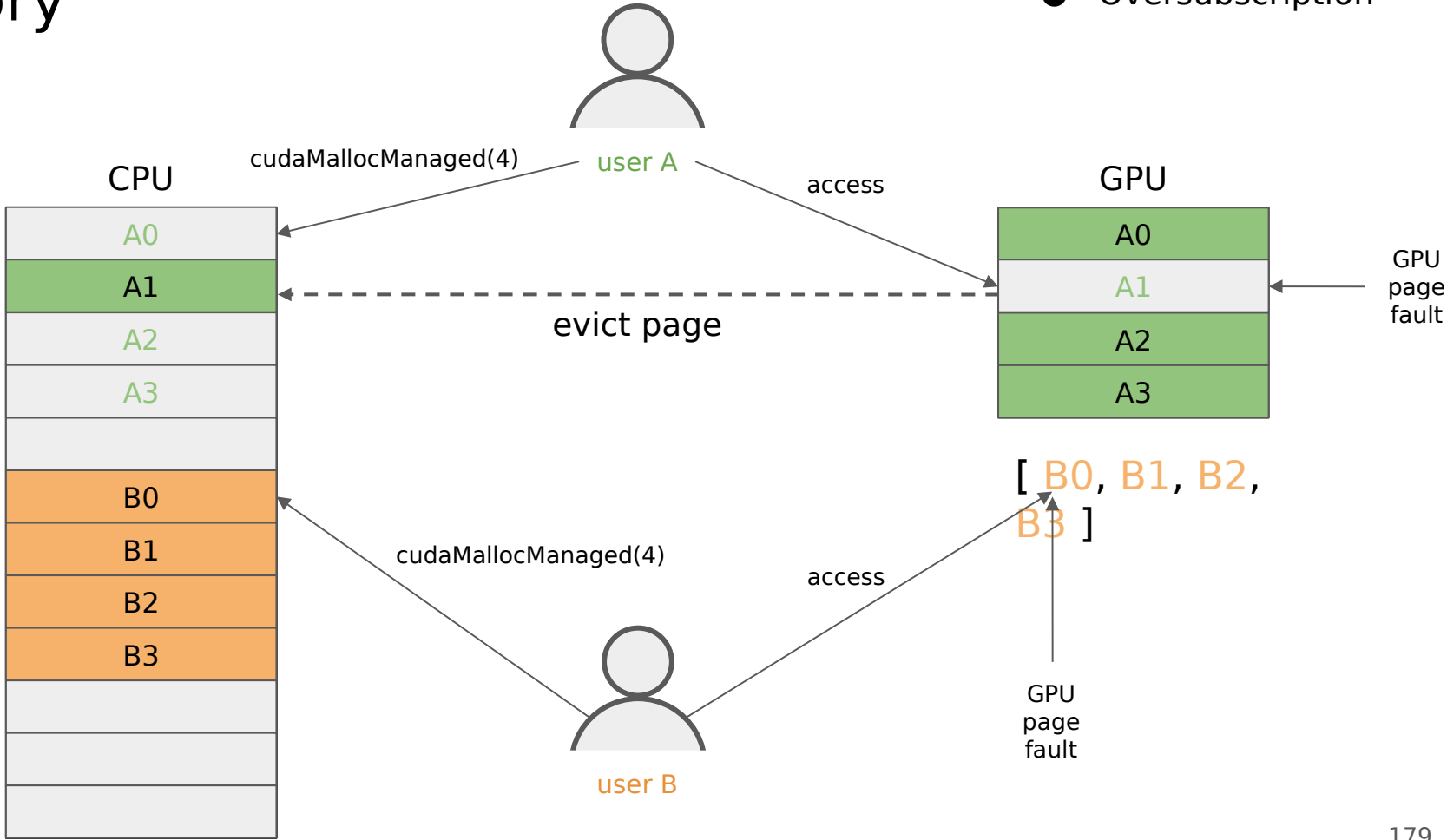
# Memory

● Oversubscription



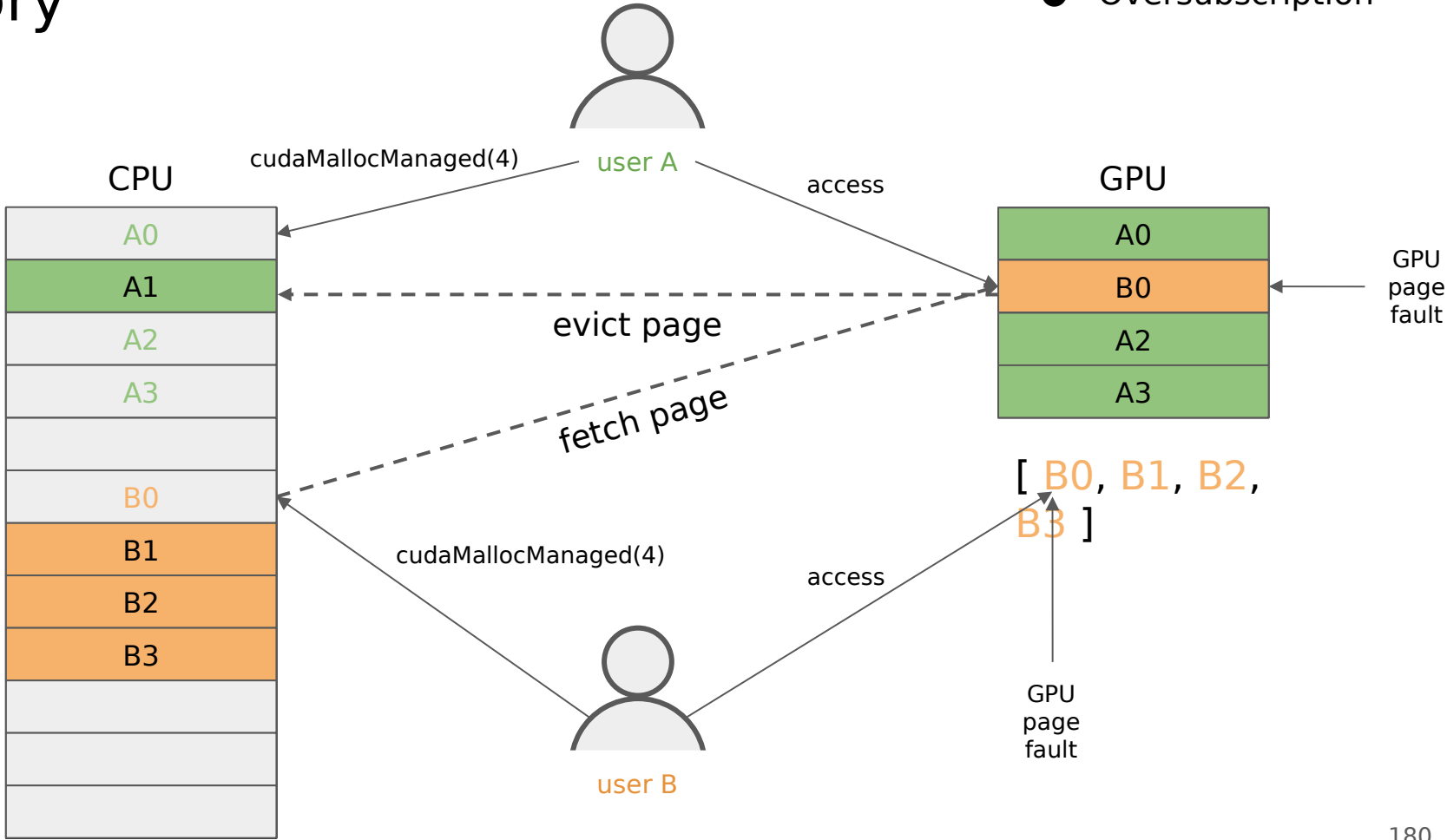
# Memory

● Oversubscription

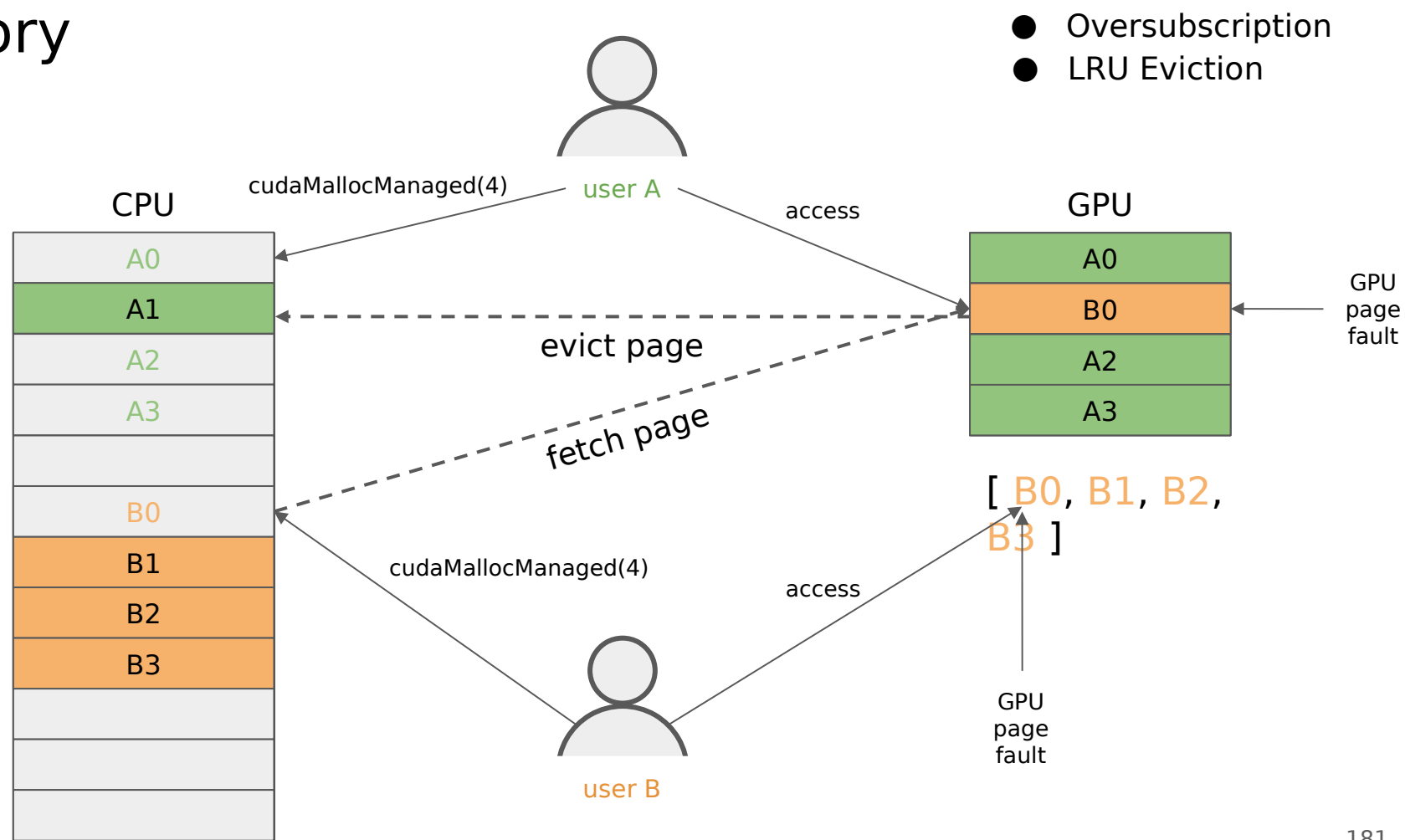


# Memory

● Oversubscription



# Memory



# Function wrappers

# Function wrappers

- *execve()* loads executable code into memory and passes control to *ld.so*

# Function wrappers

- *execve()* loads executable code into memory and passes control to *ld.so*
- *ld.so* maps shared objects to program address space and **resolves symbols**



# Function wrappers

- *execve()* loads executable code into memory and passes control to *ld.so*
- *ld.so* maps shared objects to program address space and **resolves symbols**
- application begins execution

# Function wrappers

- *execve()* loads executable code into memory and passes control to *ld.so*
- *ld.so* maps shared objects to program address space and **resolves symbols**
- application begins execution

```
$ nm -D ./matrixMul
```

```
U cudaLaunchKernel  
U cudaMalloc  
U cudaMemcpy
```

# Function wrappers

- *execve()* loads executable code into memory and passes control to *ld.so*
- *ld.so* maps shared objects to program address space and **resolves symbols**
- application begins execution

print dynamic symbol table entries

```
$ nm -D /matrixMul
```

```
U cudaLaunchKernel  
U cudaMalloc  
U cudaMemcpy
```

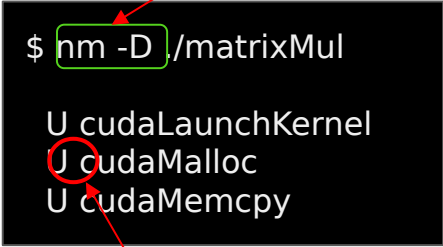
# Function wrappers

- *execve()* loads executable code into memory and passes control to *ld.so*
- *ld.so* maps shared objects to program address space and **resolves symbols**
- application begins execution

print dynamic symbol table entries

```
$ nm -D /matrixMul
U cudaLaunchKernel
U cudaMalloc
U cudaMemcpy
```

Undefined Symbol



# Function wrappers

- *execve()* loads executable code into memory and passes control to *ld.so*
- *ld.so* maps shared objects to program address space and **resolves symbols**
- application begins execution

print dynamic symbol table entries

```
$ nm -D /matrixMul
U cudaLaunchKernel
U cudaMalloc
U cudaMemcpy
```

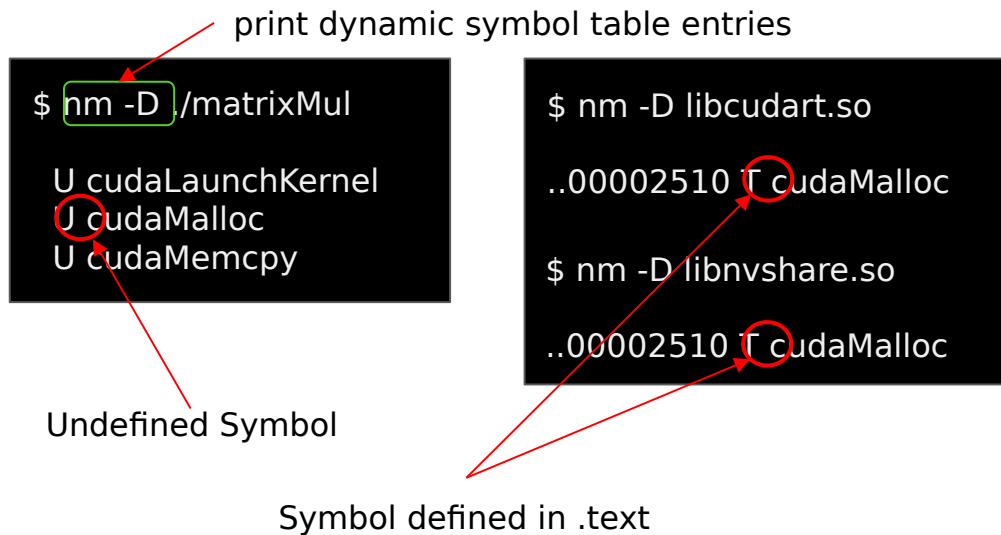
Undefined Symbol

```
$ nm -D libcudart.so
..00002510 T cudaMalloc

$ nm -D libnvshare.so
..00002510 T cudaMalloc
```

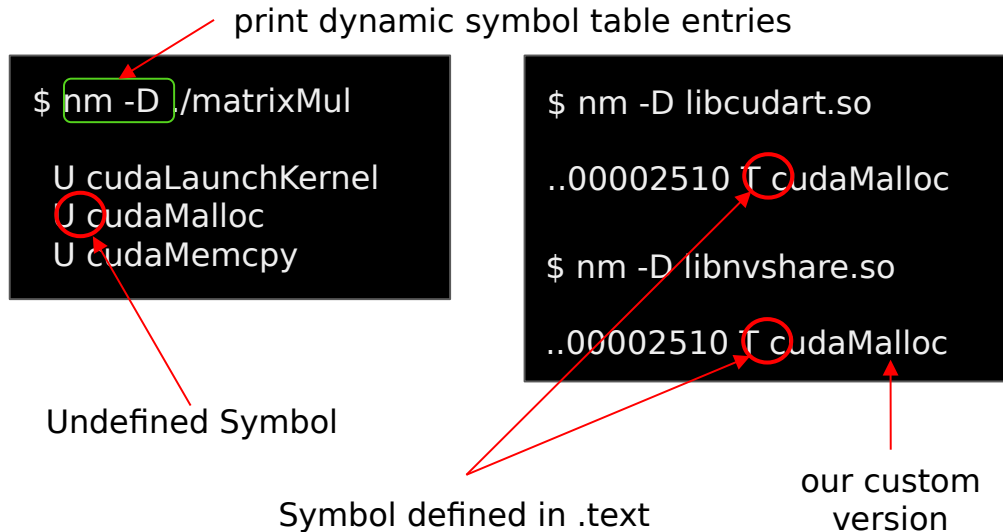
# Function wrappers

- `execve()` loads executable code into memory and passes control to `ld.so`
- `ld.so` maps shared objects to program address space and **resolves symbols**
- application begins execution



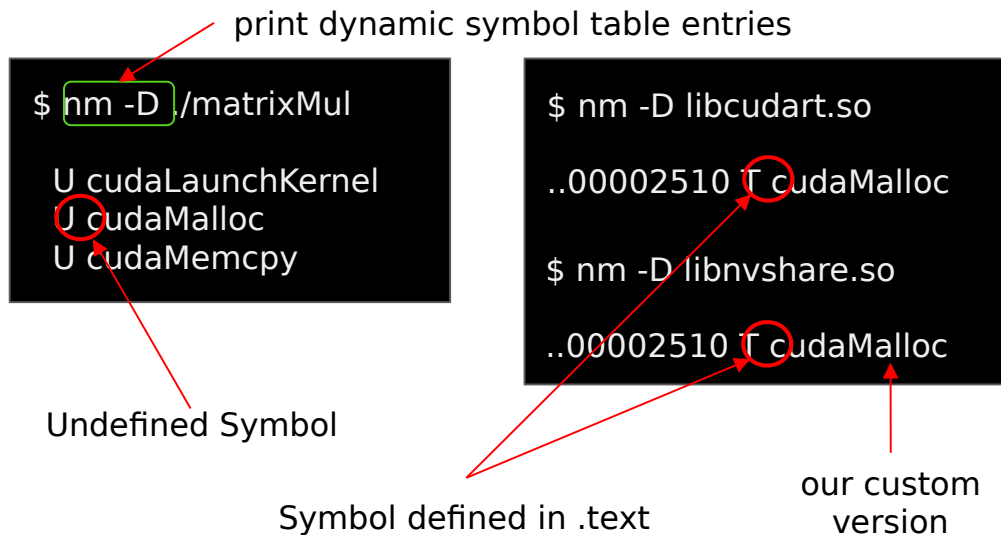
# Function wrappers

- *execve()* loads executable code into memory and passes control to *ld.so*
- *ld.so* maps shared objects to program address space and **resolves symbols**
- application begins execution



# Function wrappers

- *execve()* loads executable code into memory and passes control to *ld.so*
- *ld.so* maps shared objects to program address space and **resolves symbols**
- application begins execution

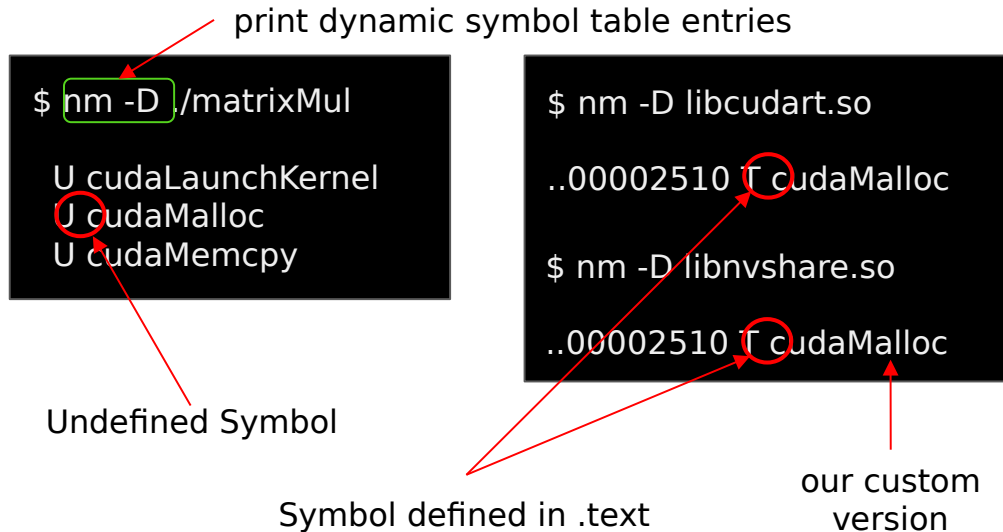




# Function wrappers

- `execve()` loads executable code into memory and passes control to `ld.so`
- `ld.so` maps shared objects to program address space and **resolves symbols**
- application begins execution

Normally:



# Function wrappers

- `execve()` loads executable code into memory and passes control to `ld.so`
- `ld.so` maps shared objects to program address space and **resolves symbols**
- application begins execution

Normally:



print dynamic symbol table entries

```
$ nm -D /matrixMul
U cudaLaunchKernel
U cudaMalloc
U cudaMemcpy
```

Undefined Symbol

```
$ nm -D libcudart.so
..00002510 T cudaMalloc

$ nm -D libnvshare.so
..00002510 T cudaMalloc
```

Symbol defined in .text

our custom version

# Function wrappers

- `execve()` loads executable code into memory and passes control to `ld.so`
- `ld.so` maps shared objects to program address space and **resolves symbols**
- application begins execution

print dynamic symbol table entries

```
$ nm -D /matrixMul
U cudaLaunchKernel
U cudaMalloc
U cudaMemcpy
```

Undefined Symbol

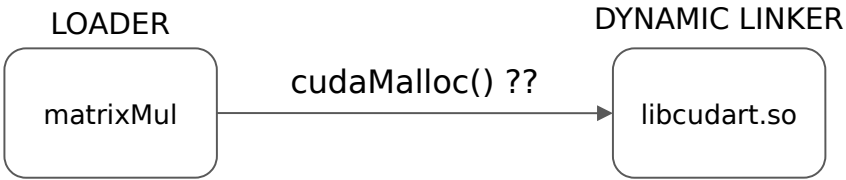
```
$ nm -D libcudart.so
..00002510 T cudaMalloc

$ nm -D libnvshare.so
..00002510 T cudaMalloc
```

Symbol defined in .text

our custom version

Normally:



# Function wrappers

- `execve()` loads executable code into memory and passes control to `ld.so`
- `ld.so` maps shared objects to program address space and **resolves symbols**
- application begins execution

print dynamic symbol table entries

```
$ nm -D /matrixMul
U cudaLaunchKernel
U cudaMalloc
U cudaMemcpy
```

Undefined Symbol

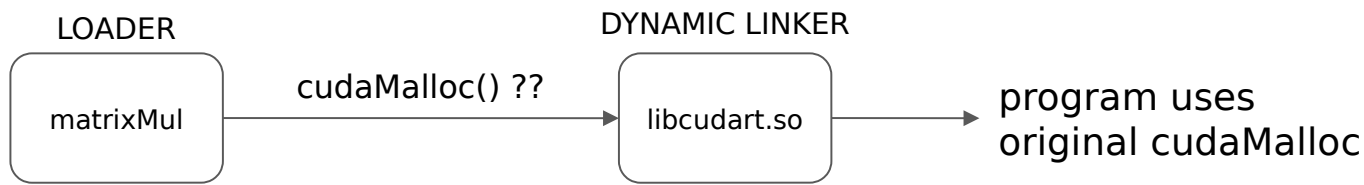
```
$ nm -D libcudart.so
..00002510 T cudaMalloc

$ nm -D libnvshare.so
..00002510 T cudaMalloc
```

Symbol defined in .text

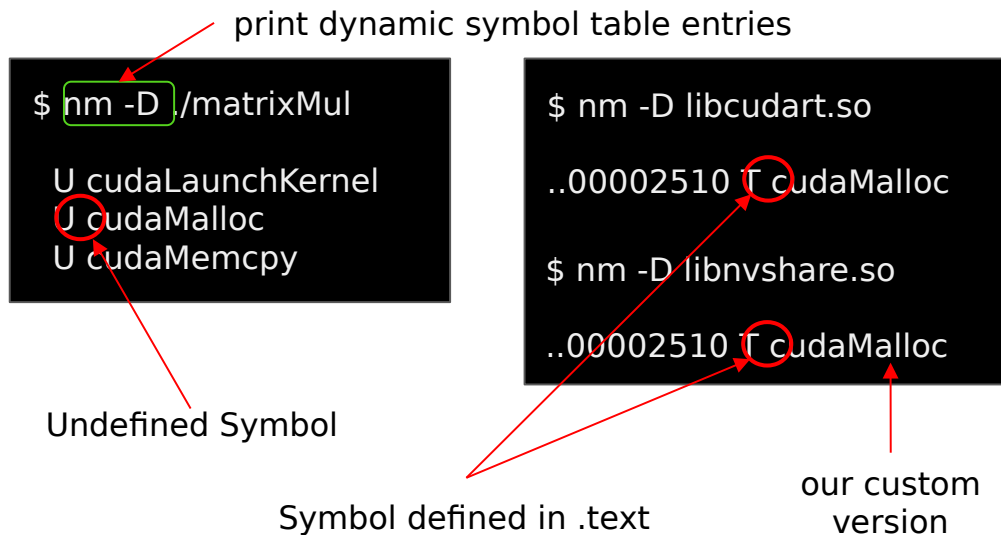
our custom version

Normally:

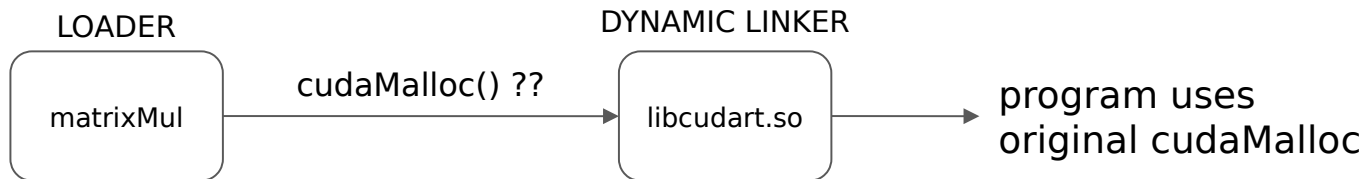


# Function wrappers

- `execve()` loads executable code into memory and passes control to `ld.so`
- `ld.so` maps shared objects to program address space and **resolves symbols**
- application begins execution



Normally:



# Function wrappers

- `execve()` loads executable code into memory and passes control to `ld.so`
- `ld.so` maps shared objects to program address space and **resolves symbols**
- application begins execution

print dynamic symbol table entries

```
$ nm -D /matrixMul
U cudaLaunchKernel
U cudaMalloc
U cudaMemcpy
```

Undefined Symbol

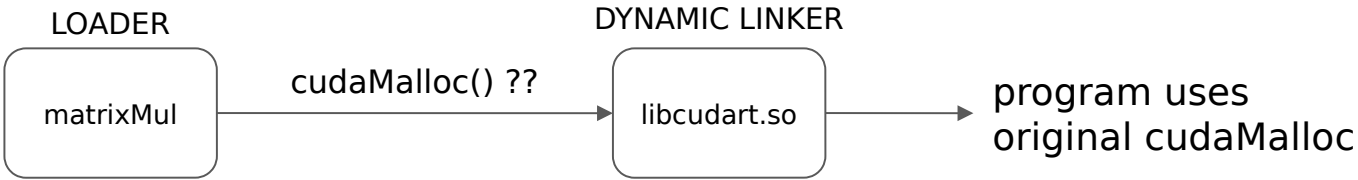
```
$ nm -D libcudart.so
..00002510 T cudaMalloc

$ nm -D libnvshare.so
..00002510 T cudaMalloc
```

Symbol defined in .text

our custom version

Normally:



LD\_PRELOAD=libnvshare.so

# Function wrappers

- `execve()` loads executable code into memory and passes control to `ld.so`
- `ld.so` maps shared objects to program address space and **resolves symbols**
- application begins execution

print dynamic symbol table entries

```
$ nm -D /matrixMul
```

U cudaLaunchKernel  
U **cudaMalloc**  
U cudaMemcpy

Undefined Symbol

```
$ nm -D libcudart.so
```

..00002510 T **cudaMalloc**

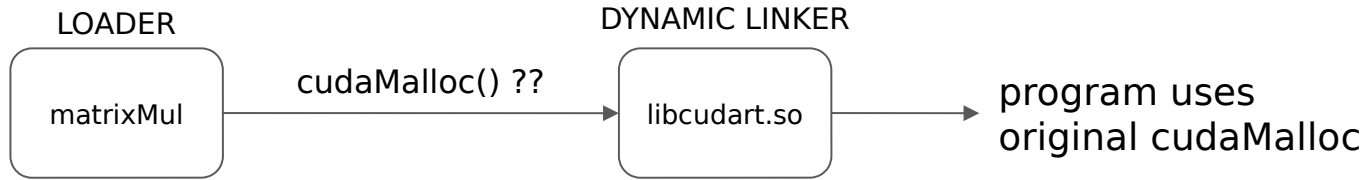
```
$ nm -D libnvshare.so
```

..00002510 T **cudaMalloc**

Symbol defined in .text

our custom version

Normally:



LD\_PRELOAD=libnvshare.so



# Function wrappers

- `execve()` loads executable code into memory and passes control to `ld.so`
- `ld.so` maps shared objects to program address space and **resolves symbols**
- application begins execution

print dynamic symbol table entries

```
$ nm -D /matrixMul
```

U cudaLaunchKernel  
J **cudaMalloc**  
U cudaMemcpy

Undefined Symbol

```
$ nm -D libcudart.so
```

..00002510 T **cudaMalloc**

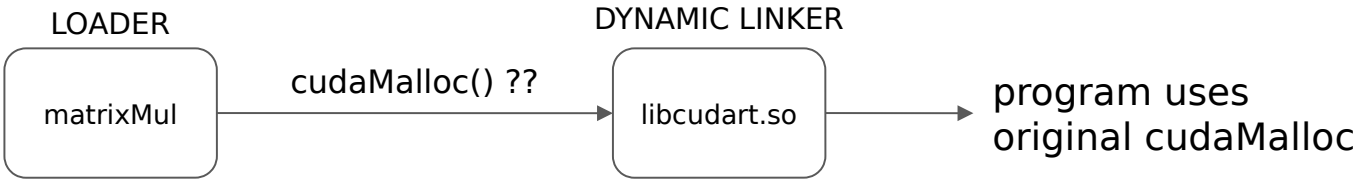
```
$ nm -D libnvshare.so
```

..00002510 T **cudaMalloc**

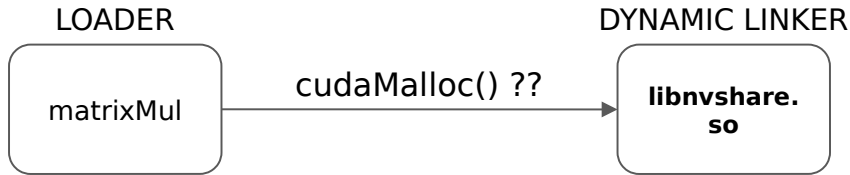
Symbol defined in .text

our custom version

Normally:



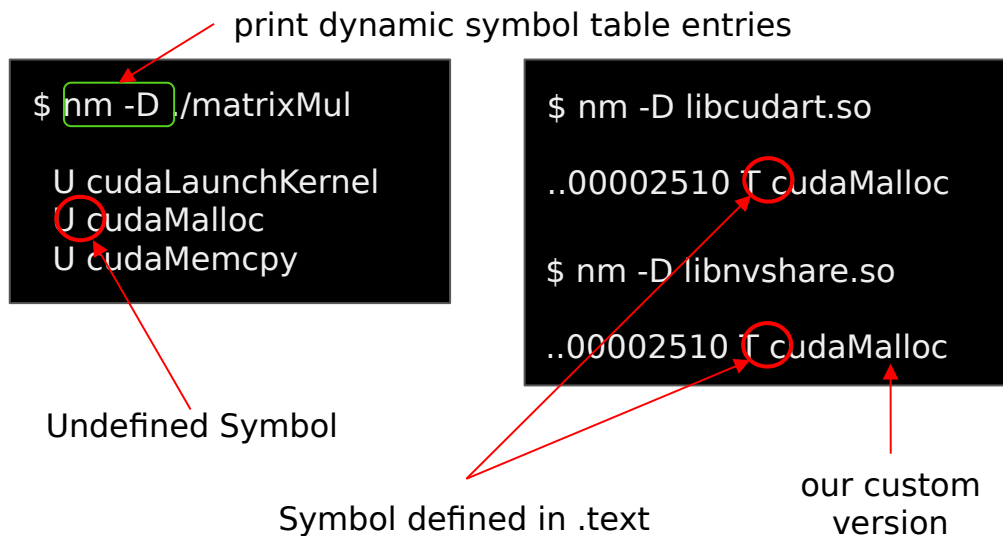
LD\_PRELOAD=libnvshare.so



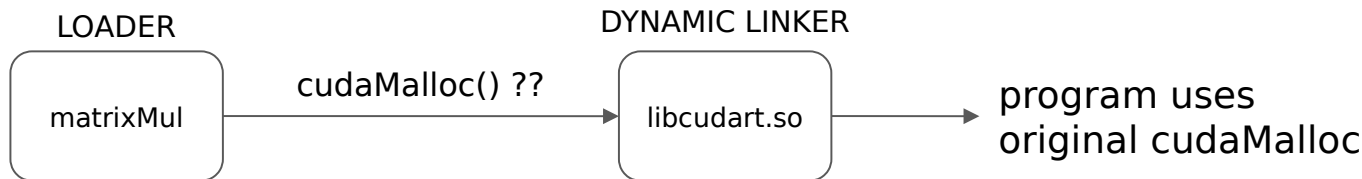


# Function wrappers

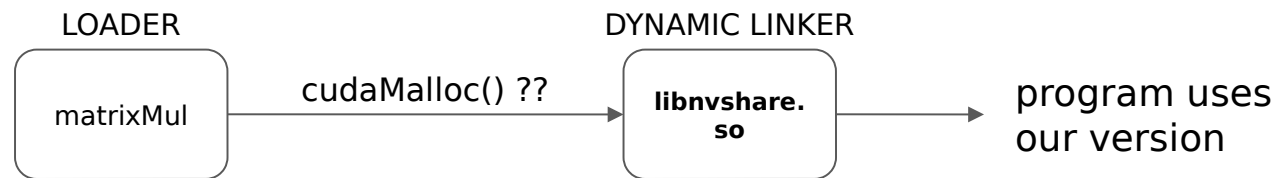
- `execve()` loads executable code into memory and passes control to `ld.so`
- `ld.so` maps shared objects to program address space and **resolves symbols**
- application begins execution



Normally:



LD\_PRELOAD=libnvshare.so



# Function wrappers

- `execve()` loads executable code into memory and passes control to `ld.so`
- `ld.so` maps shared objects to program address space and **resolves symbols**
- application begins execution

```
$ nm -D /matrixMul
U cudaLaunchKernel
U cudaMalloc
U cudaMemcpy
```

print dynamic symbol table entries

Undefined Symbol

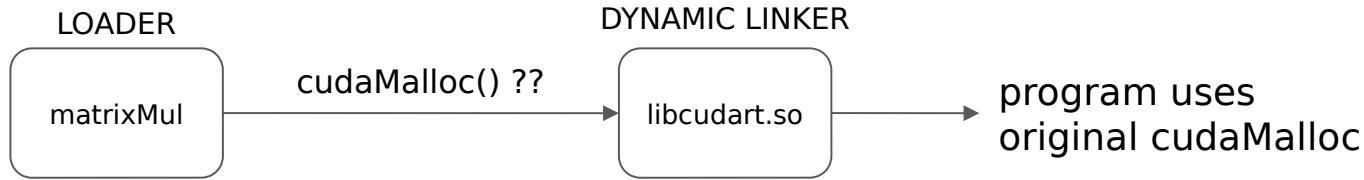
```
$ nm -D libcudart.so
..00002510 T cudaMalloc

$ nm -D libnvshare.so
..00002510 T cudaMalloc
```

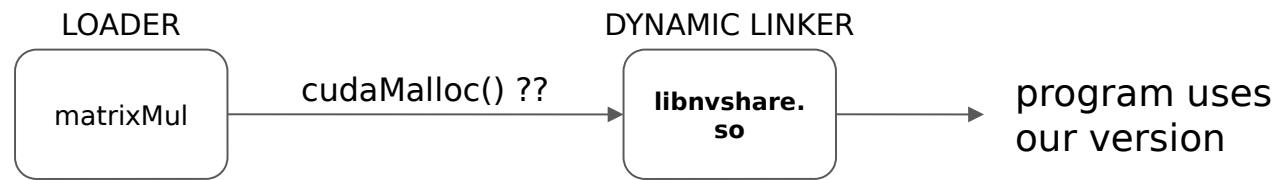
Symbol defined in .text

our custom version

Normally:



LD\_PRELOAD=libnvshare.so



- which can **internally** call the original

# Function wrappers

- `execve()` loads executable code into memory and passes control to `ld.so`
- `ld.so` maps shared objects to program address space and **resolves symbols**
- application begins execution

```
$ nm -D /matrixMul
U cudaLaunchKernel
U cudaMalloc
U cudaMemcpy
```

print dynamic symbol table entries

Undefined Symbol

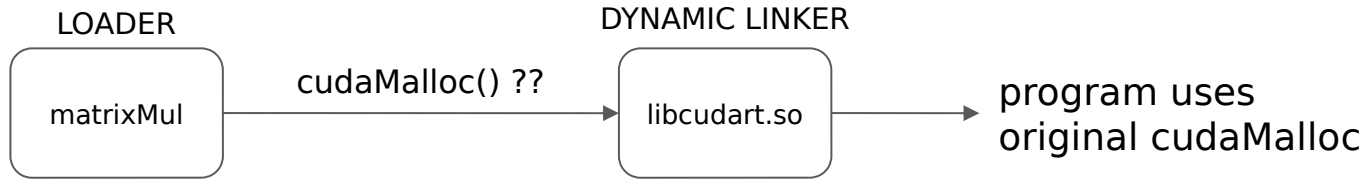
```
$ nm -D libcudart.so
..00002510 T cudaMalloc

$ nm -D libnvshare.so
..00002510 T cudaMalloc
```

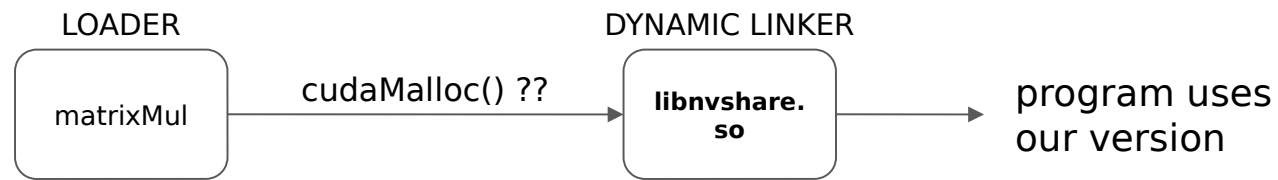
Symbol defined in .text

our custom version

Normally:



LD\_PRELOAD=libnvshare.so



- which can **internally** call the original
- ... if we will it

# Key Idea

# Key Idea



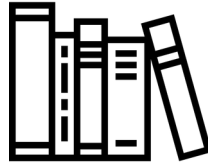
user

# Key Idea



user

libnvshare.so



# Key Idea

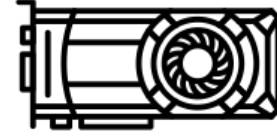


user

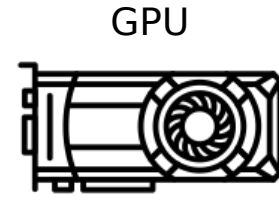
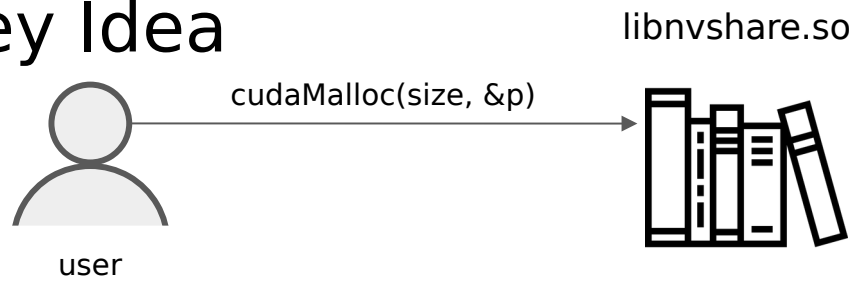
libnvshare.so



GPU

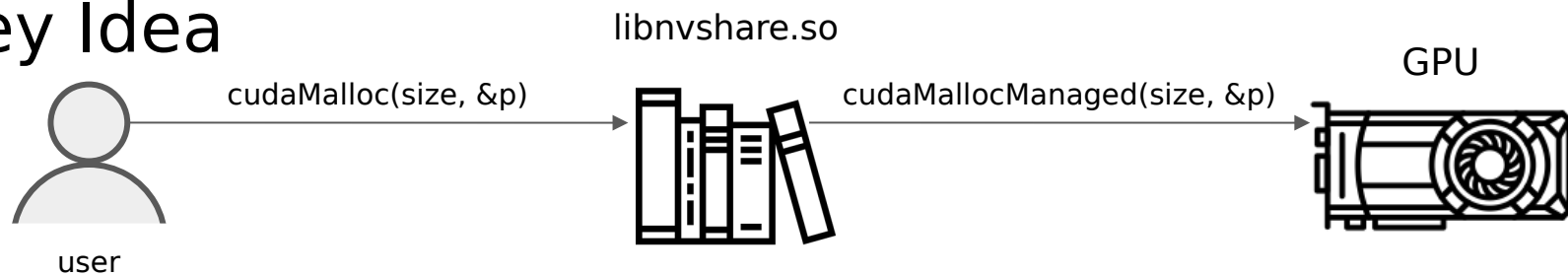


# Key Idea

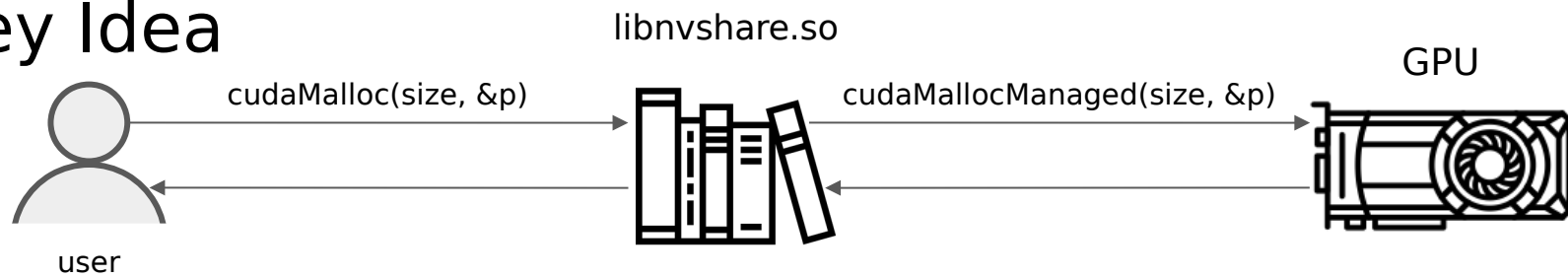




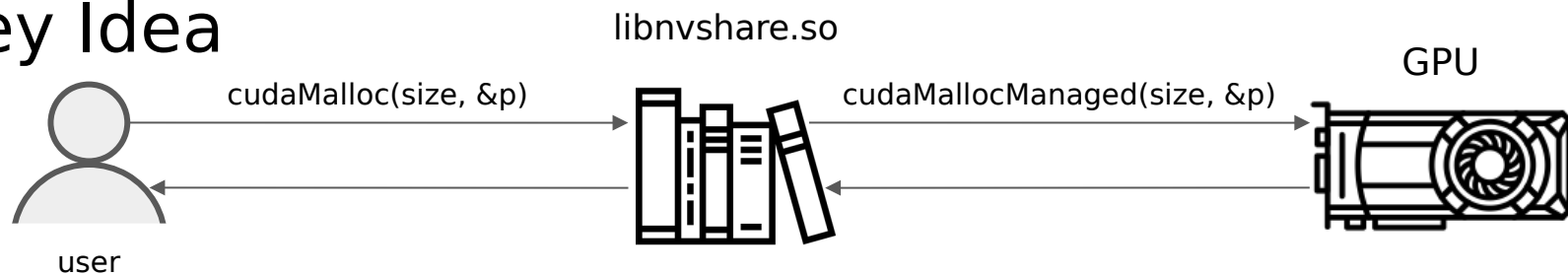
# Key Idea



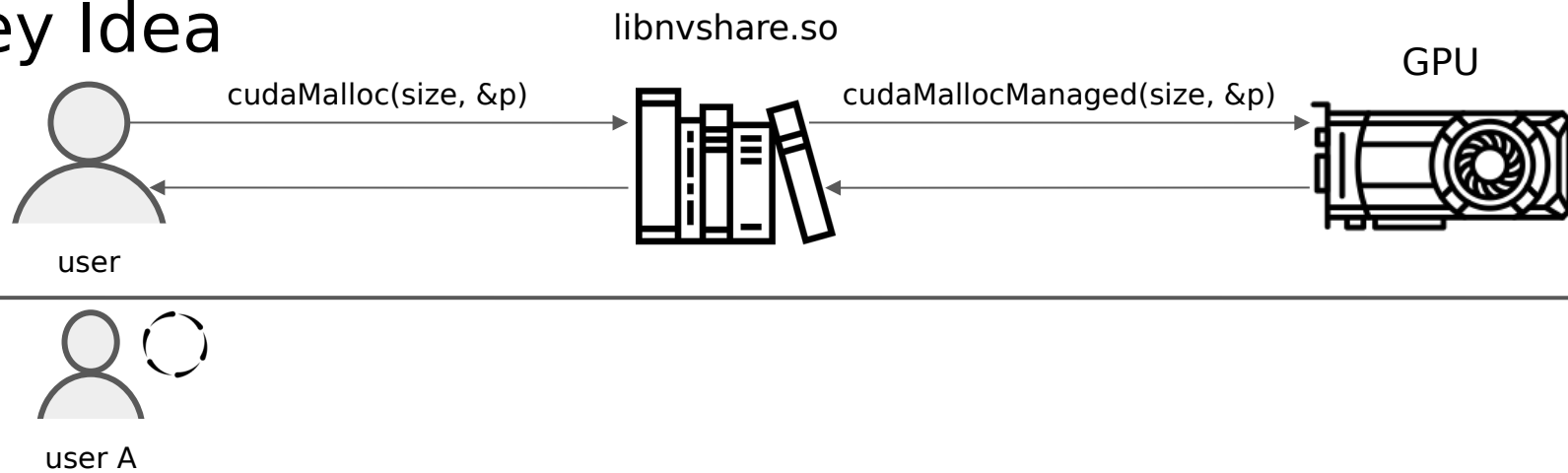
# Key Idea



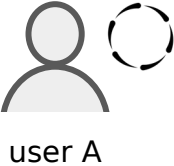
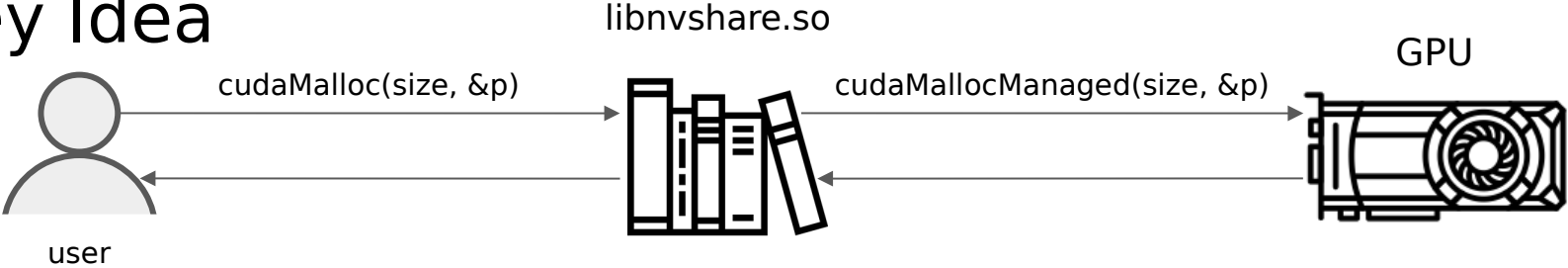
# Key Idea



# Key Idea



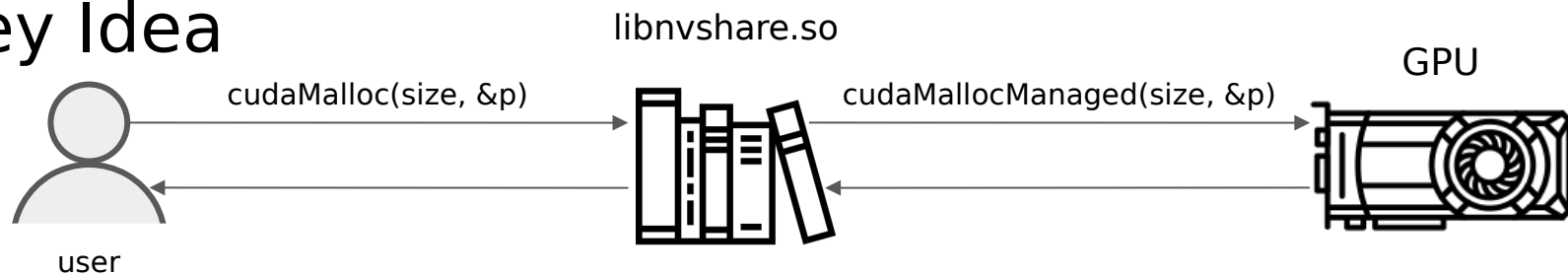
# Key Idea



GPU

A0
A1
A2
A3

# Key Idea



user A

GPU

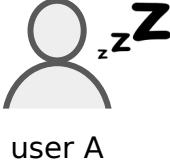
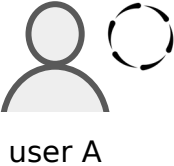
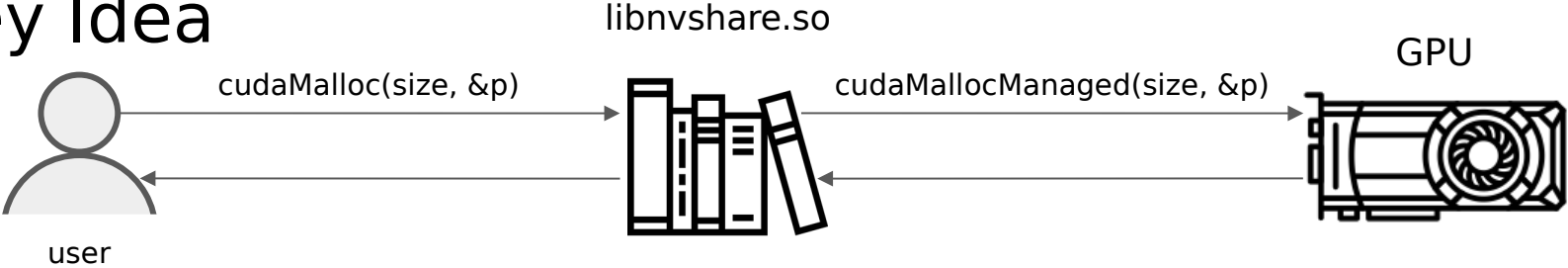
A0

A1

A2

A3

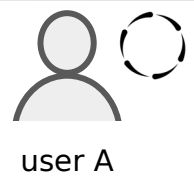
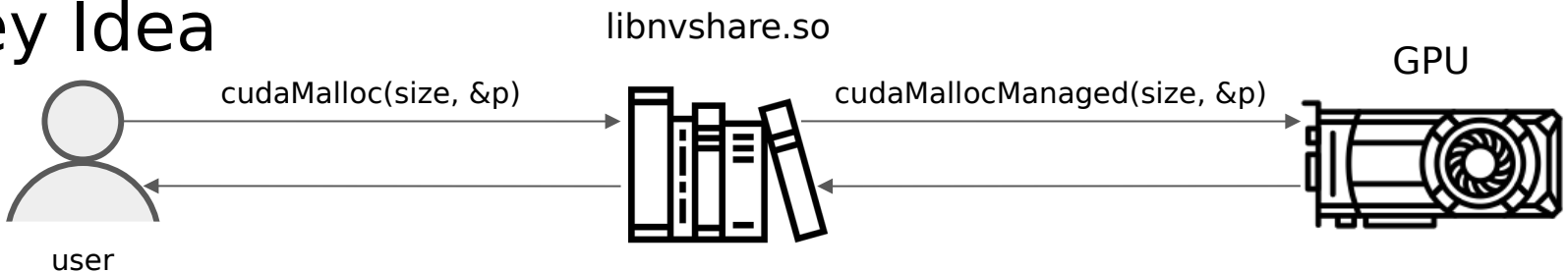
# Key Idea



GPU

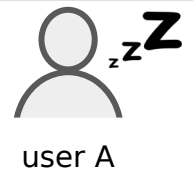
A0
A1
A2
A3

# Key Idea



GPU

A0
A1
A2
A3

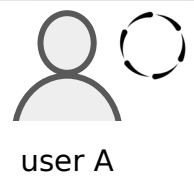
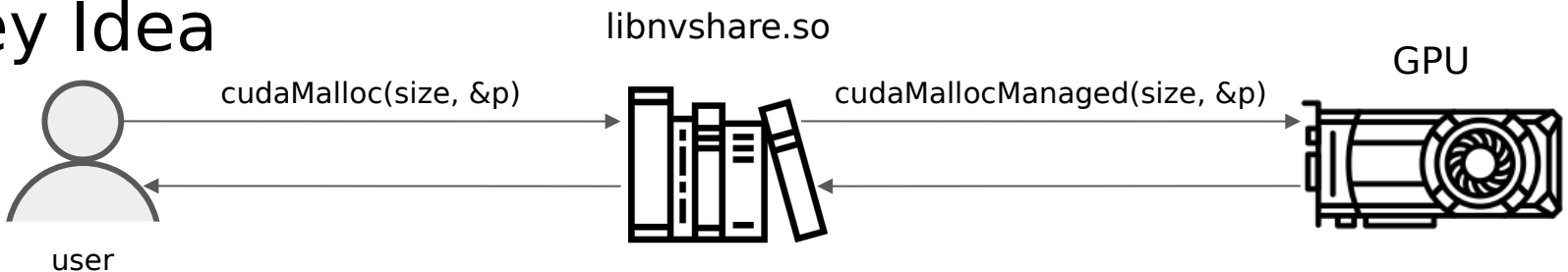


GPU

A0
A1
A2
A3

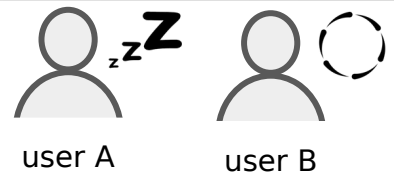


# Key Idea



GPU

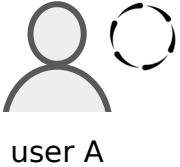
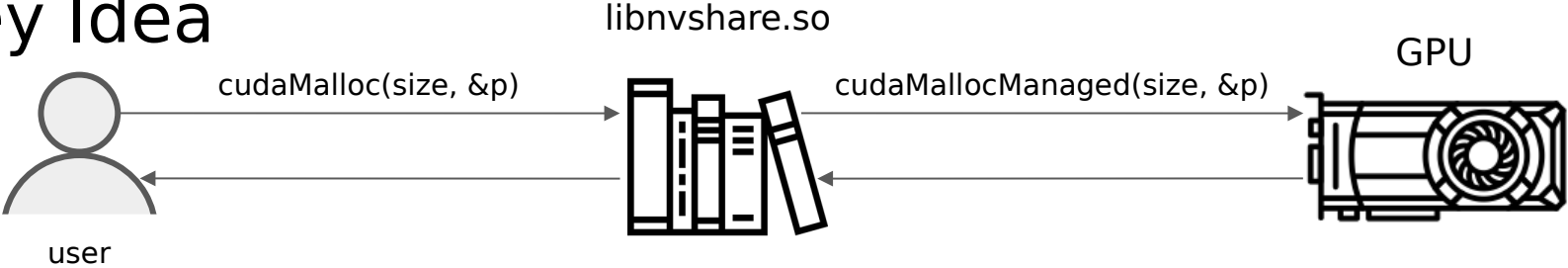
A0
A1
A2
A3



GPU

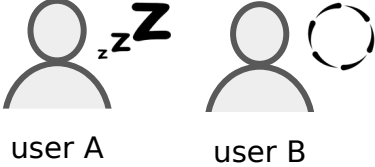
A0
A1
A2
A3

# Key Idea



GPU

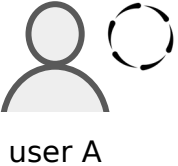
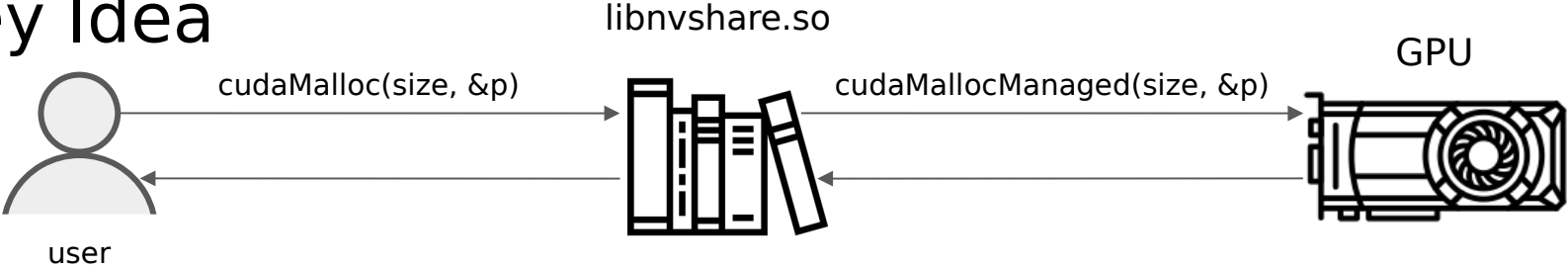
A0
A1
A2
A3



GPU

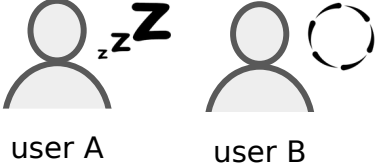
B0
A1
A2
A3

# Key Idea



GPU

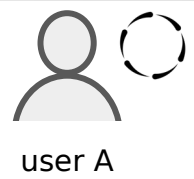
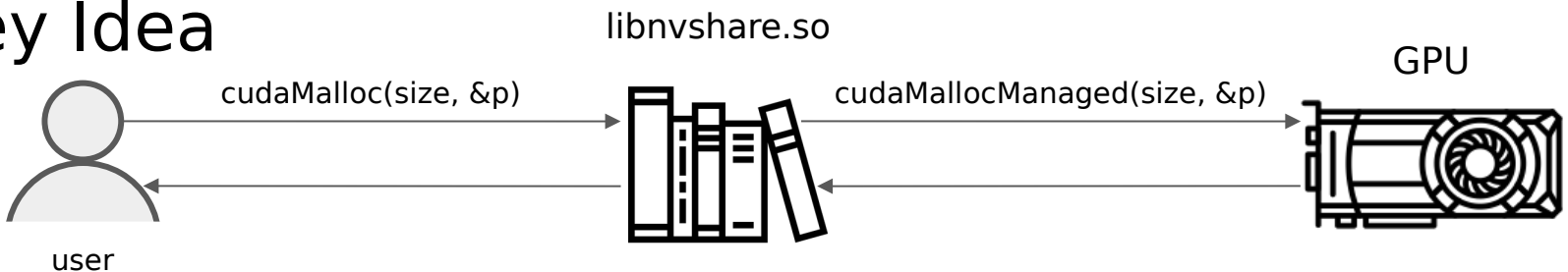
A0
A1
A2
A3



GPU

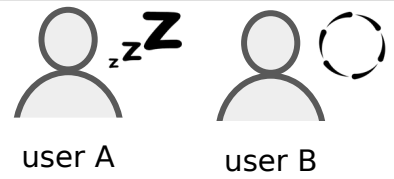
B0
B1
A2
A3

# Key Idea



GPU

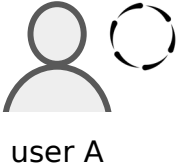
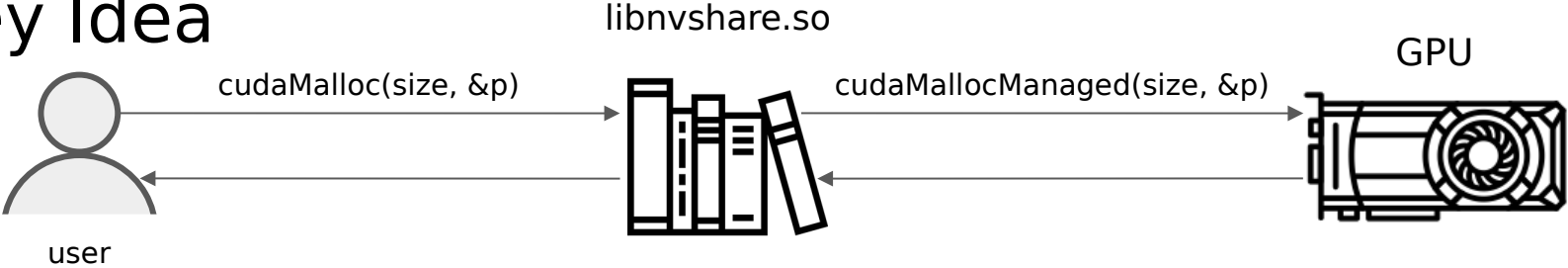
A0
A1
A2
A3



GPU

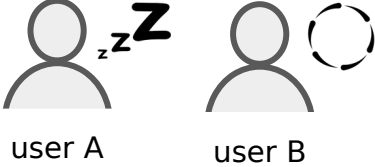
B0
B1
B2
A3

# Key Idea



GPU

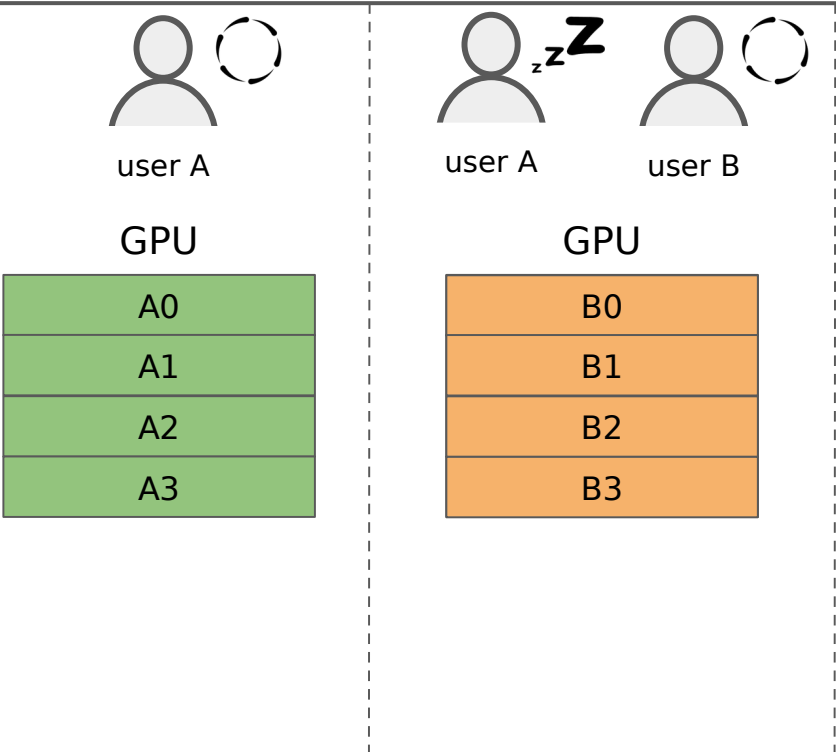
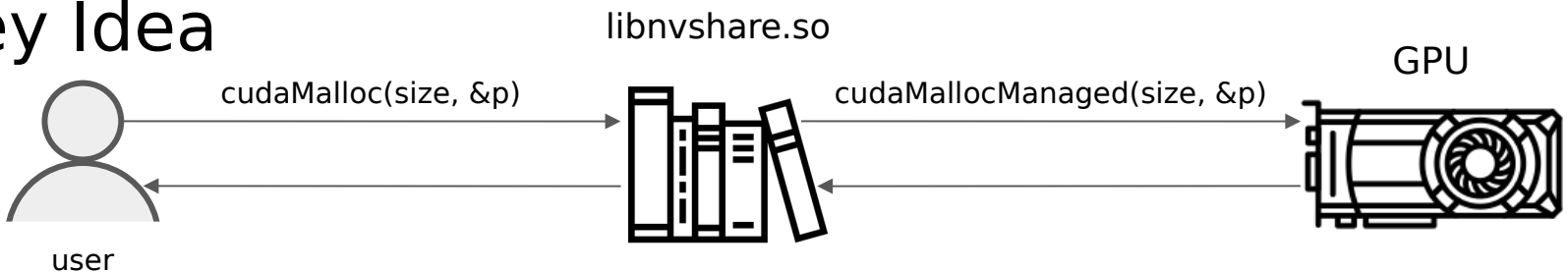
A0
A1
A2
A3



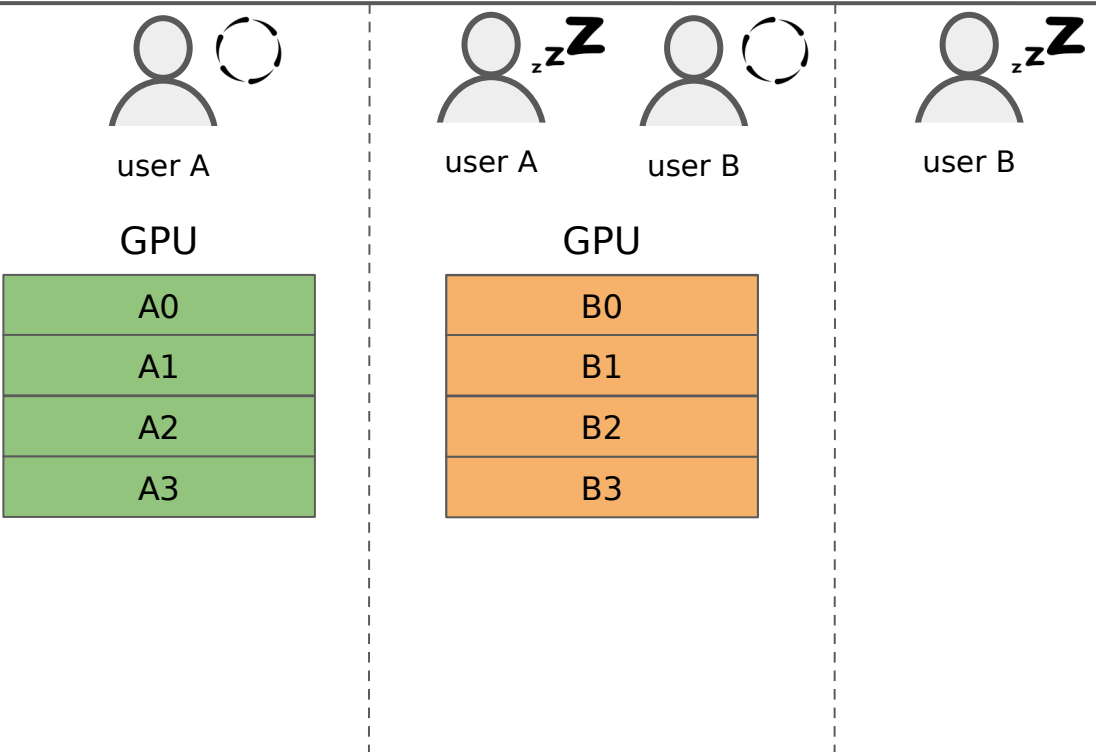
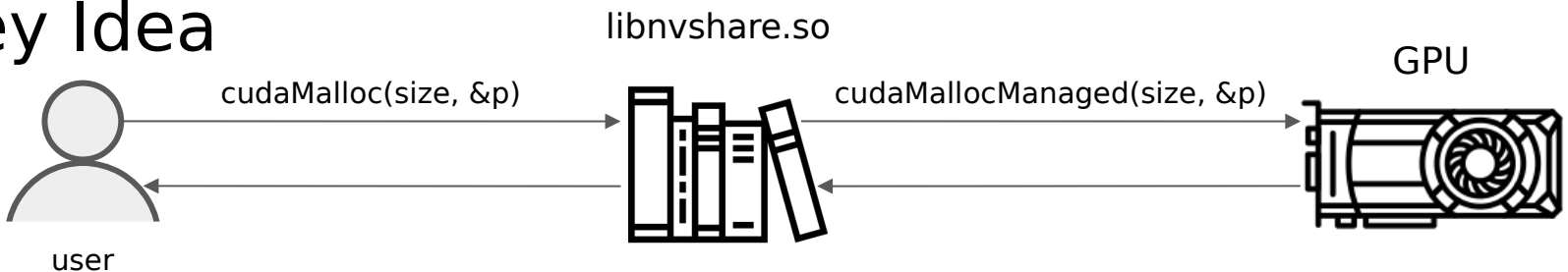
GPU

B0
B1
B2
B3

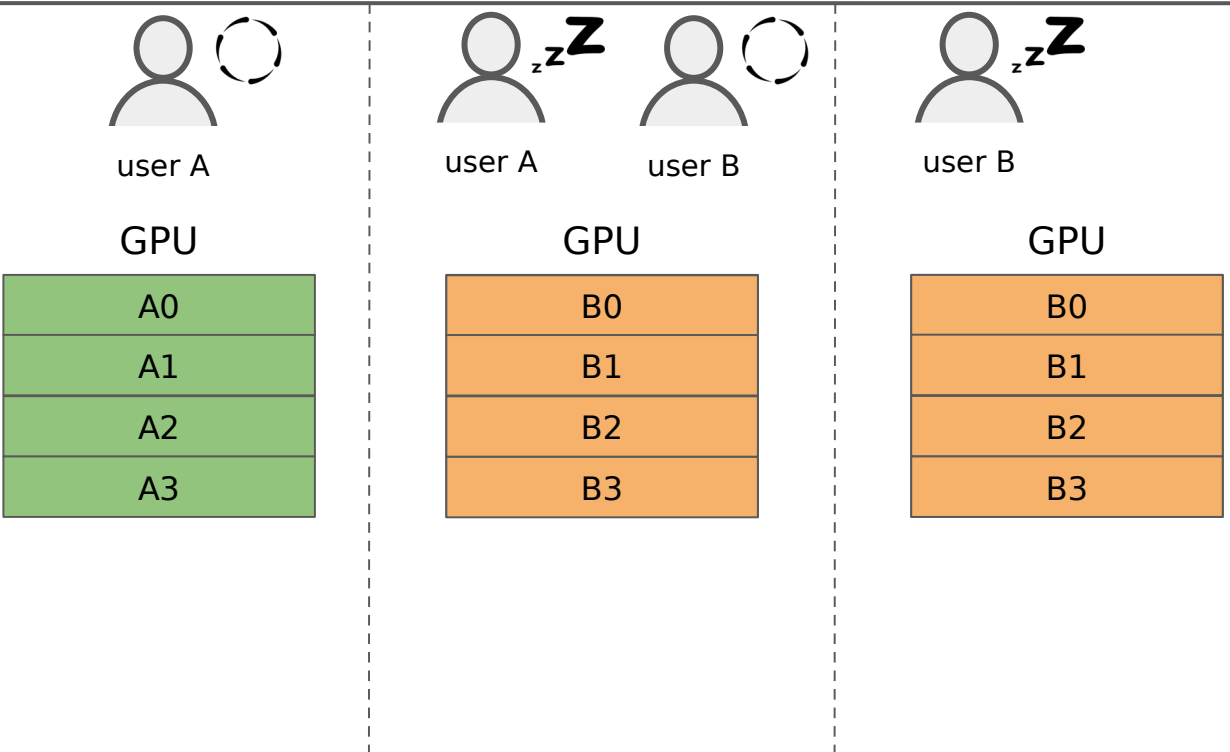
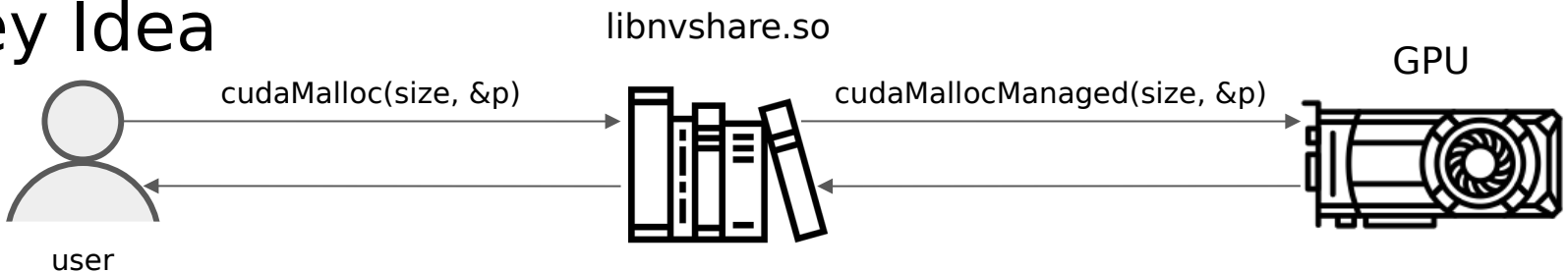
# Key Idea



# Key Idea

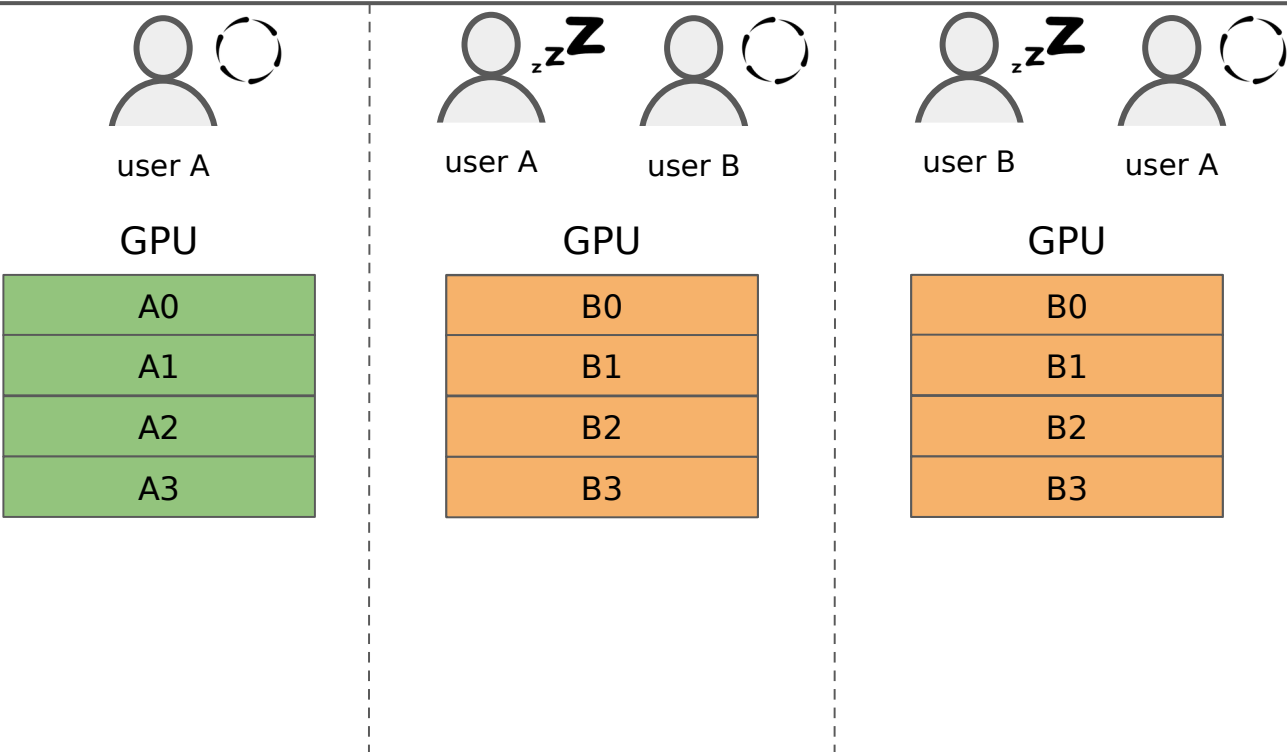
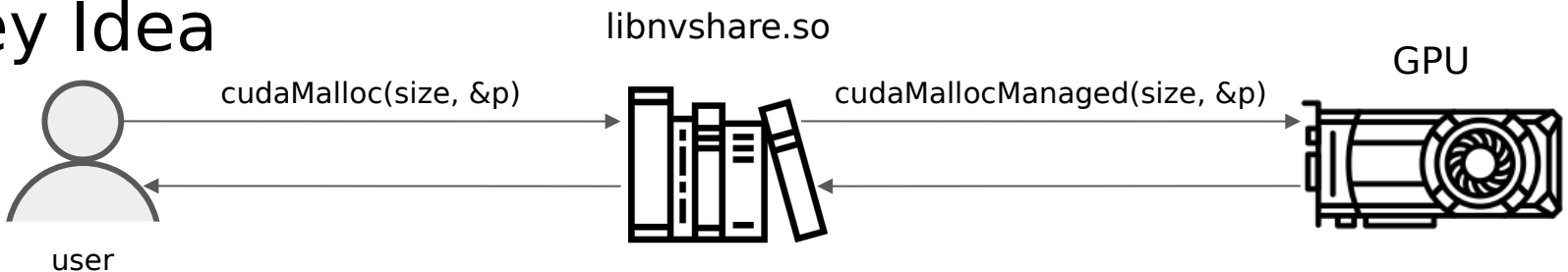


# Key Idea

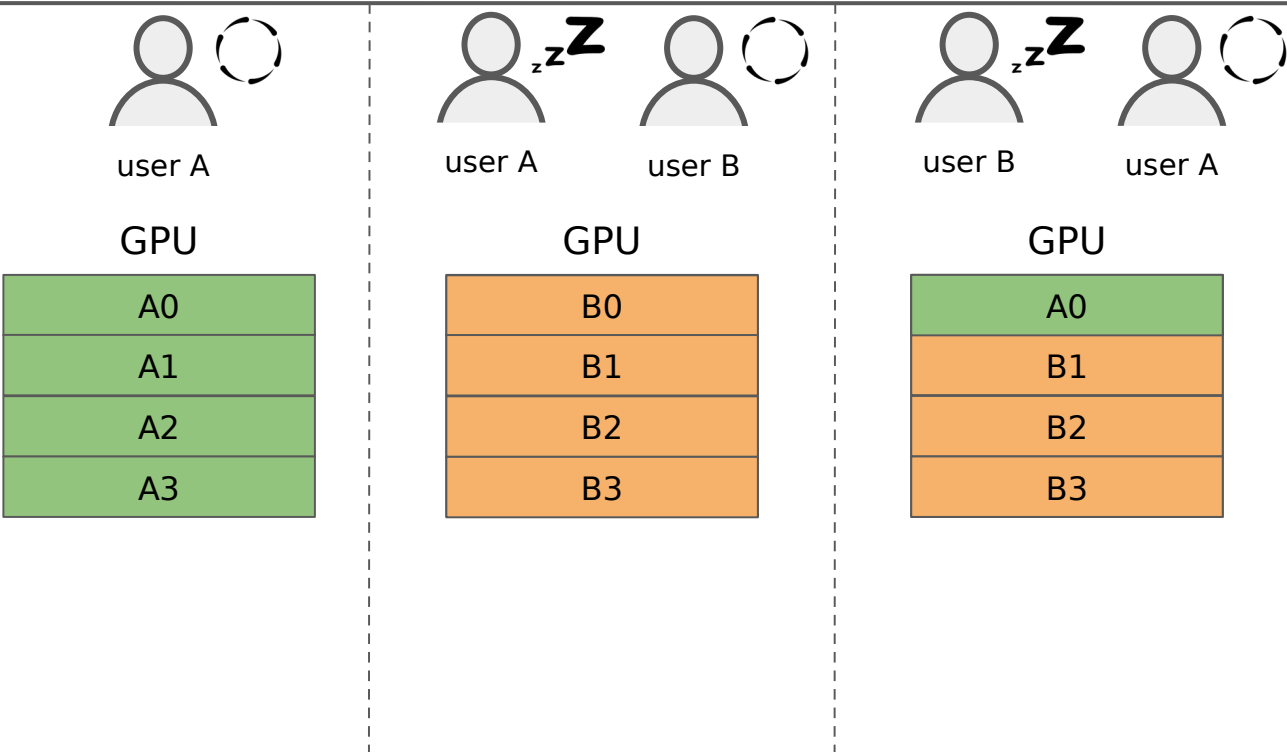
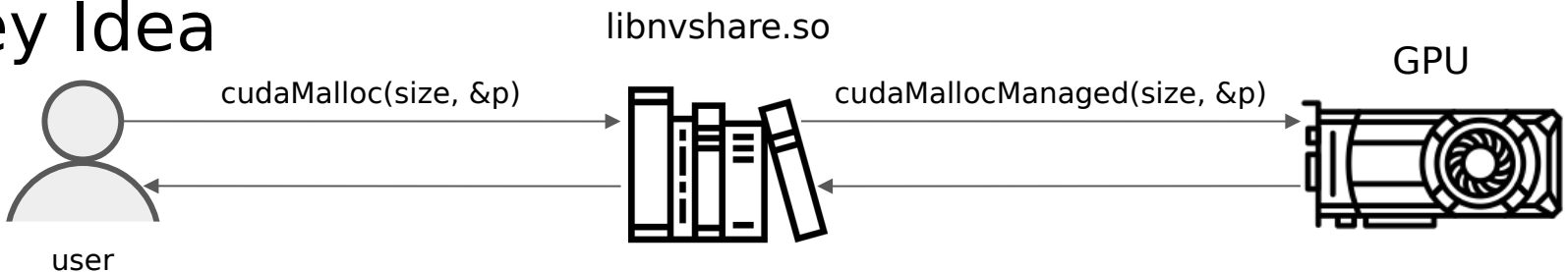




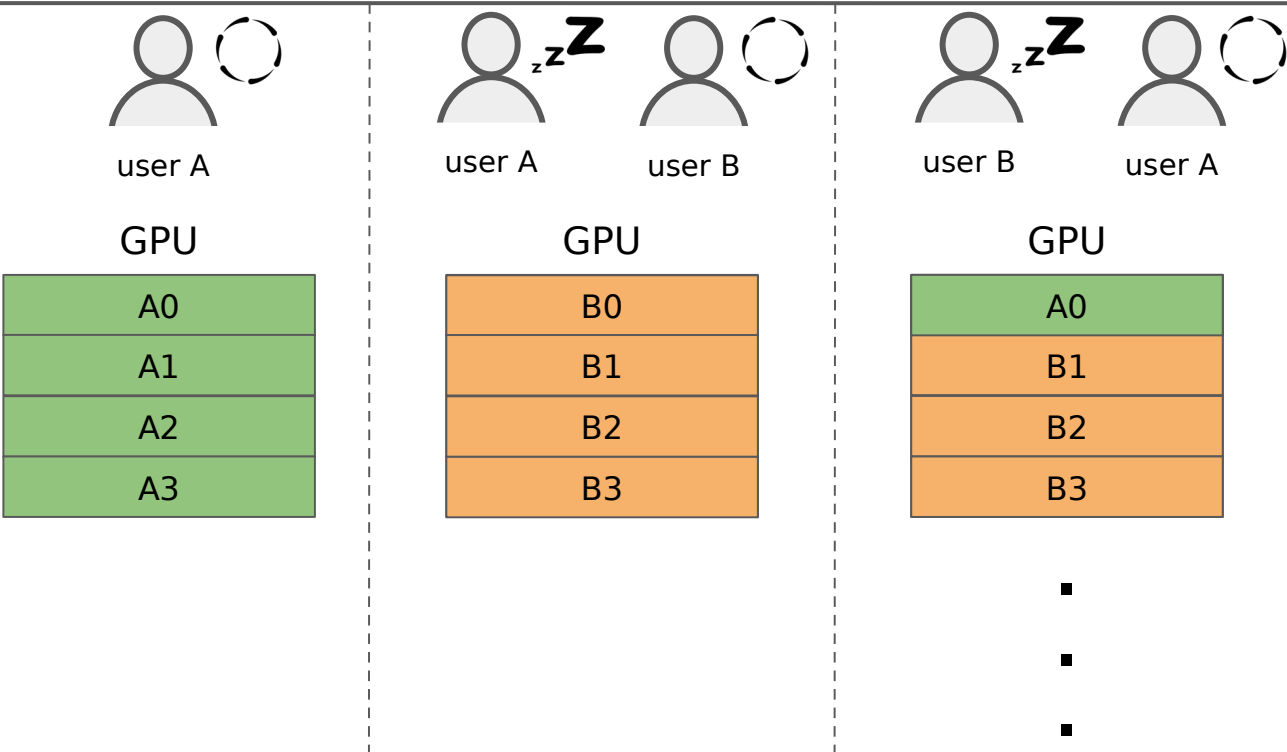
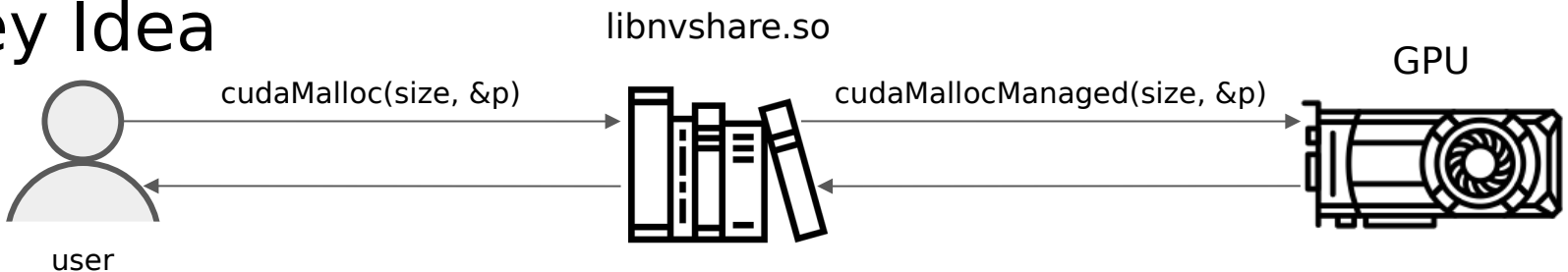
# Key Idea



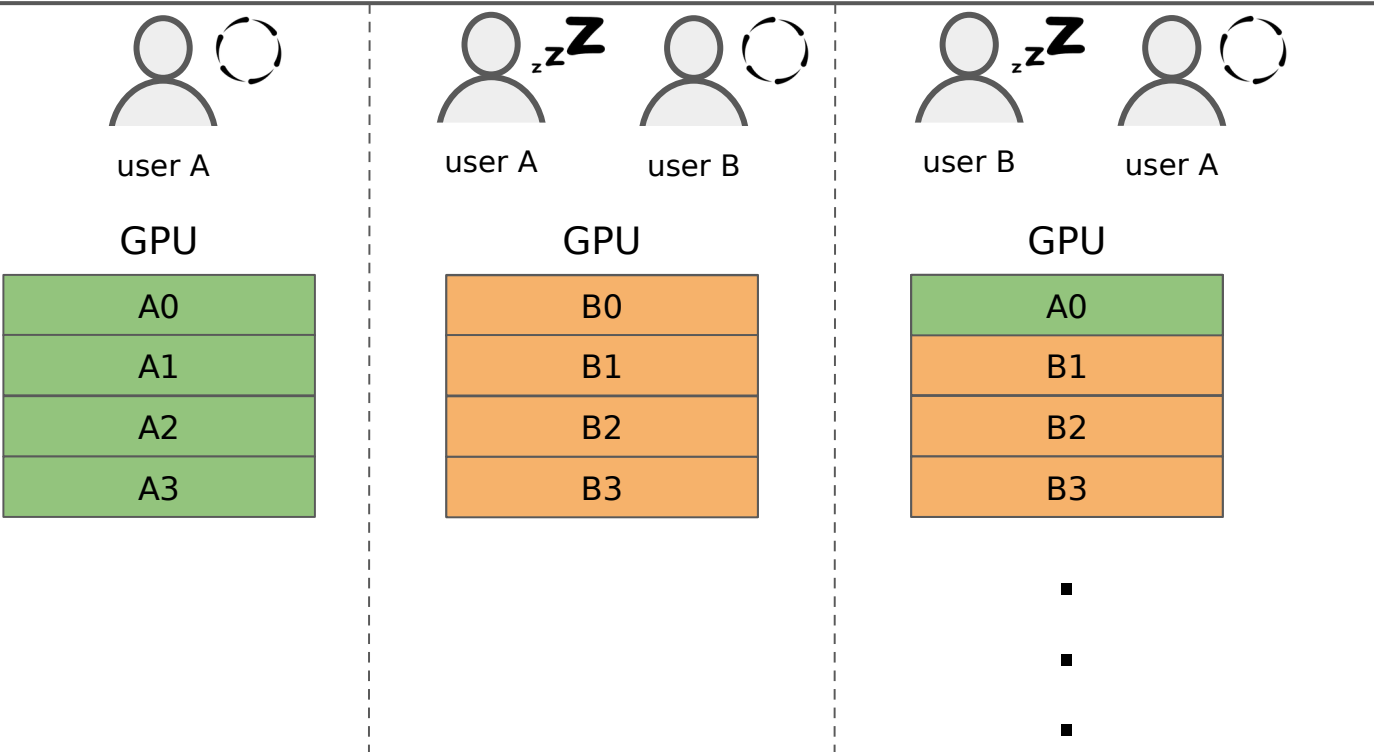
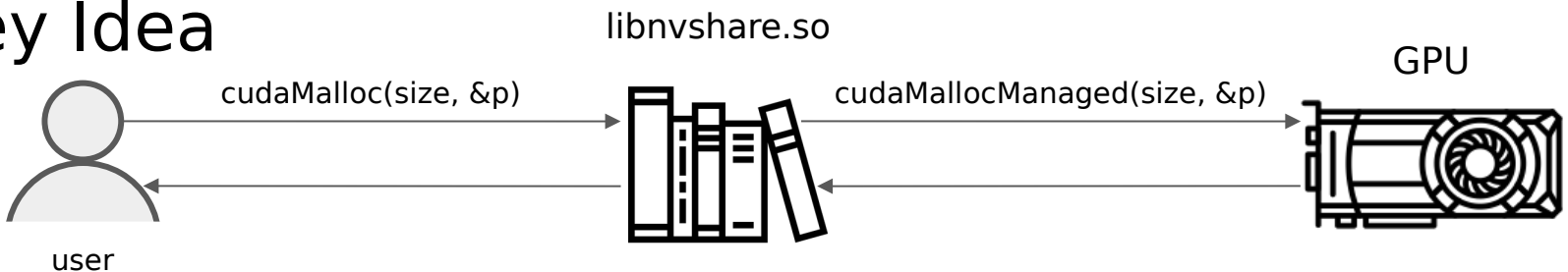
# Key Idea



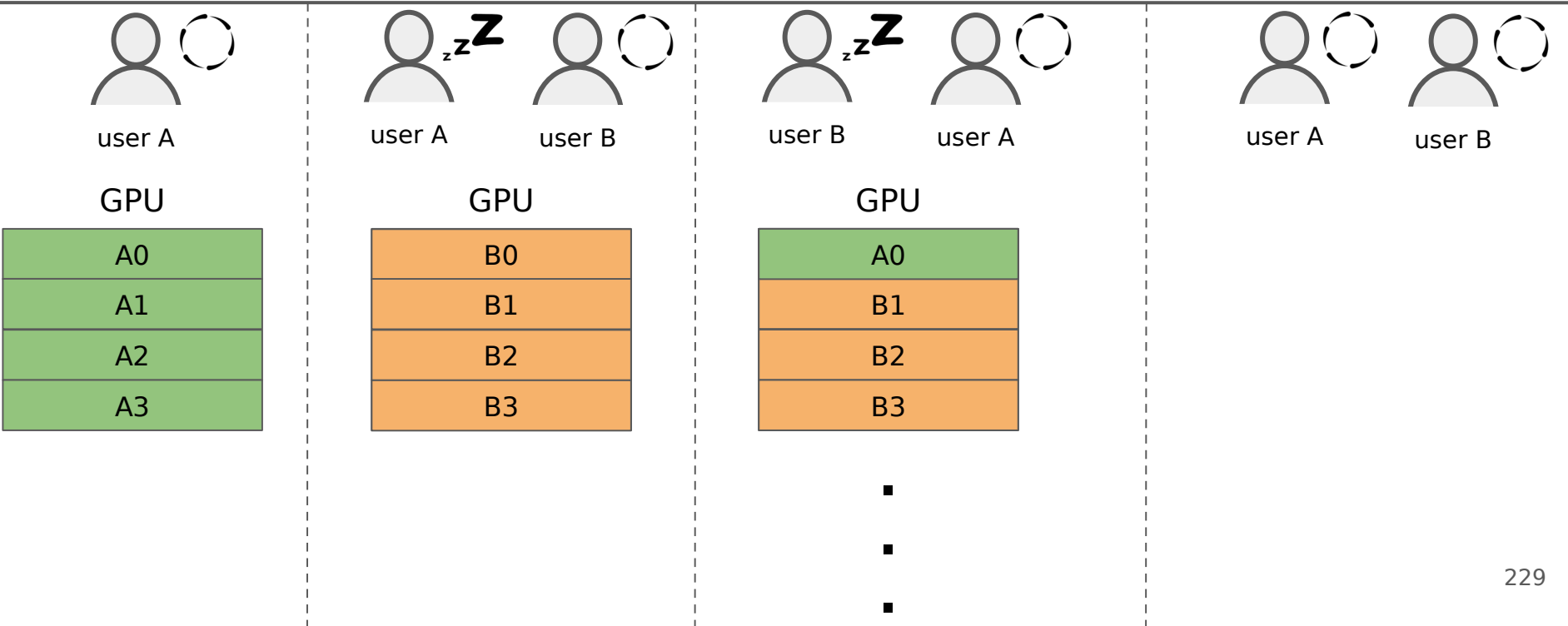
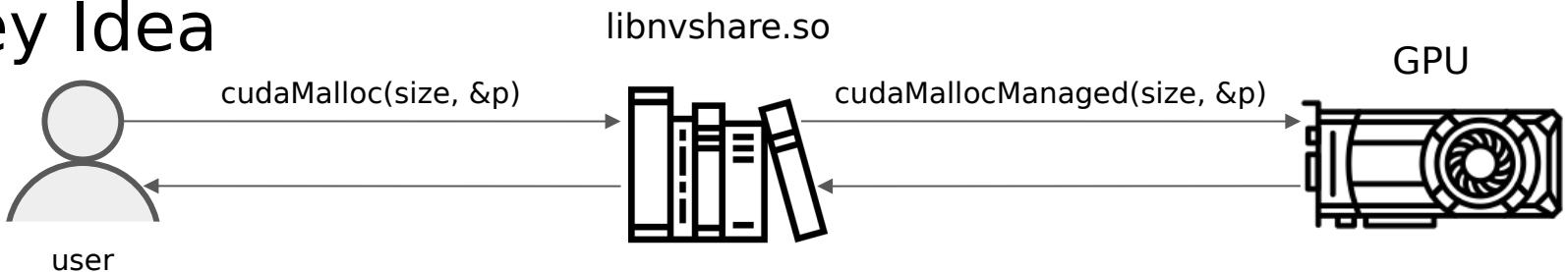
# Key Idea



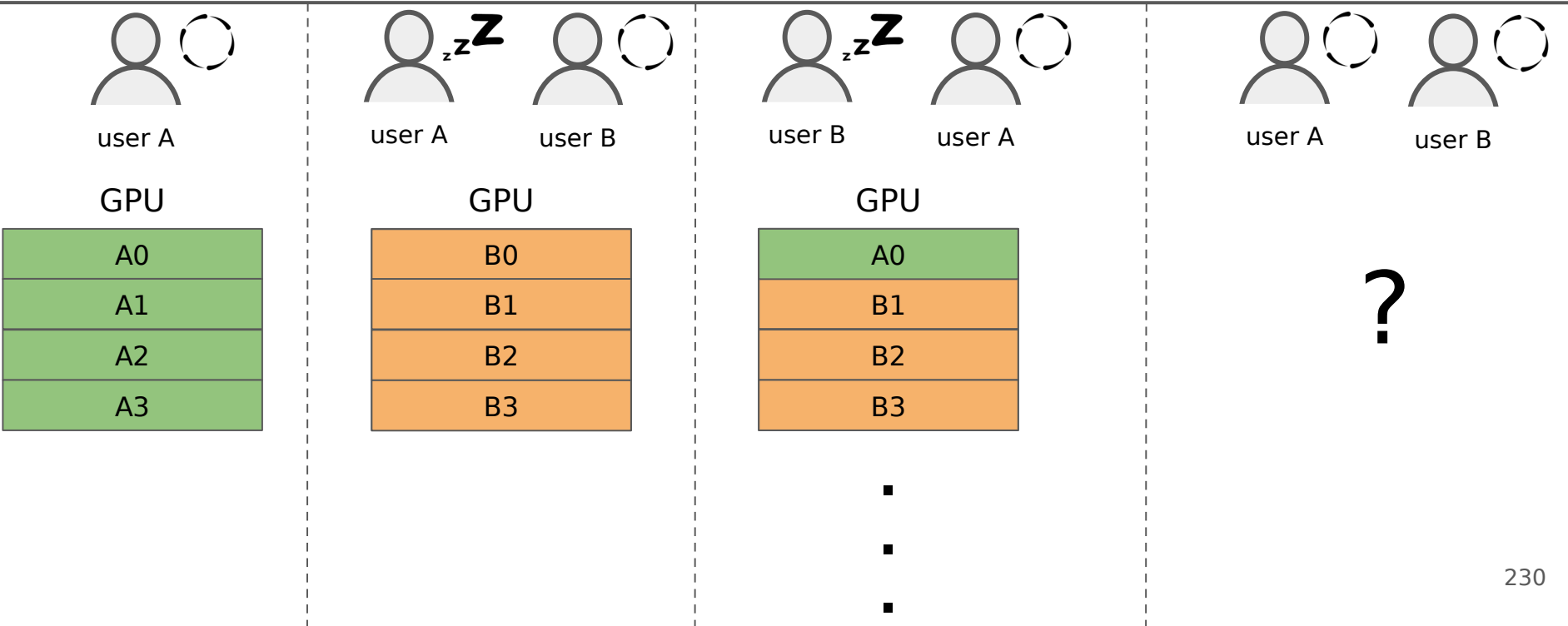
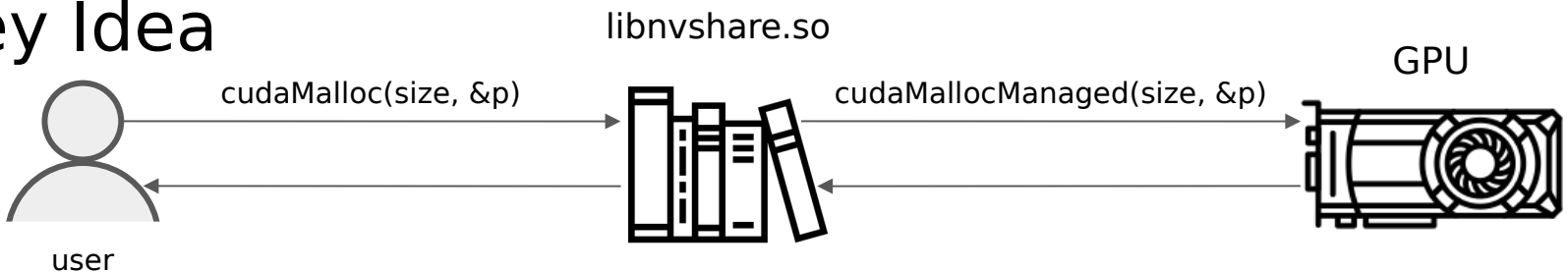
# Key Idea



# Key Idea



# Key Idea



# When GPU bursts overlap

# When GPU bursts overlap

---



# When GPU bursts overlap

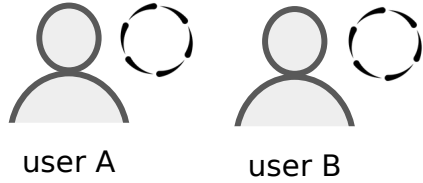
Sometimes

---

# When GPU bursts overlap

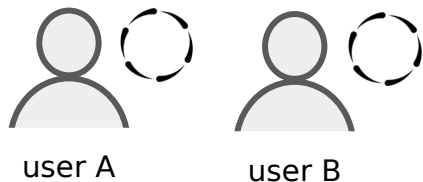
Sometimes

---

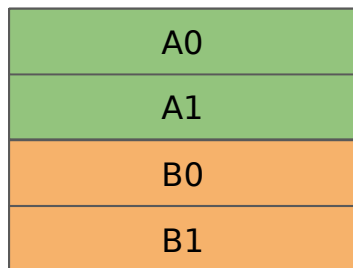


# When GPU bursts overlap

Sometimes

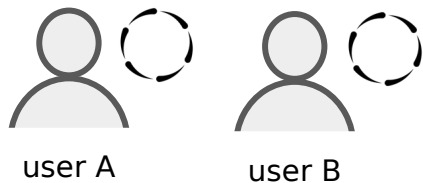


GPU

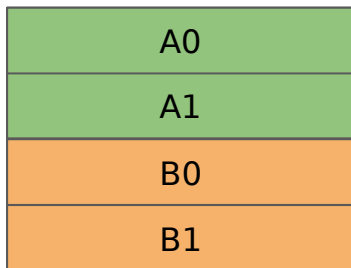


# When GPU bursts overlap

Sometimes



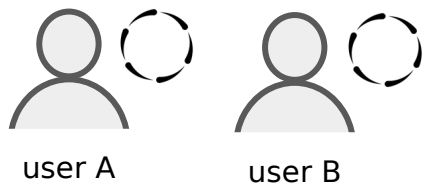
GPU



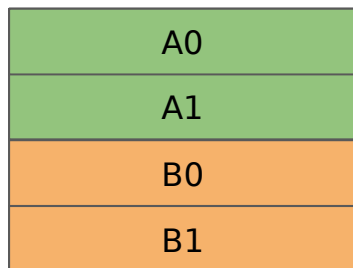
$$\sum W_{ss_i} < \text{GPU memory}$$

# When GPU bursts overlap

Sometimes



GPU

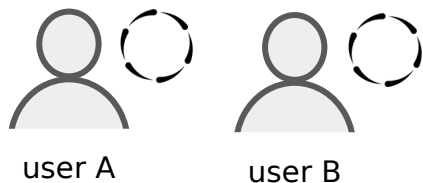


$$\sum W_{ss_i} < \text{GPU memory}$$

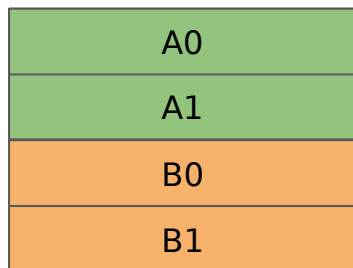
No page faults

# When GPU bursts overlap

Sometimes



GPU

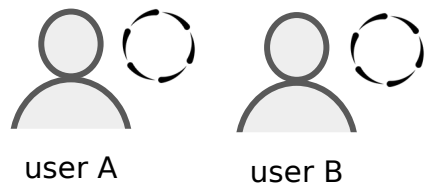


$$\sum W_{ss_i} < \text{GPU memory}$$

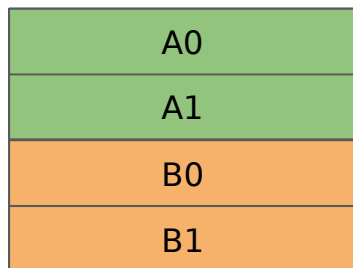
No page faults

# When GPU bursts overlap

Sometimes



GPU



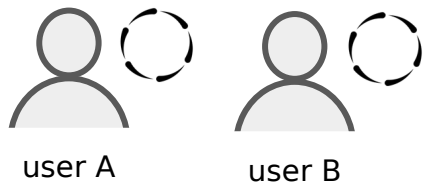
$\sum W_{ss_i} < \text{GPU memory}$

No page faults

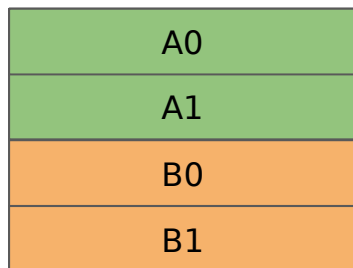
Other times

# When GPU bursts overlap

Sometimes



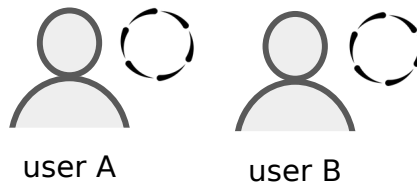
GPU



$\sum W_{ss_i} < \text{GPU memory}$

No page faults

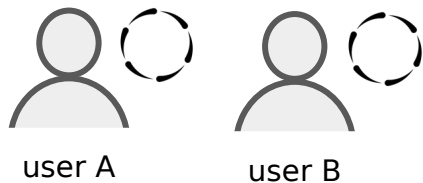
Other times



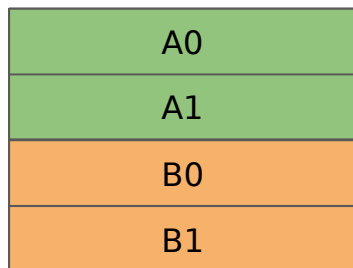


# When GPU bursts overlap

Sometimes



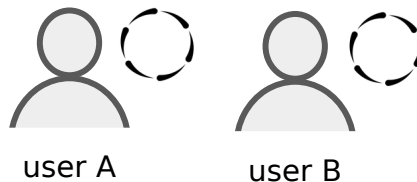
GPU



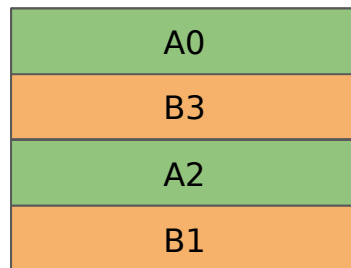
$$\sum W_{ss_i} < \text{GPU memory}$$

No page faults

Other times

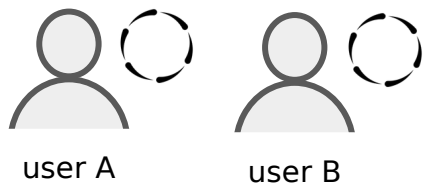


GPU

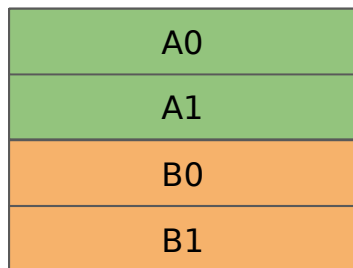


# When GPU bursts overlap

Sometimes



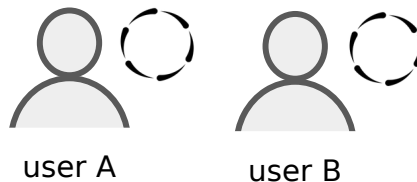
GPU



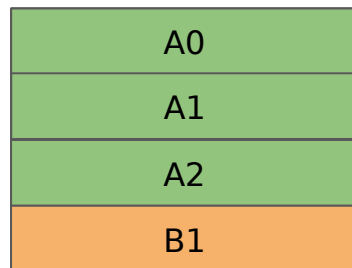
$\sum W_{ss_i} < \text{GPU memory}$

No page faults

Other times

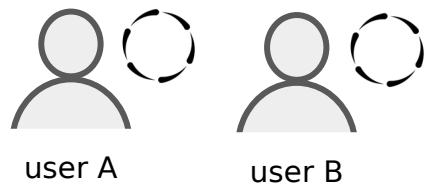


GPU

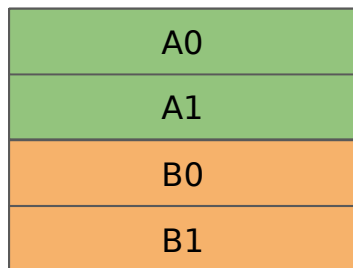


# When GPU bursts overlap

Sometimes



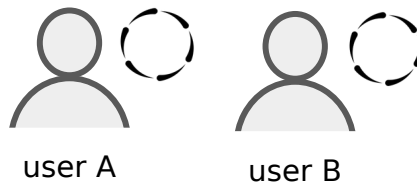
GPU



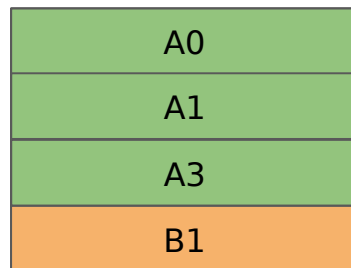
$\sum W_{ss_i} < \text{GPU memory}$

No page faults

Other times

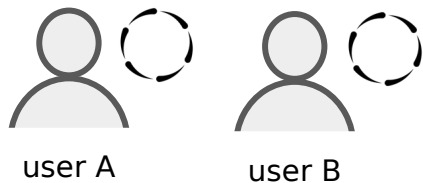


GPU

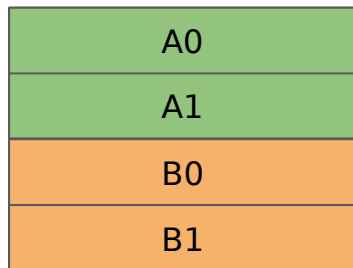


# When GPU bursts overlap

Sometimes



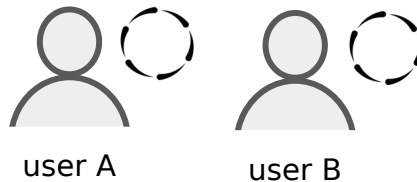
GPU



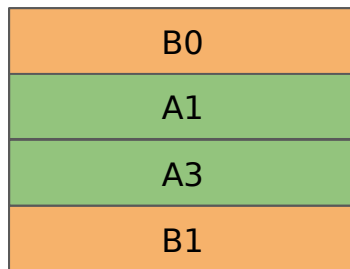
$$\sum W_{ss_i} < \text{GPU memory}$$

No page faults

Other times

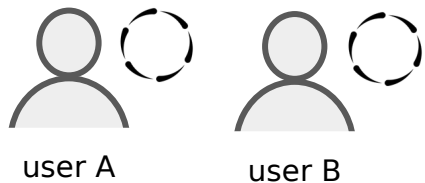


GPU

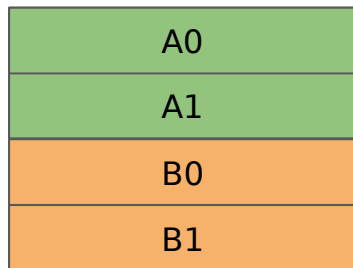


# When GPU bursts overlap

Sometimes



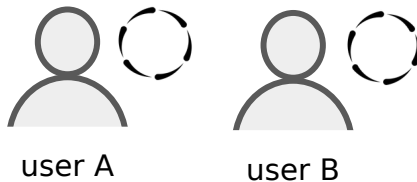
GPU



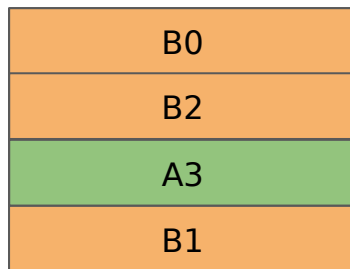
$\sum W_{ss_i} < \text{GPU memory}$

No page faults

Other times

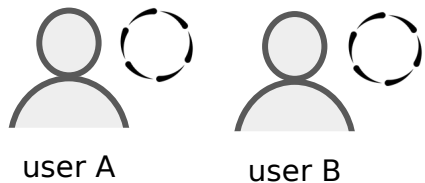


GPU

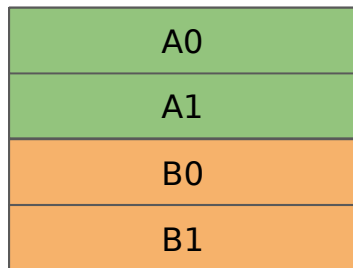


# When GPU bursts overlap

Sometimes



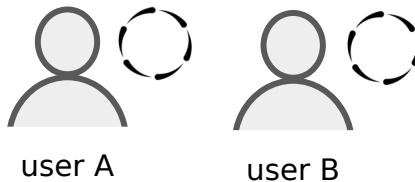
GPU



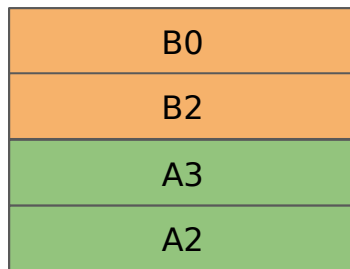
$\sum W_{ss_i} < \text{GPU memory}$

No page faults

Other times

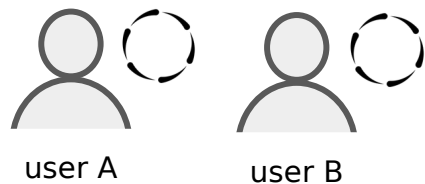


GPU

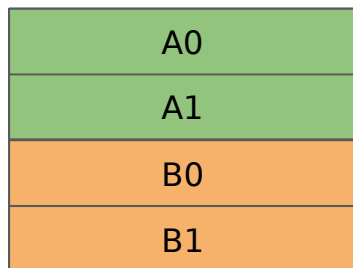


# When GPU bursts overlap

Sometimes



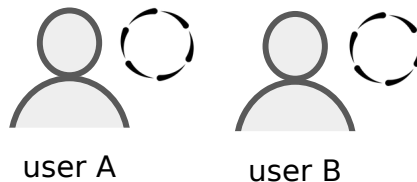
GPU



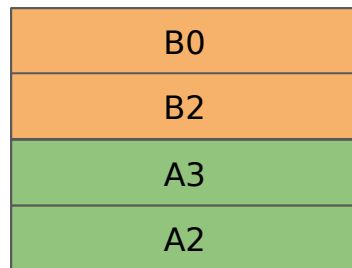
$\sum W_{ss_i} < \text{GPU memory}$

No page faults

Other times



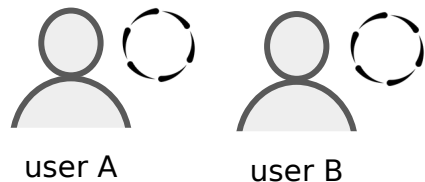
GPU



many  
page  
faults  
!

# When GPU bursts overlap

Sometimes



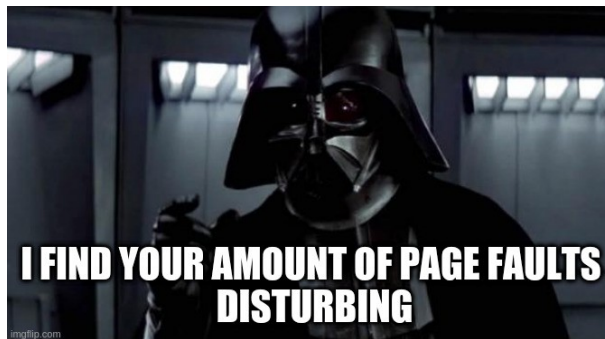
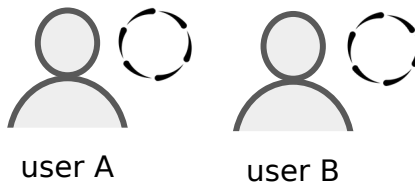
GPU

A0
A1
B0
B1

$$\sum W_{ss_i} < \text{GPU memory}$$

No page faults

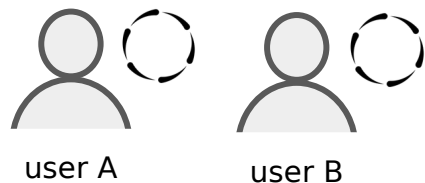
Other times



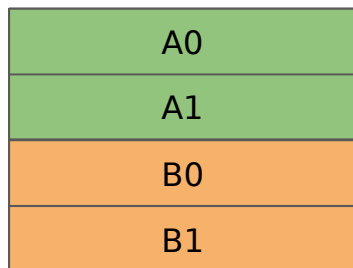


# When GPU bursts overlap

Sometimes



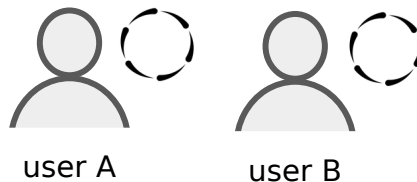
GPU



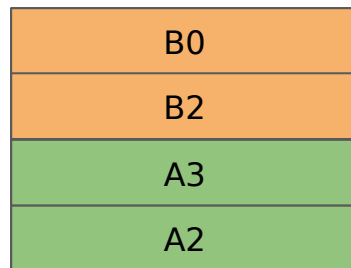
$\sum W_{ss_i} < \text{GPU memory}$

No page faults

Other times



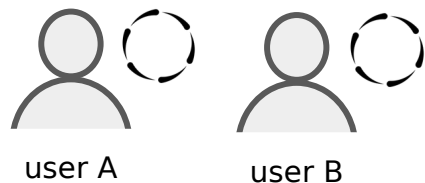
GPU



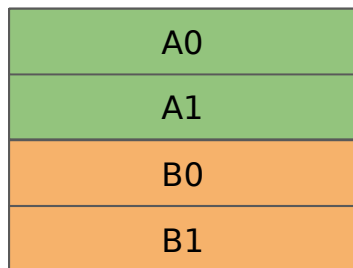
many  
page  
faults  
!

# When GPU bursts overlap

Sometimes



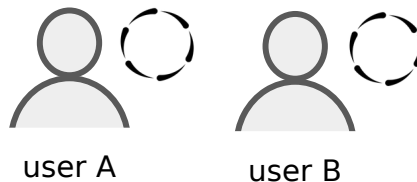
GPU



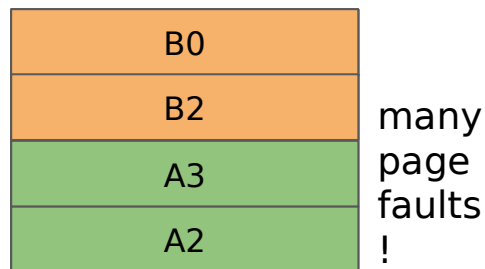
$$\sum W_{ss_i} < \text{GPU memory}$$

No page faults

Other times



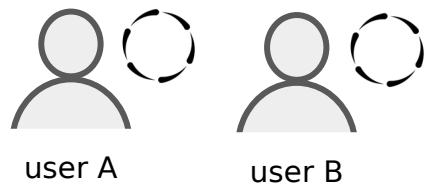
GPU



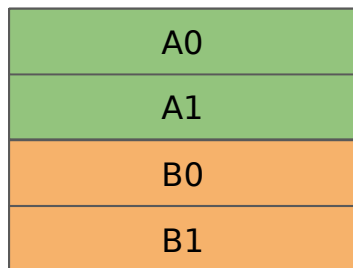
$$\sum W_{ss_i} > \text{GPU memory}$$

# When GPU bursts overlap

Sometimes



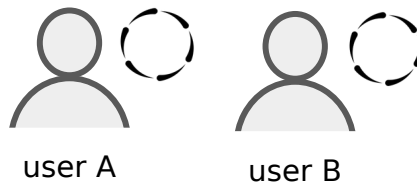
GPU



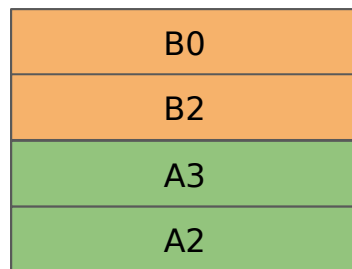
$$\sum W_{ss_i} < \text{GPU memory}$$

No page faults

Other times



GPU



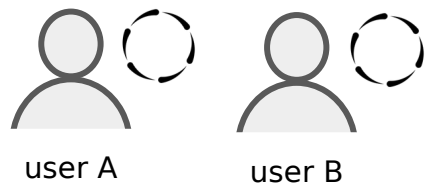
many  
page  
faults  
!

$$\sum W_{ss_i} > \text{GPU memory}$$

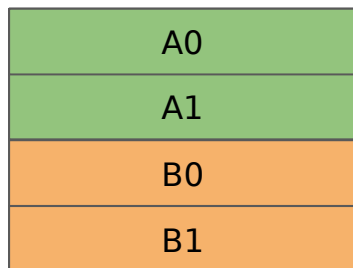
Thrashing!

# When GPU bursts overlap

Sometimes



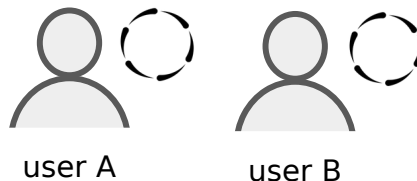
GPU



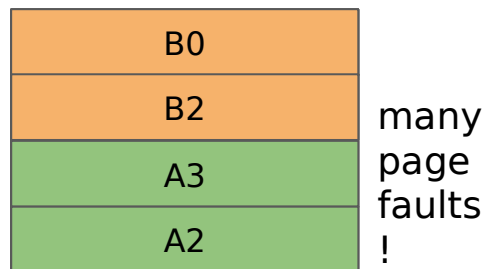
$$\sum W_{ss_i} < \text{GPU memory}$$

No page faults

Other times



GPU

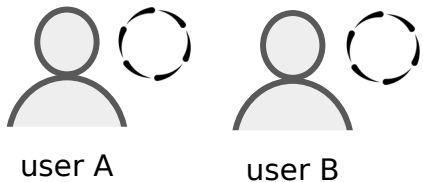


$$\sum W_{ss_i} > \text{GPU memory}$$

Thrashing!

# When GPU bursts overlap

Sometimes

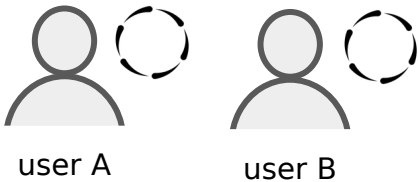


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



GPU

B0
B2
A3
A2

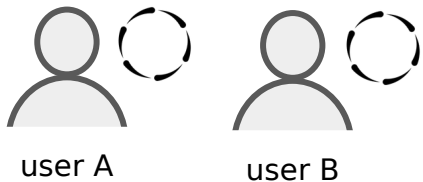
many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution

# When GPU bursts overlap

Sometimes

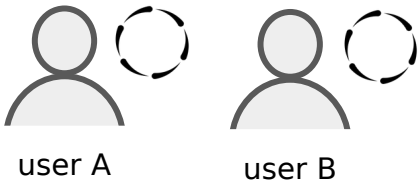


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



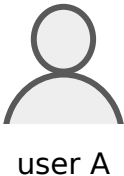
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



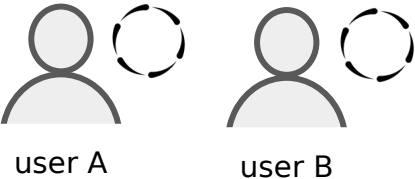
# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

Other times

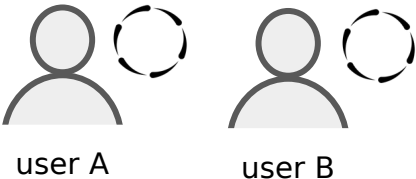
Solution



GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

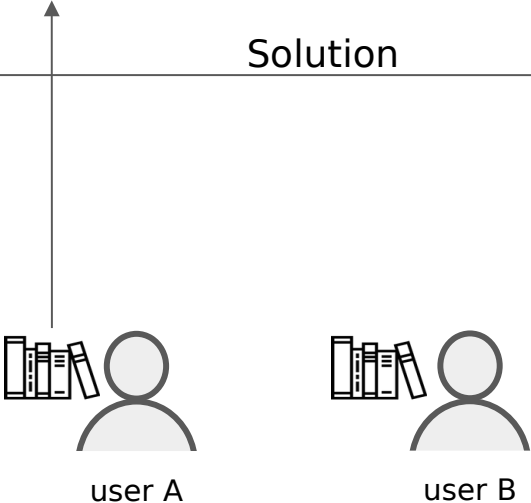


GPU

B0
B2
A3
A2

many  
page  
faults  
!

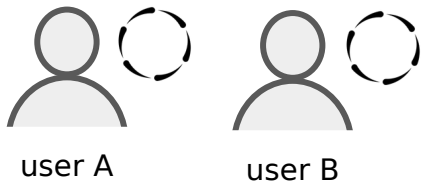
$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!



# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

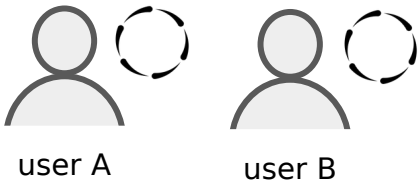


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



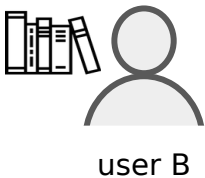
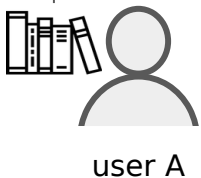
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution





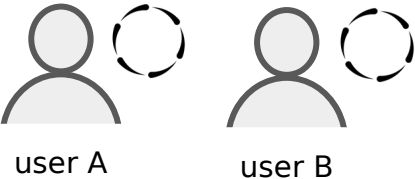
# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

Other times

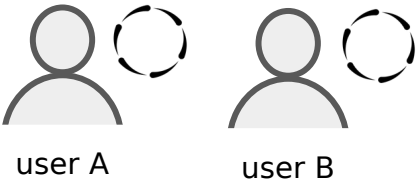
Solution



GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

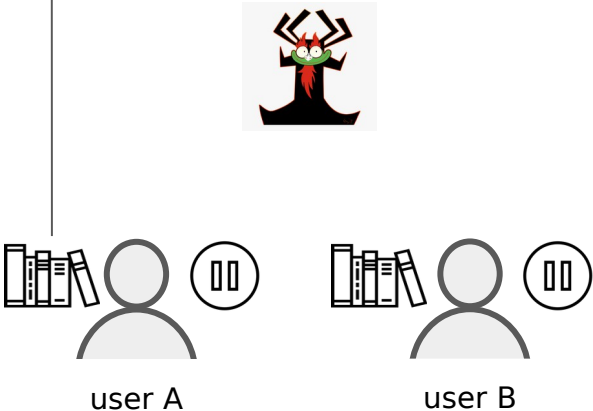


GPU

B0
B2
A3
A2

many  
page  
faults  
!

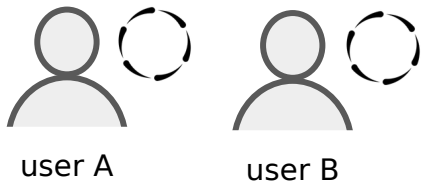
$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!



# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

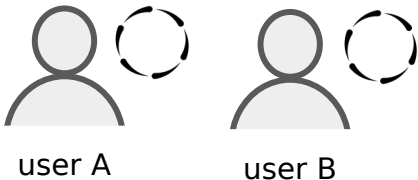


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



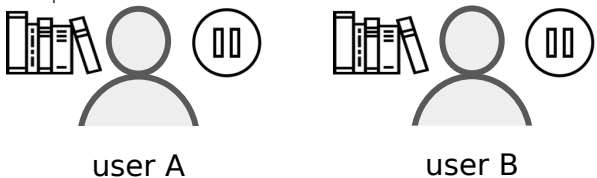
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

B0
B2
A3
A2

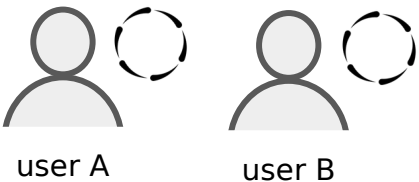
# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

Other times

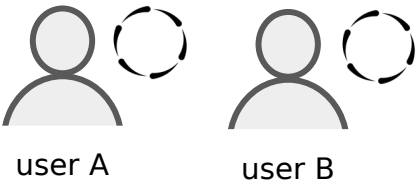
Solution



GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

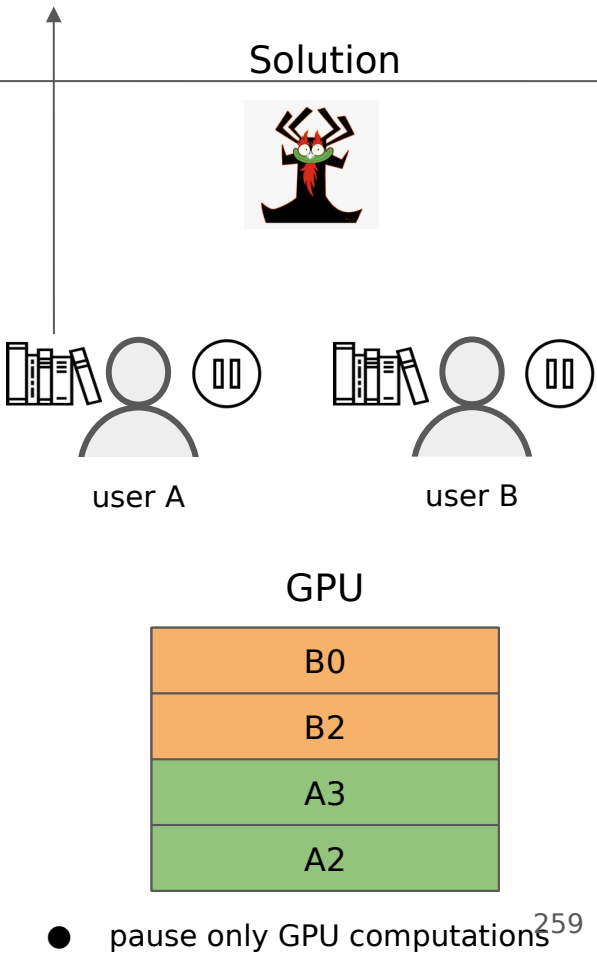


GPU

B0
B2
A3
A2

many  
page  
faults  
!

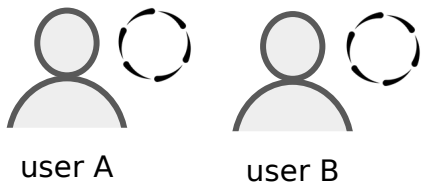
$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!



# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

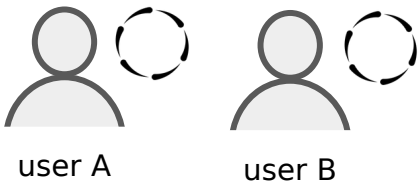


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



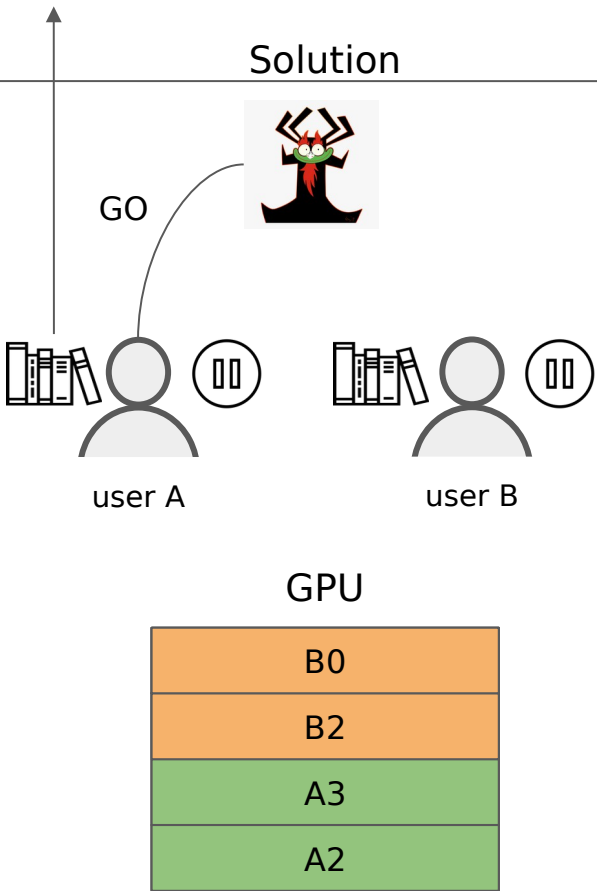
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution

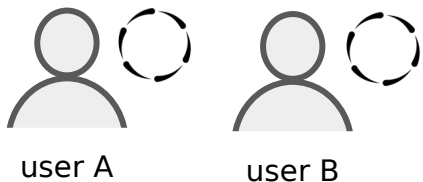


● pause only GPU computations<sup>260</sup>

# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes



GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



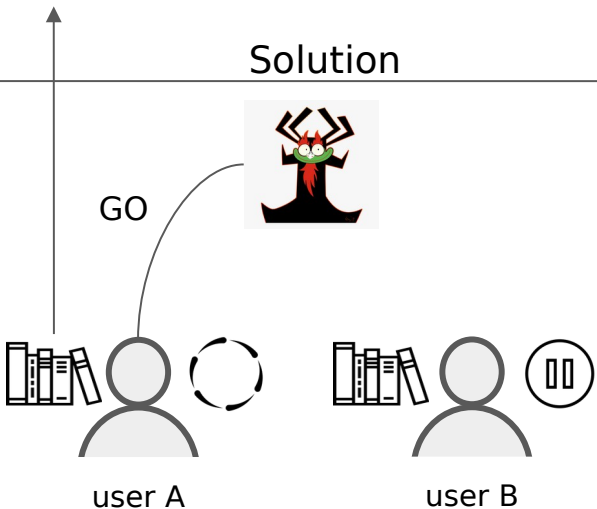
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

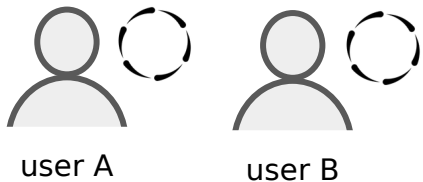
B0
B2
A3
A2

● pause only GPU computations<sup>261</sup>

# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes



GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



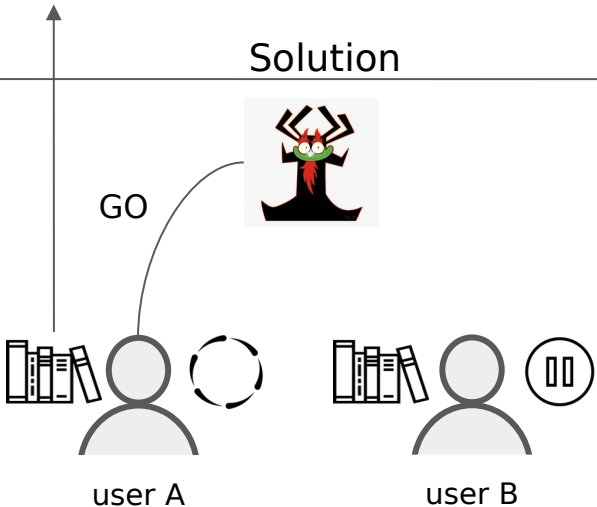
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

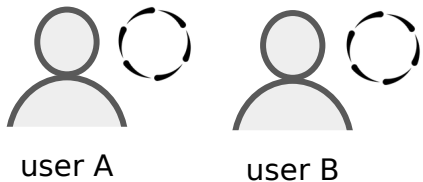
A0
B2
A3
A2

● pause only GPU computations<sup>262</sup>

# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

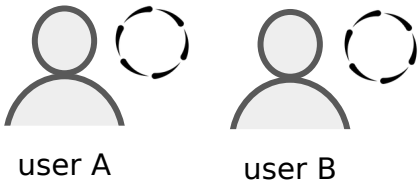


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



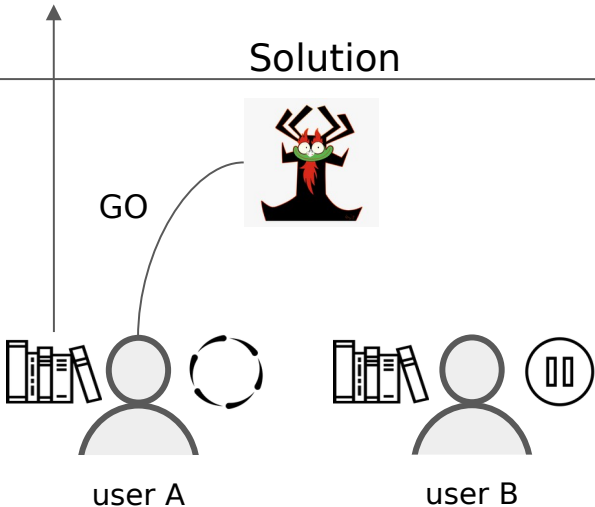
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

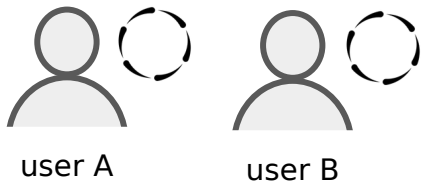
A0
A1
A3
A2

● pause only GPU computations<sup>263</sup>

# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

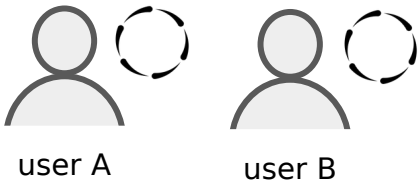


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



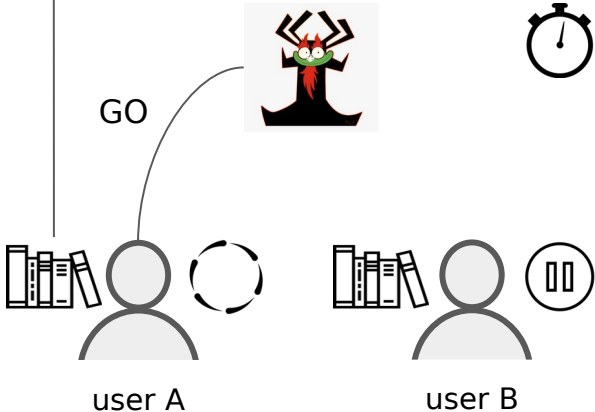
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

A0
A1
A3
A2

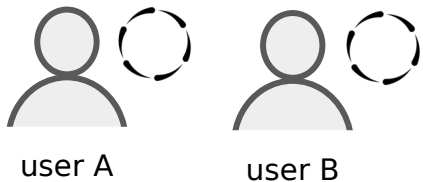
● pause only GPU computations<sup>264</sup>



# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

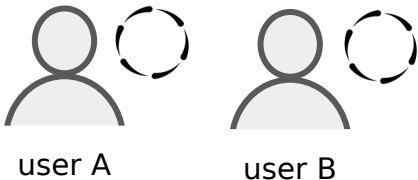


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



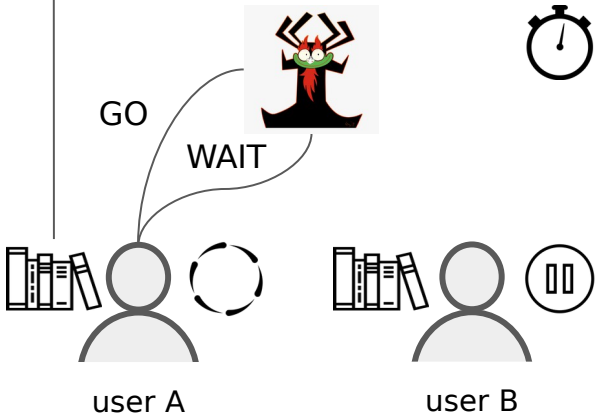
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

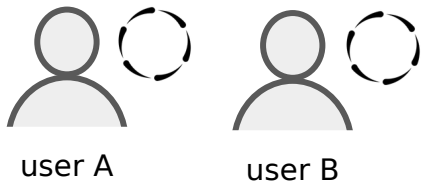
A0
A1
A3
A2

● pause only GPU computations<sup>265</sup>

# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

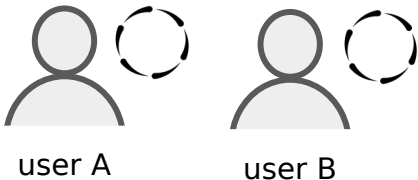


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



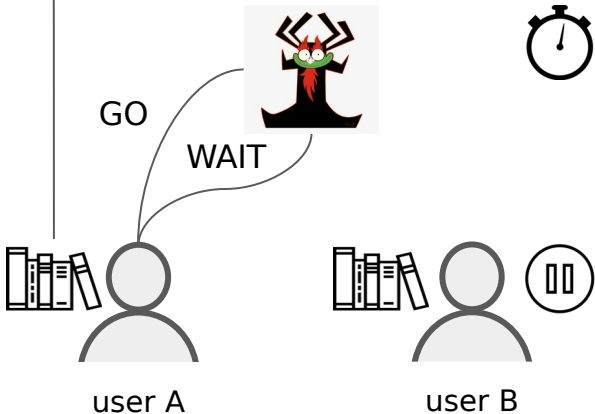
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

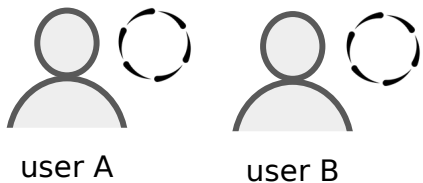
A0
A1
A3
A2

● pause only GPU computations<sup>266</sup>

# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes



GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



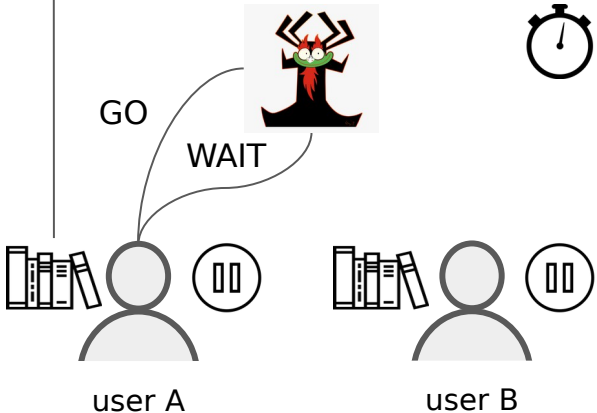
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

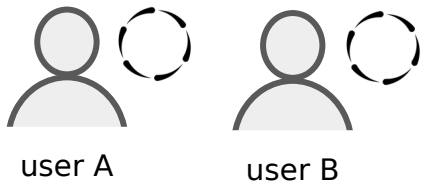
A0
A1
A3
A2

● pause only GPU computations<sup>267</sup>

# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

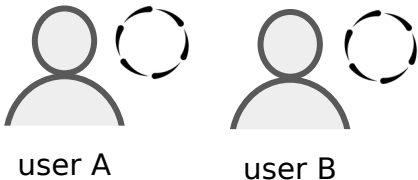


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



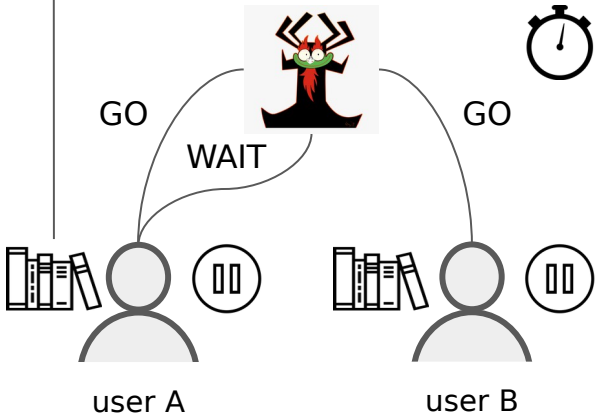
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

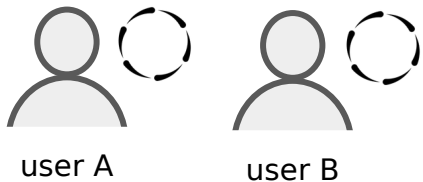
A0
A1
A3
A2

● pause only GPU computations<sup>268</sup>

# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

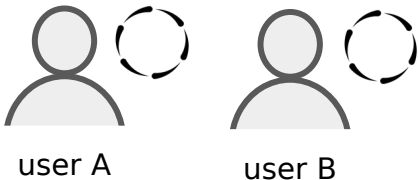


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



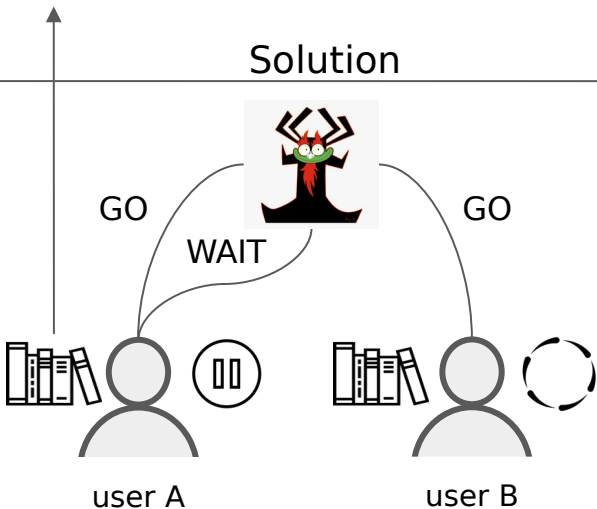
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

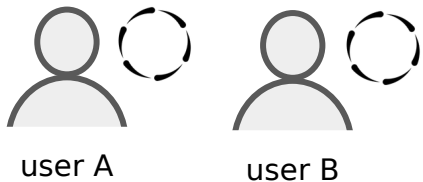
A0
A1
A3
A2

● pause only GPU computations<sup>269</sup>

# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

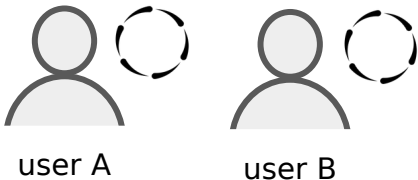


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



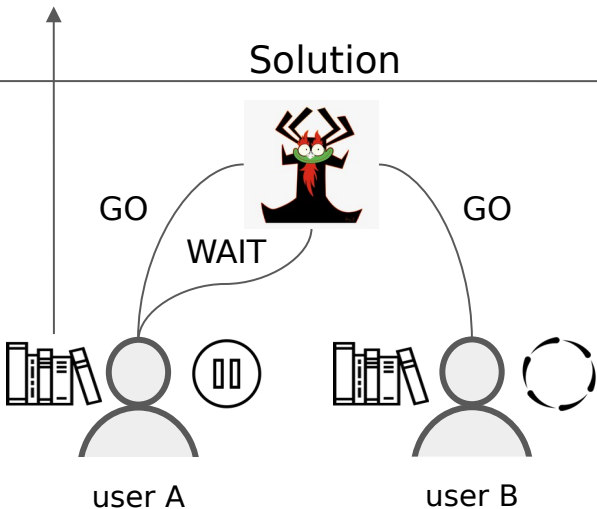
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

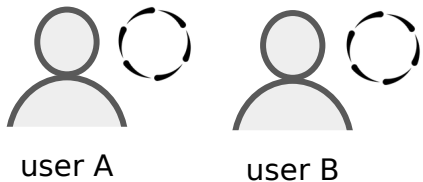
A0
A1
B0
A2

● pause only GPU computations<sup>270</sup>

# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

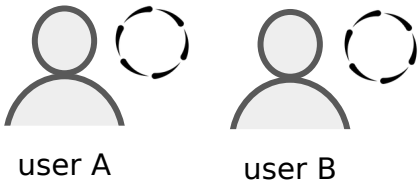


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



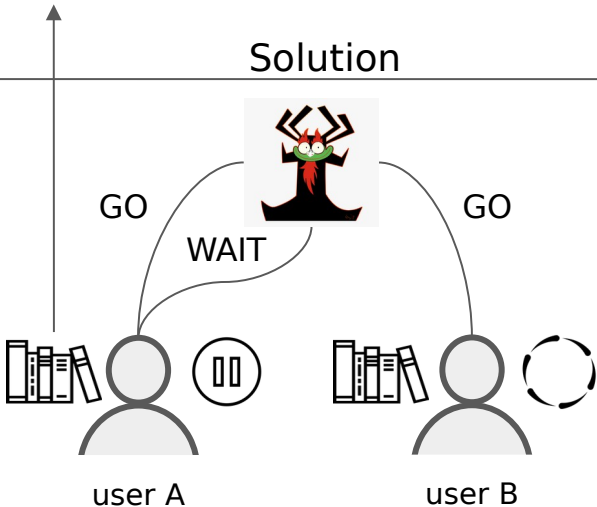
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

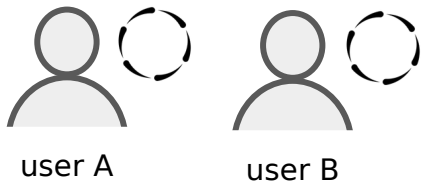
A0
A1
B0
B1

● pause only GPU computations<sup>271</sup>

# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes



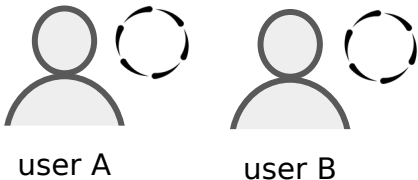
GPU

A0
A1
B0
B1

$$\sum W_{ss_i} < \text{GPU memory}$$

No page faults

Other times



GPU

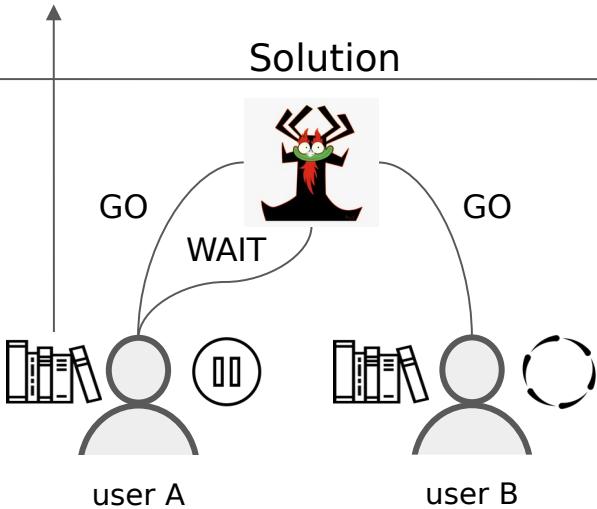
B0
B2
A3
A2

many  
page  
faults  
!

$$\sum W_{ss_i} > \text{GPU memory}$$

Thrashing!

Solution



GPU

B2
A1
B0
B1

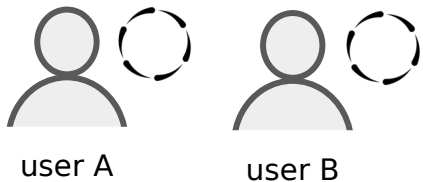
● pause only GPU computations<sup>272</sup>



# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

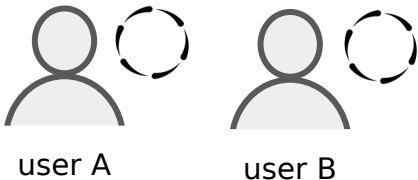


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



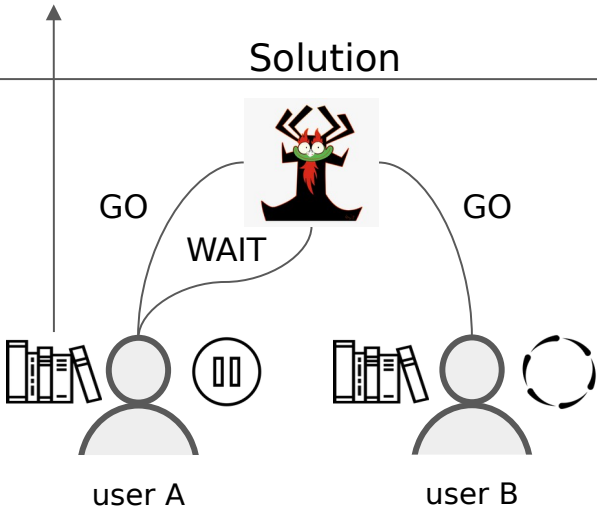
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

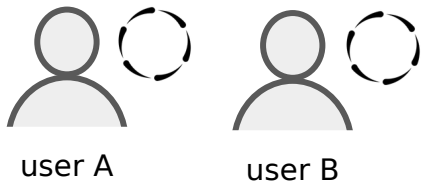
B2
B3
B0
B1

● pause only GPU computations<sup>273</sup>

# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

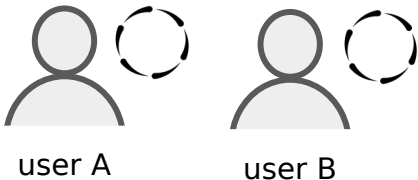


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



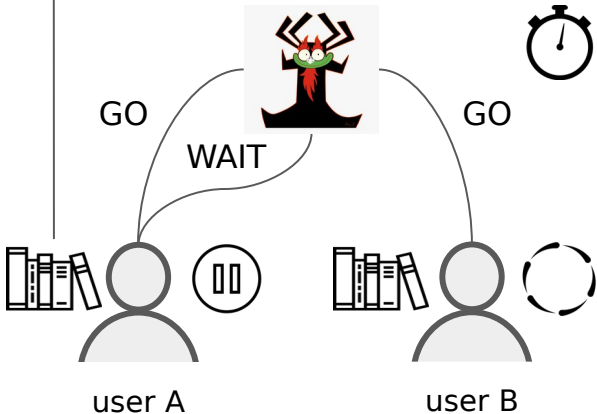
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

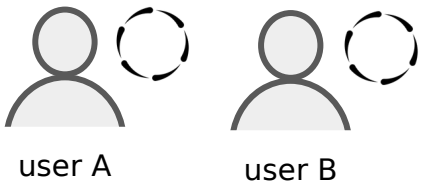
B2
B3
B0
B1

● pause only GPU computations<sup>274</sup>

# When GPU bursts overlap

LD\_PRELOAD=libnvshare.so

Sometimes

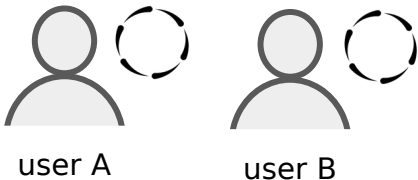


GPU

A0
A1
B0
B1

$\sum W_{ss_i} < \text{GPU memory}$   
No page faults

Other times



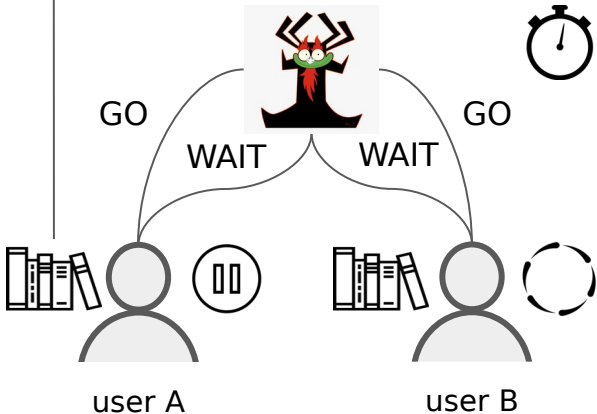
GPU

B0
B2
A3
A2

many  
page  
faults  
!

$\sum W_{ss_i} > \text{GPU memory}$   
Thrashing!

Solution



GPU

B2
B3
B0
B1

● pause only GPU computations<sup>275</sup>

# Anti-thrashing idea

# Anti-thrashing idea

- $WSS > GPU \text{ Memory}$

# Anti-thrashing idea

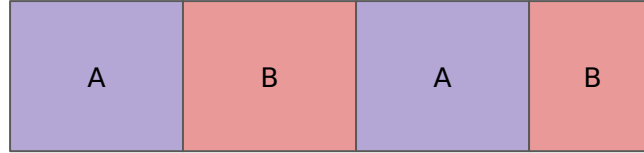
- WSS > GPU Memory

Default case:

# Anti-thrashing idea

- WSS > GPU Memory

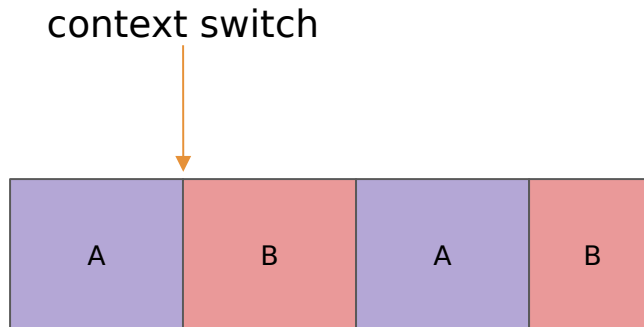
Default case:



# Anti-thrashing idea

- WSS > GPU Memory

Default case:

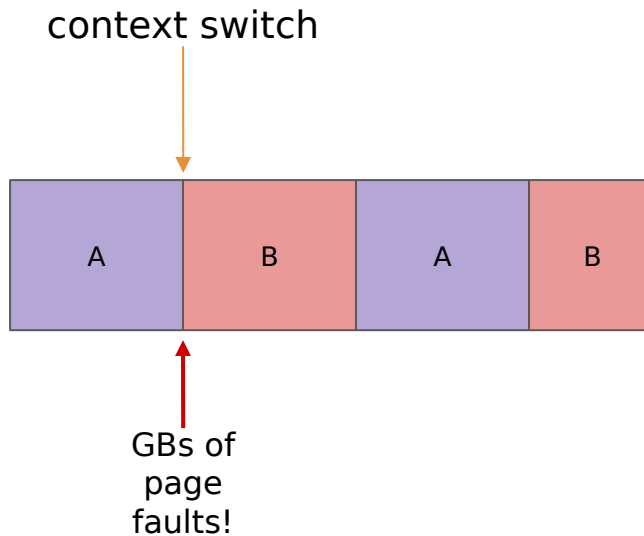




# Anti-thrashing idea

- WSS > GPU Memory

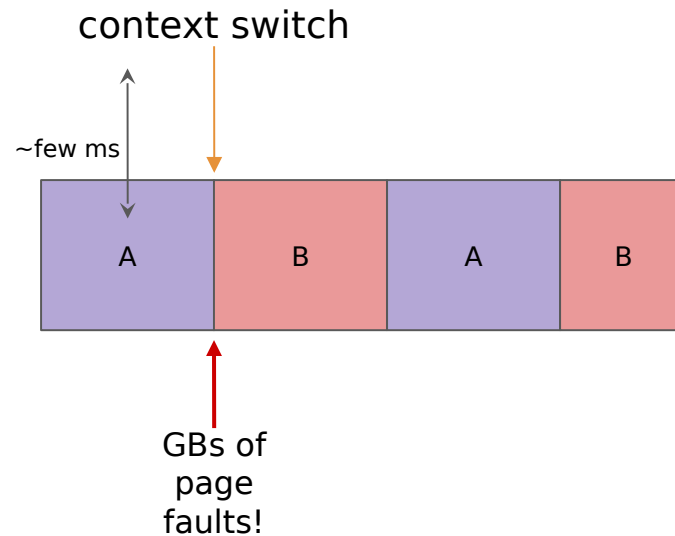
Default case:



# Anti-thrashing idea

- WSS > GPU Memory

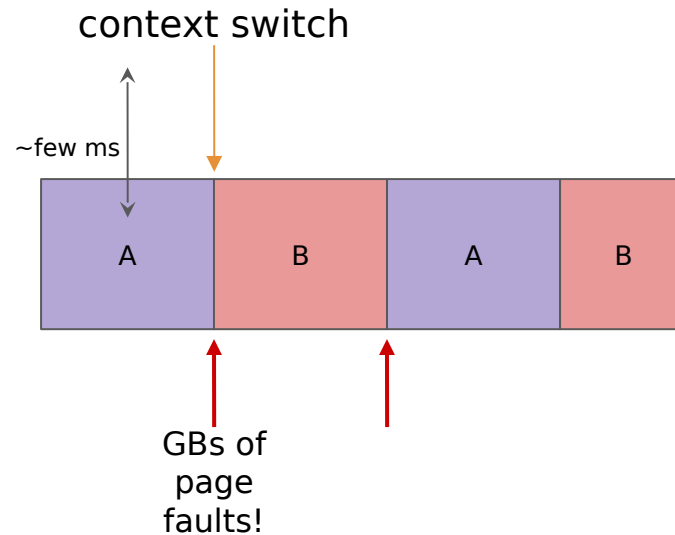
Default case:



# Anti-thrashing idea

- WSS > GPU Memory

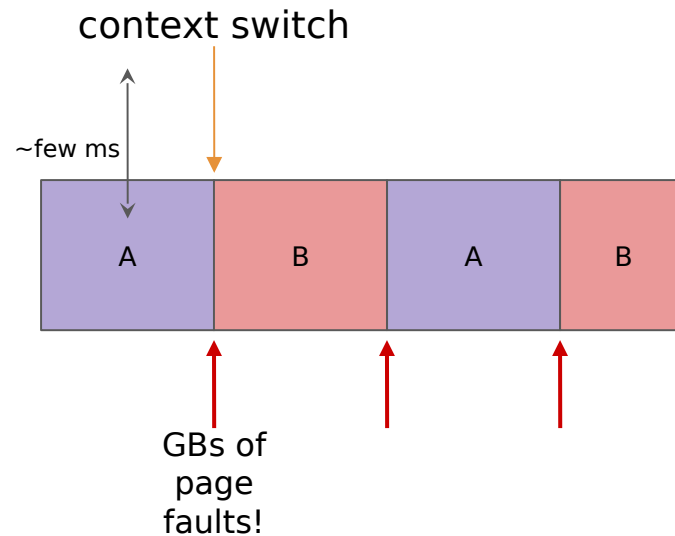
Default case:



# Anti-thrashing idea

- WSS > GPU Memory

Default case:

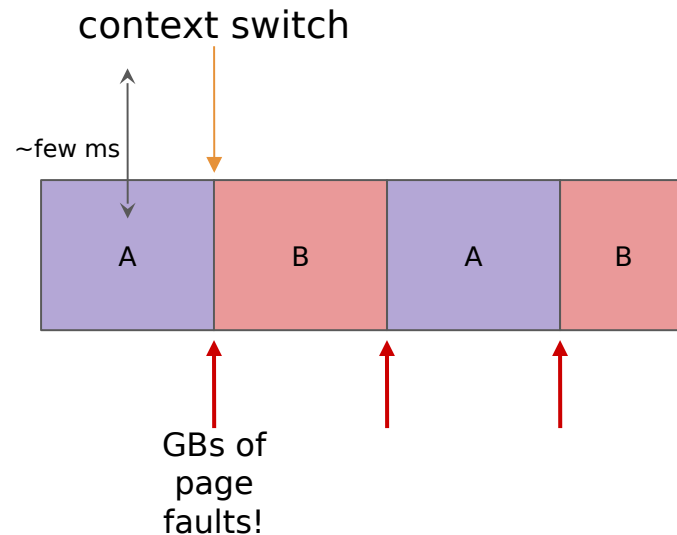


# Anti-thrashing idea

- WSS > GPU Memory

## Default case:

- memory swaps happen too frequently

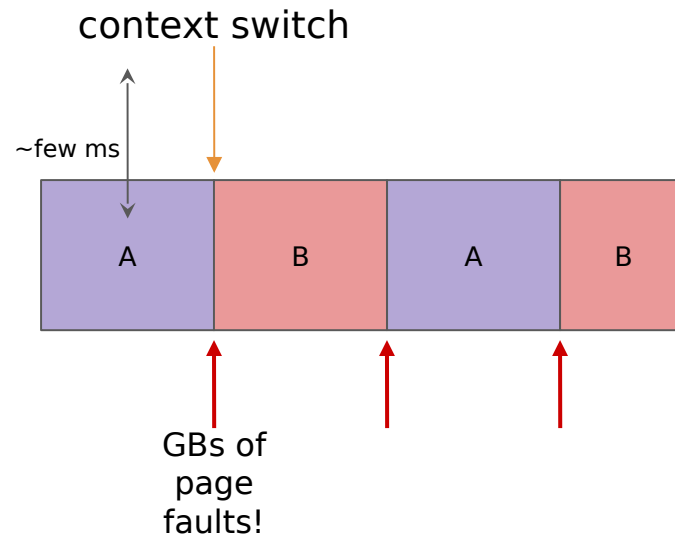


# Anti-thrashing idea

- WSS > GPU Memory

## Default case:

- memory swaps happen too frequently
- Thrashing!

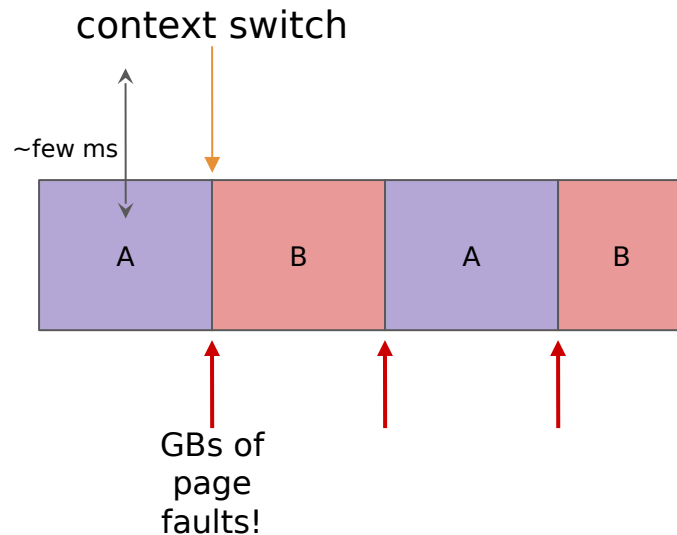


# Anti-thrashing idea

- WSS > GPU Memory

## Default case:

- memory swaps happen too frequently
- Thrashing!



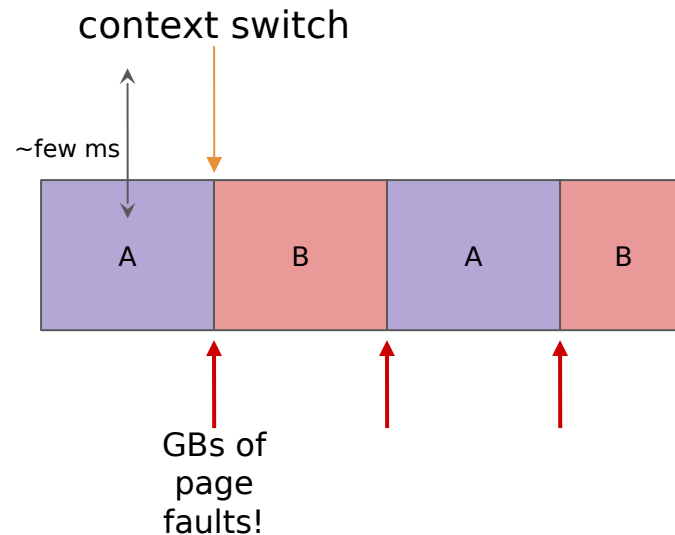
## anti-thrashing mechanism:

# Anti-thrashing idea

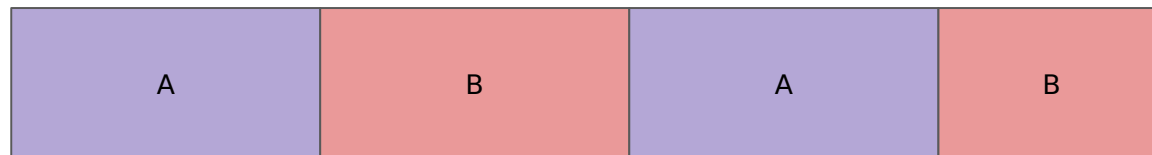
- WSS > GPU Memory

## Default case:

- memory swaps happen too frequently
- Thrashing!



## anti-thrashing mechanism:



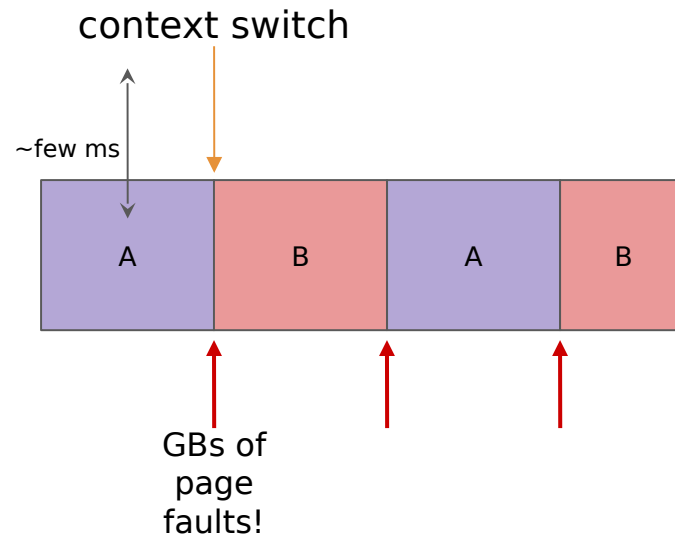


# Anti-thrashing idea

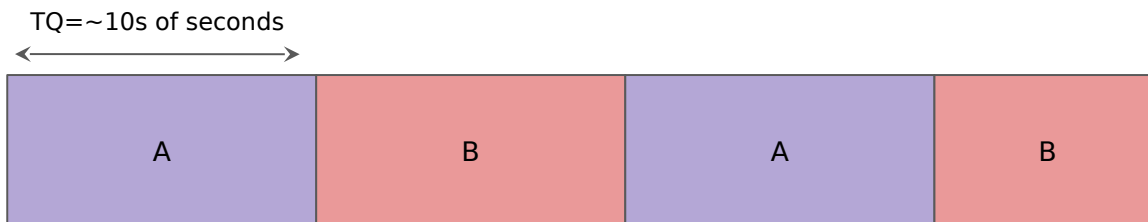
- WSS > GPU Memory

## Default case:

- memory swaps happen too frequently
- Thrashing!



## anti-thrashing mechanism:

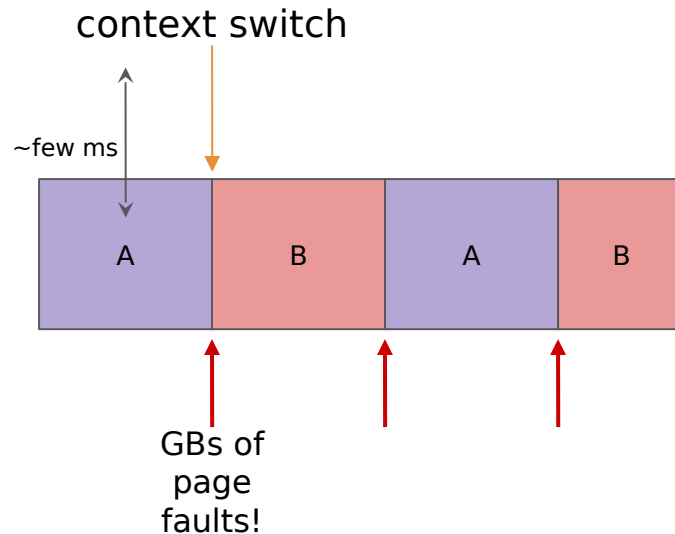


# Anti-thrashing idea

- WSS > GPU Memory

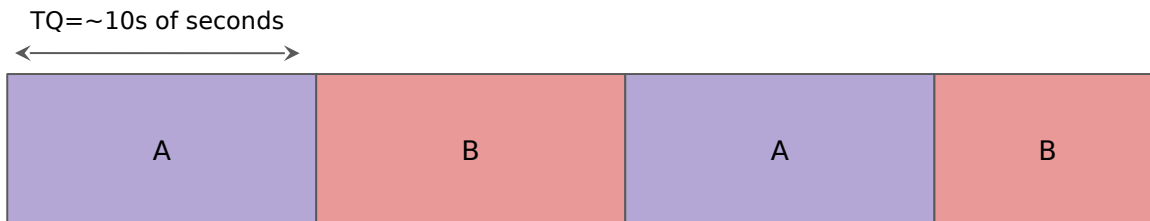
## Default case:

- memory swaps happen too frequently
- Thrashing!



## anti-thrashing mechanism:

- orders of magnitude fewer memory swaps

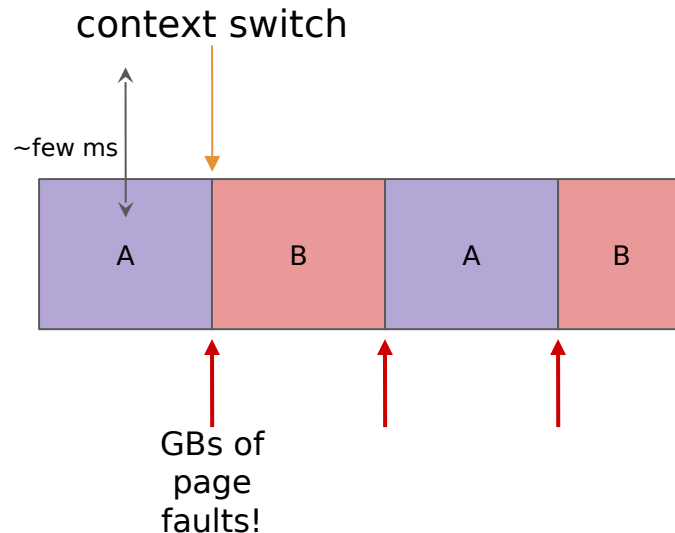


# Anti-thrashing idea

- WSS > GPU Memory

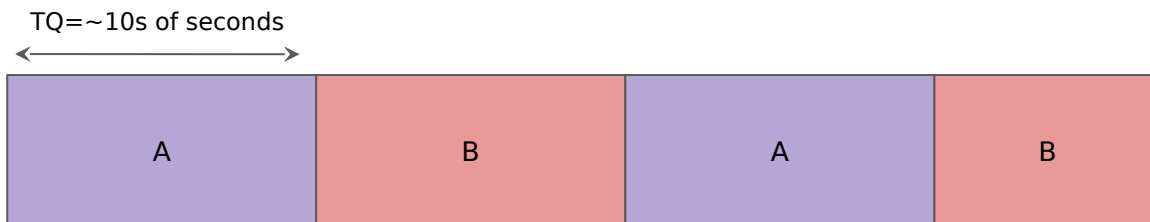
## Default case:

- memory swaps happen too frequently
- Thrashing!



## anti-thrashing mechanism:

- orders of magnitude fewer memory swaps
- Cooperative Scheduling

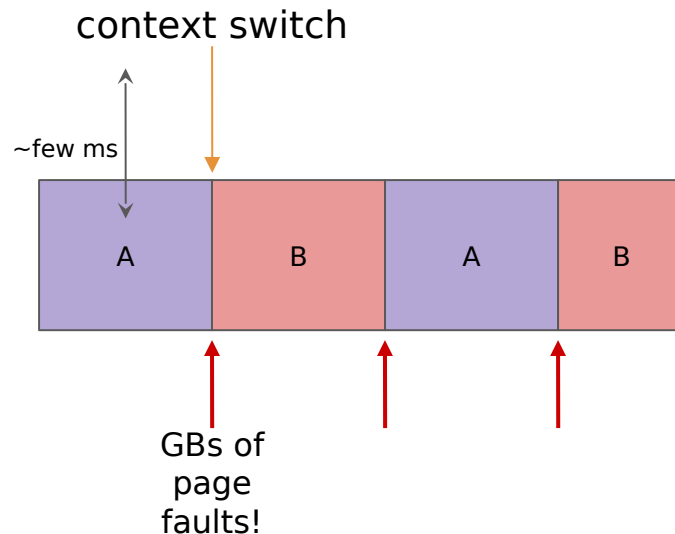


# Anti-thrashing idea

- WSS > GPU Memory

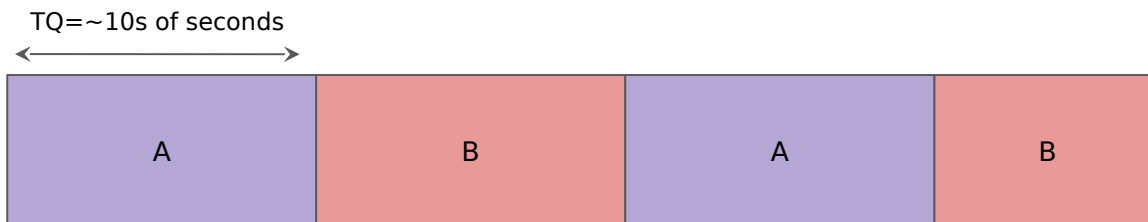
## Default case:

- memory swaps happen too frequently
- Thrashing!



## anti-thrashing mechanism:

- orders of magnitude fewer memory swaps
- Cooperative Scheduling
- No Thrashing!



Demo

# Evaluation

# Evaluation

- 16-core CPU (Intel Xeon 2.3 GHz - Haswell), 104 GB of RAM and an NVIDIA Tesla P100 GPU (16 GB).

# Evaluation

- 16-core CPU (Intel Xeon 2.3 GHz - Haswell), 104 GB of RAM and an NVIDIA Tesla P100 GPU (16 GB).
- non-interactive (conventional) tasks



# Evaluation

- 16-core CPU (Intel Xeon 2.3 GHz - Haswell), 104 GB of RAM and an NVIDIA Tesla P100 GPU (16 GB).
- non-interactive (conventional) tasks
- measure **Total Completion Time**

# Evaluation

- 16-core CPU (Intel Xeon 2.3 GHz - Haswell), 104 GB of RAM and an NVIDIA Tesla P100 GPU (16 GB).
- non-interactive (conventional) tasks
- measure **Total Completion Time**

## Programs

Name	Working Set Size	GPU/CPU ratio
small_90	7.2 GB	90/10
small_50	7.2 GB	50/50
big_90	15.3 GB	90/10
big_50	15.3 GB	50/50

# Evaluation

- 16-core CPU (Intel Xeon 2.3 GHz - Haswell), 104 GB of RAM and an NVIDIA Tesla P100 GPU (16 GB).
- non-interactive (conventional) tasks
- measure **Total Completion Time**

## Programs

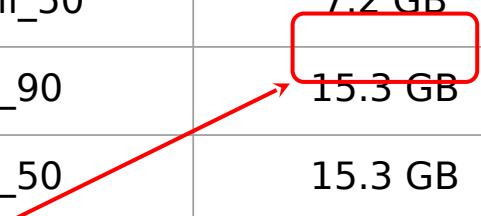
Name	Working Set Size	GPU/CPU ratio
small_90	7.2 GB	90/10
small_50	7.2 GB	50/50
big_90	15.3 GB	90/10
big_50	15.3 GB	50/50

# Evaluation

- 16-core CPU (Intel Xeon 2.3 GHz - Haswell), 104 GB of RAM and an NVIDIA Tesla P100 GPU (16 GB).
- non-interactive (conventional) tasks
- measure **Total Completion Time**

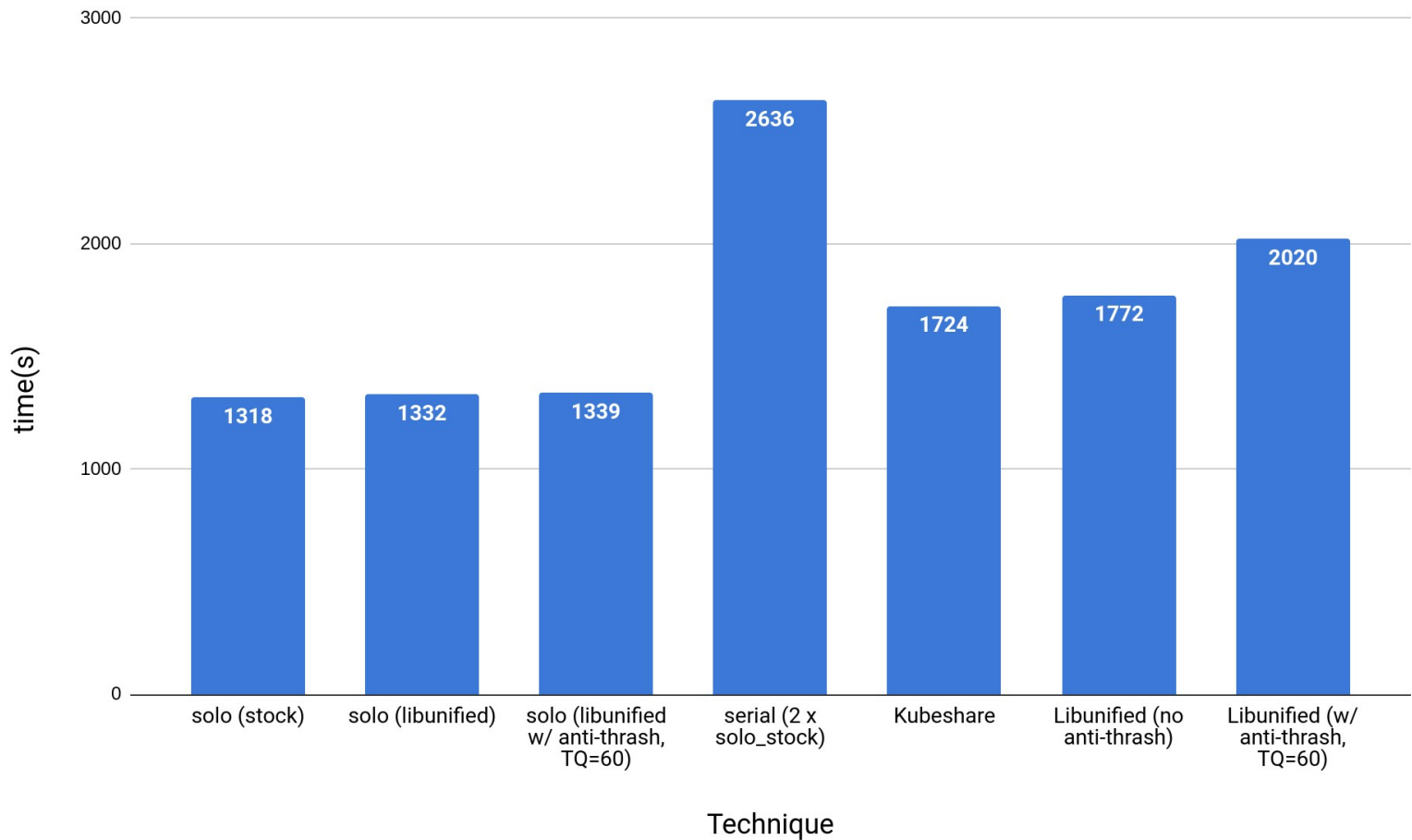
## Programs

Name	Working Set Size	GPU/CPU ratio
small_90	7.2 GB	90/10
small_50	7.2 GB	50/50
big_90	15.3 GB	90/10
big_50	15.3 GB	50/50

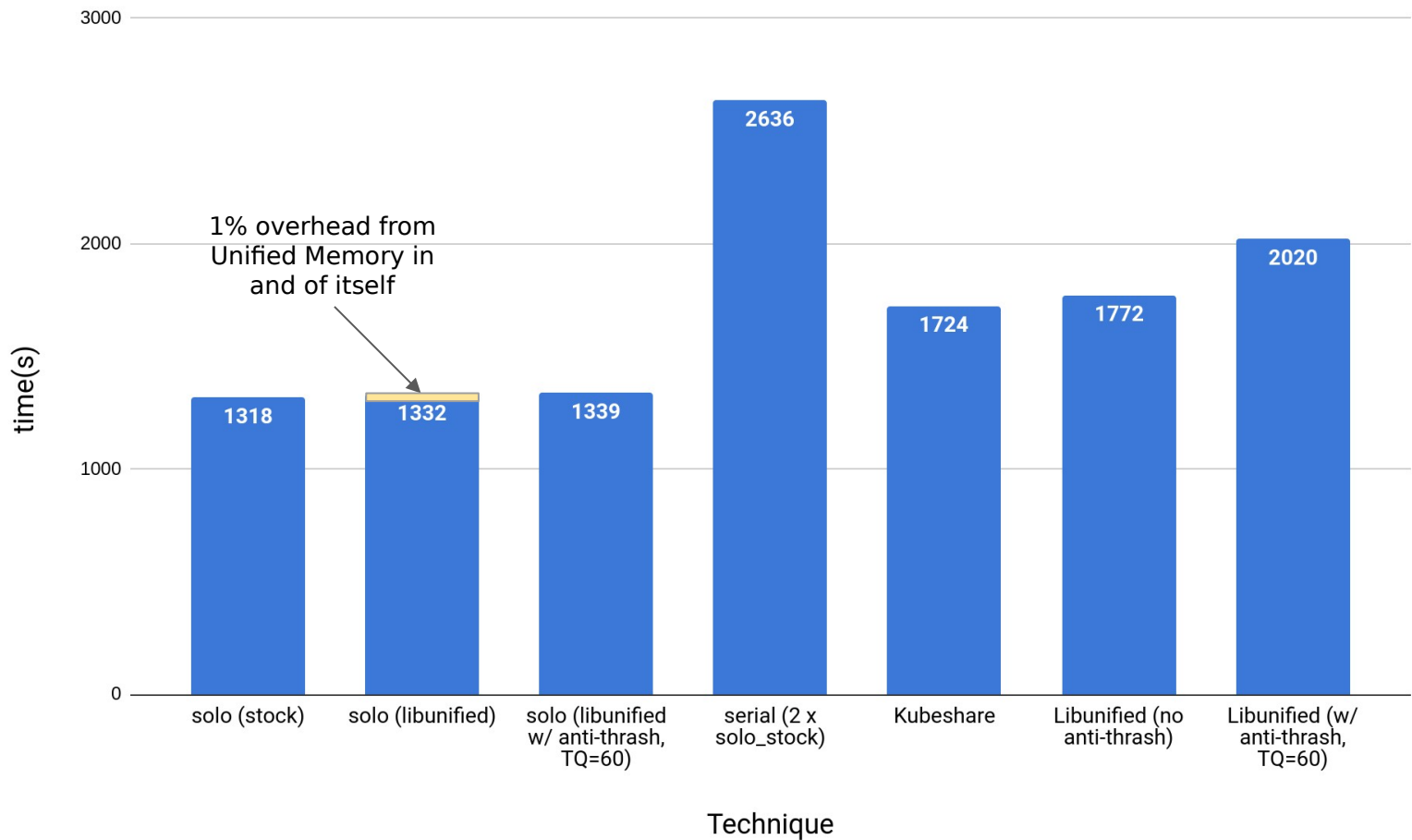


- Running 2 x big\_{50,90} in parallel causes **thrashing!**

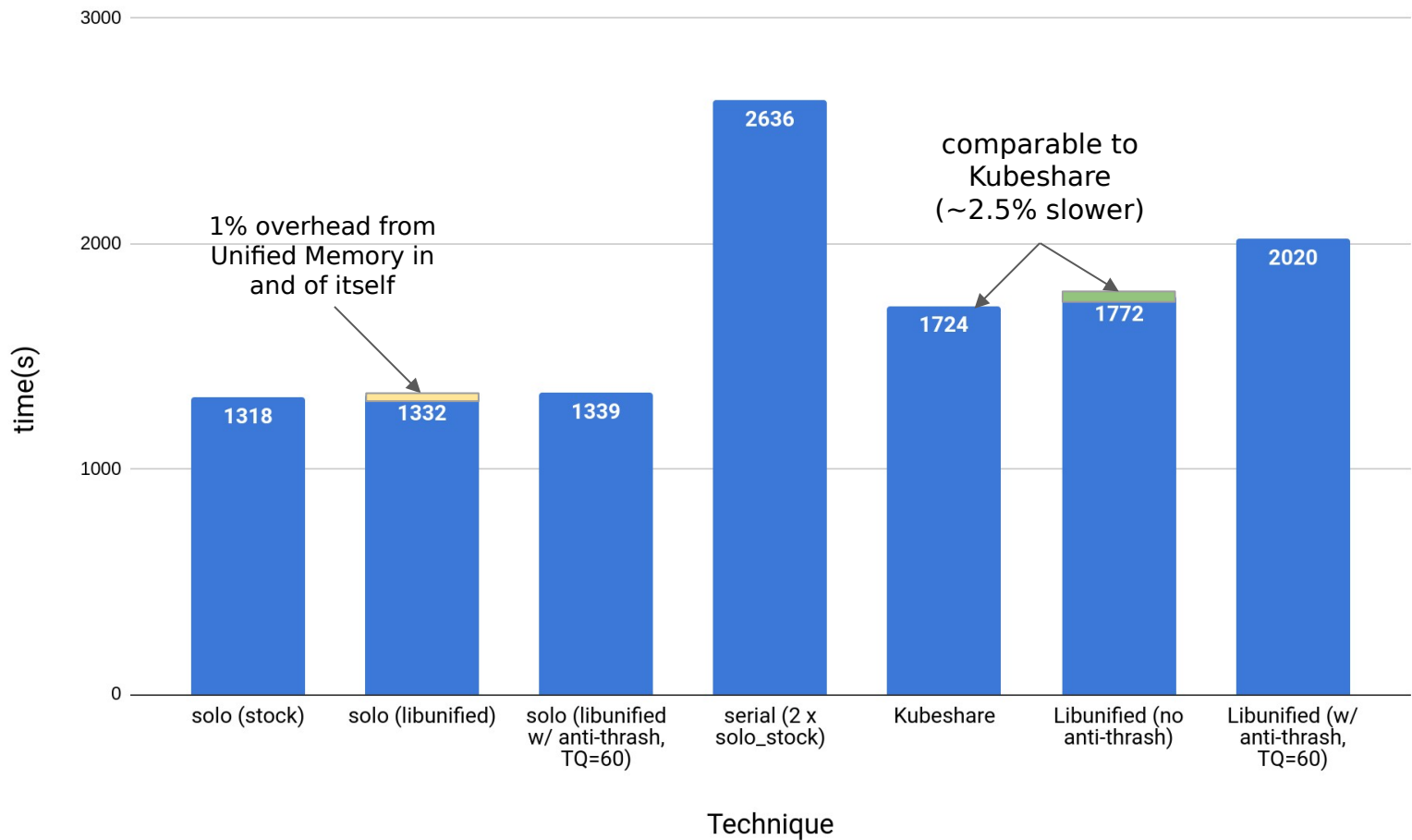
# Results for small\_50



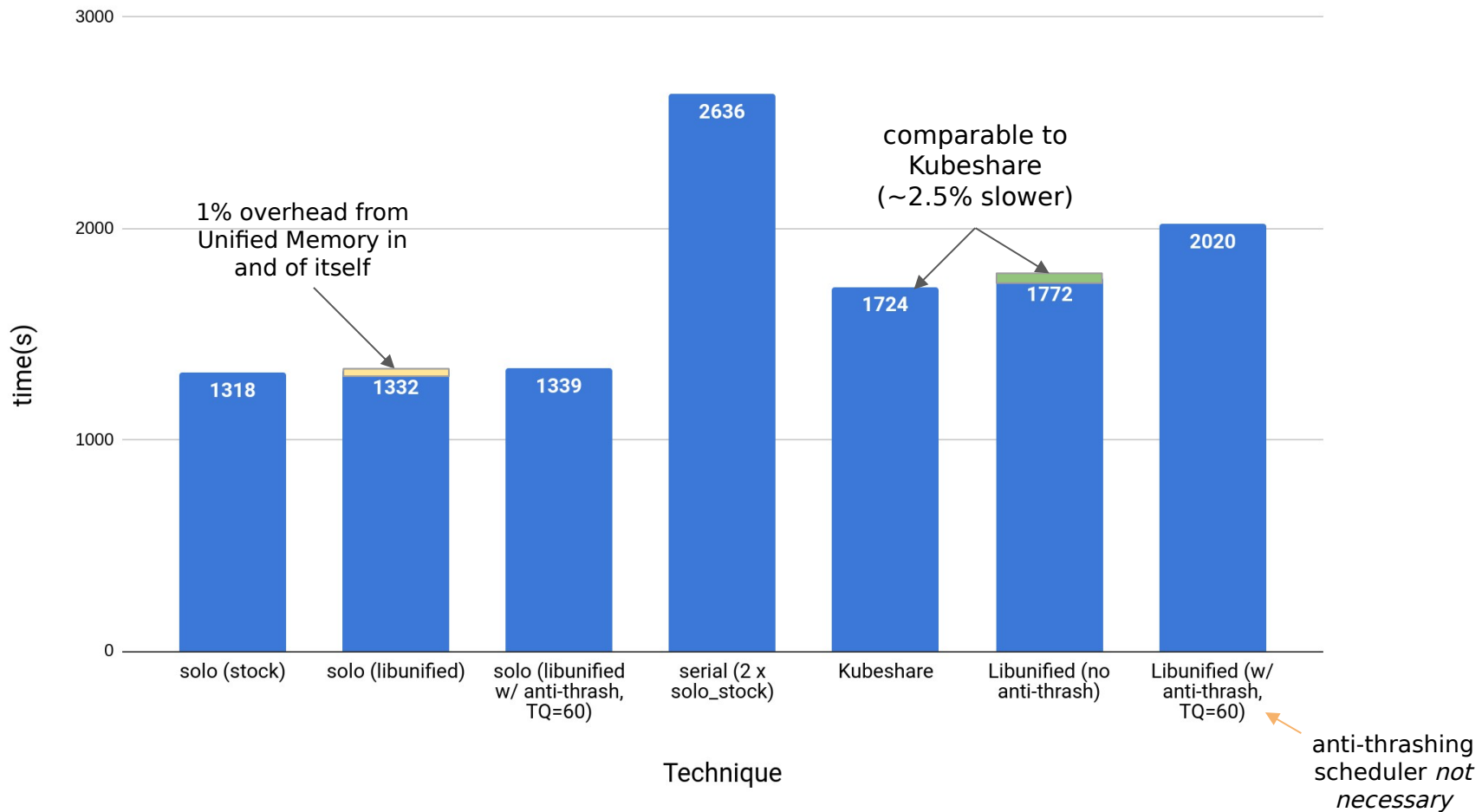
# Results for small\_50



# Results for small\_50

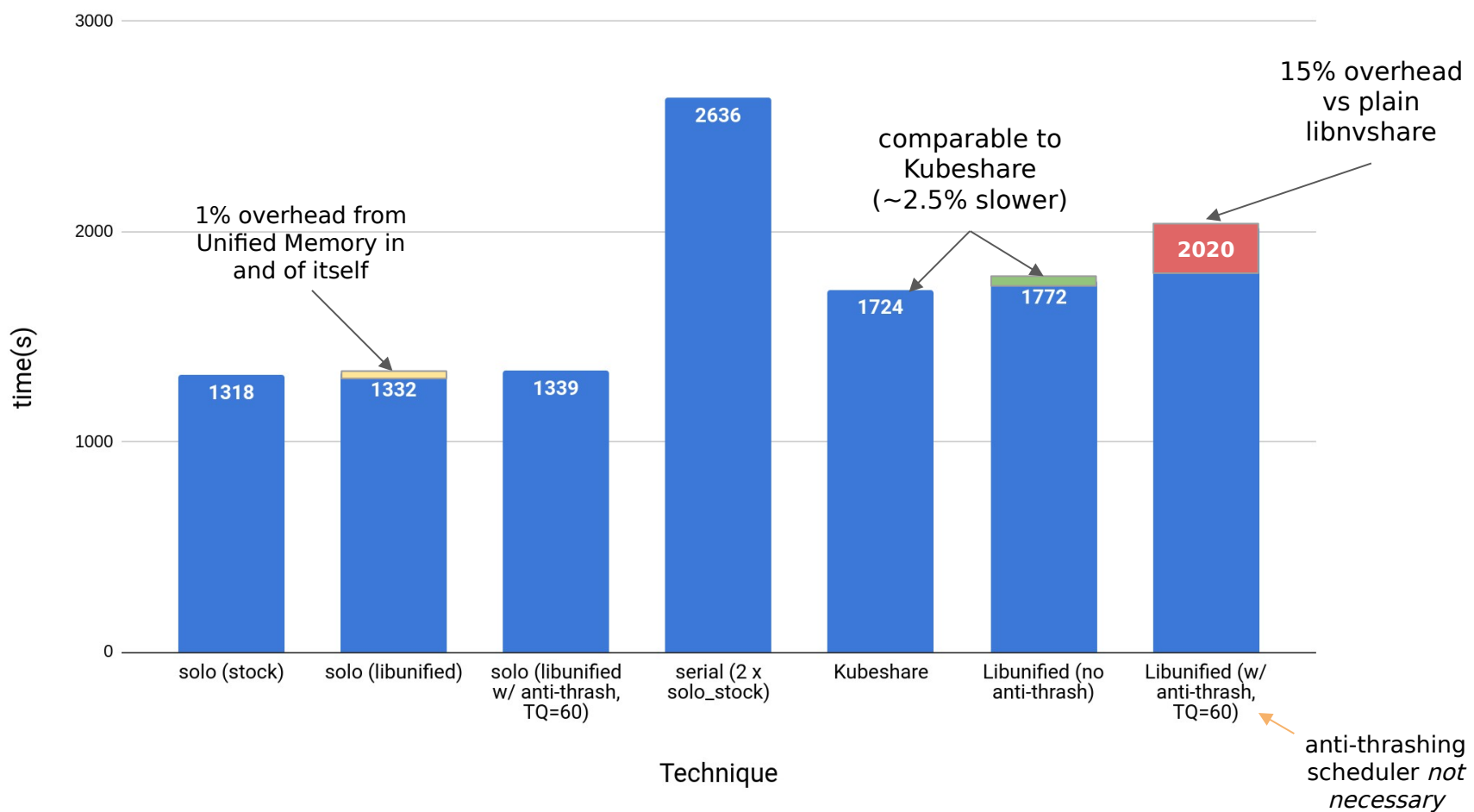


# Results for small\_50

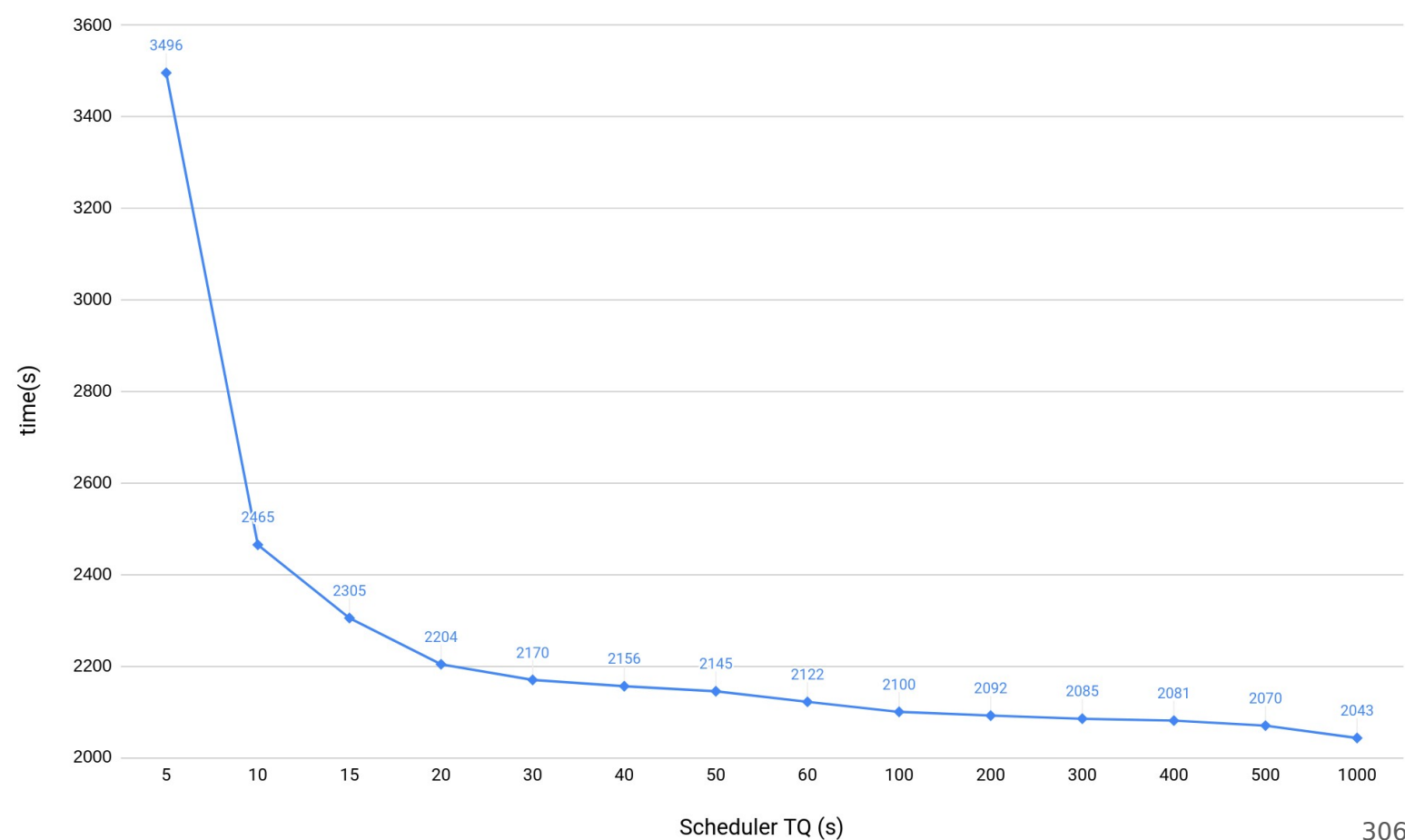




# Results for small\_50

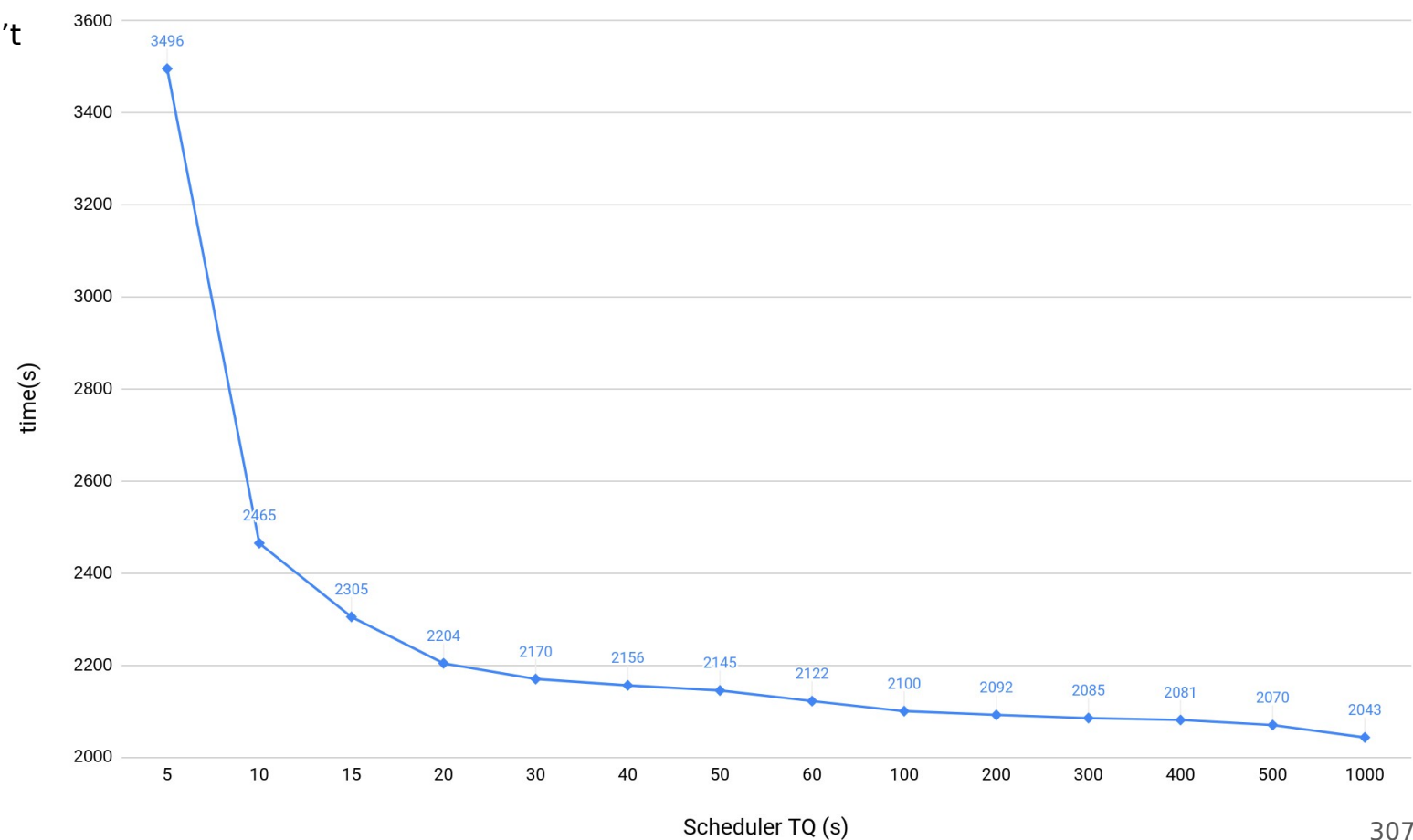


# Results for big\_50



# Results for big\_50

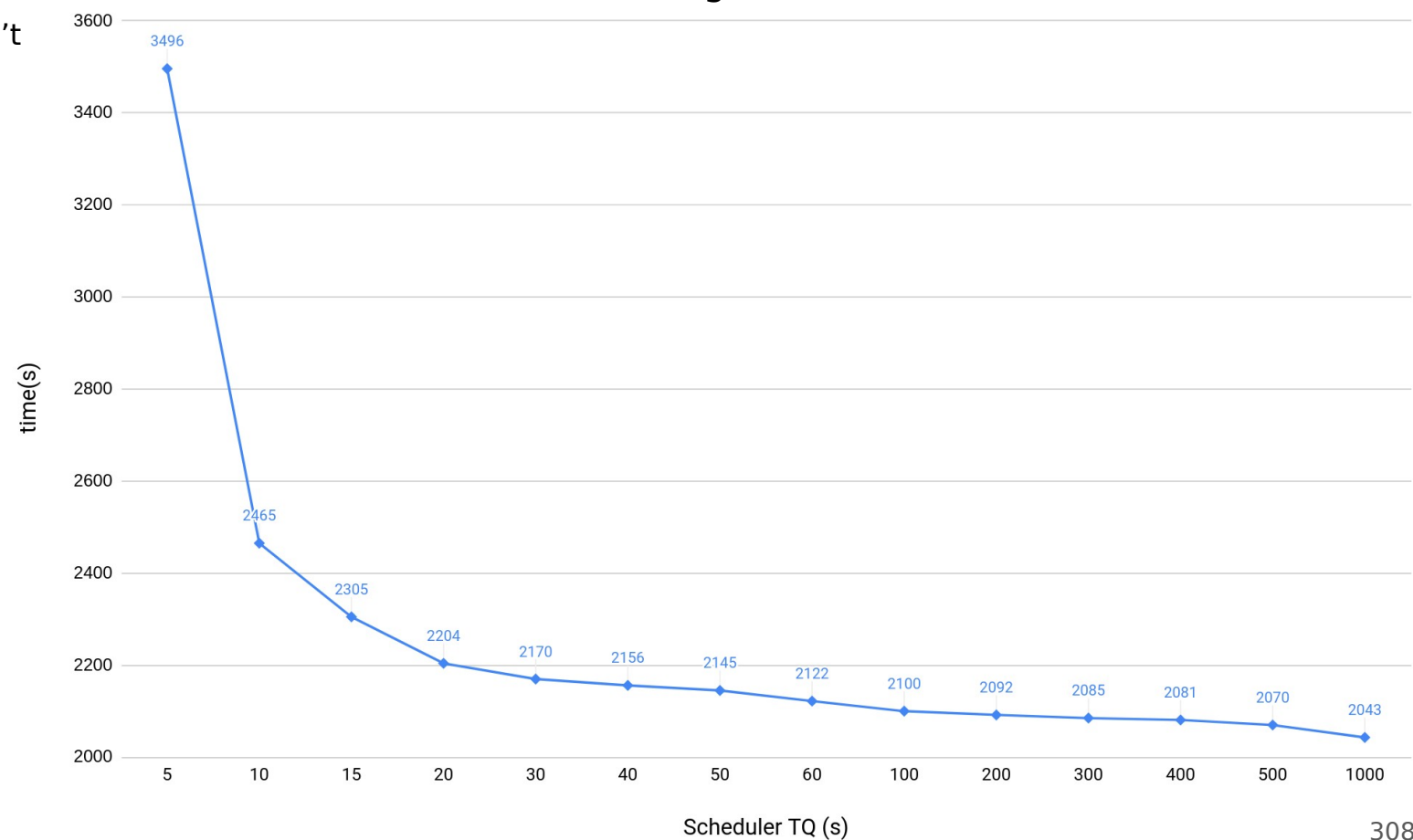
● Kubeshare can't run 2 x big !



# Results for big\_50

T(libnvshare, no anti-thrash) = 11757

● Kubeshare can't run 2 x big !

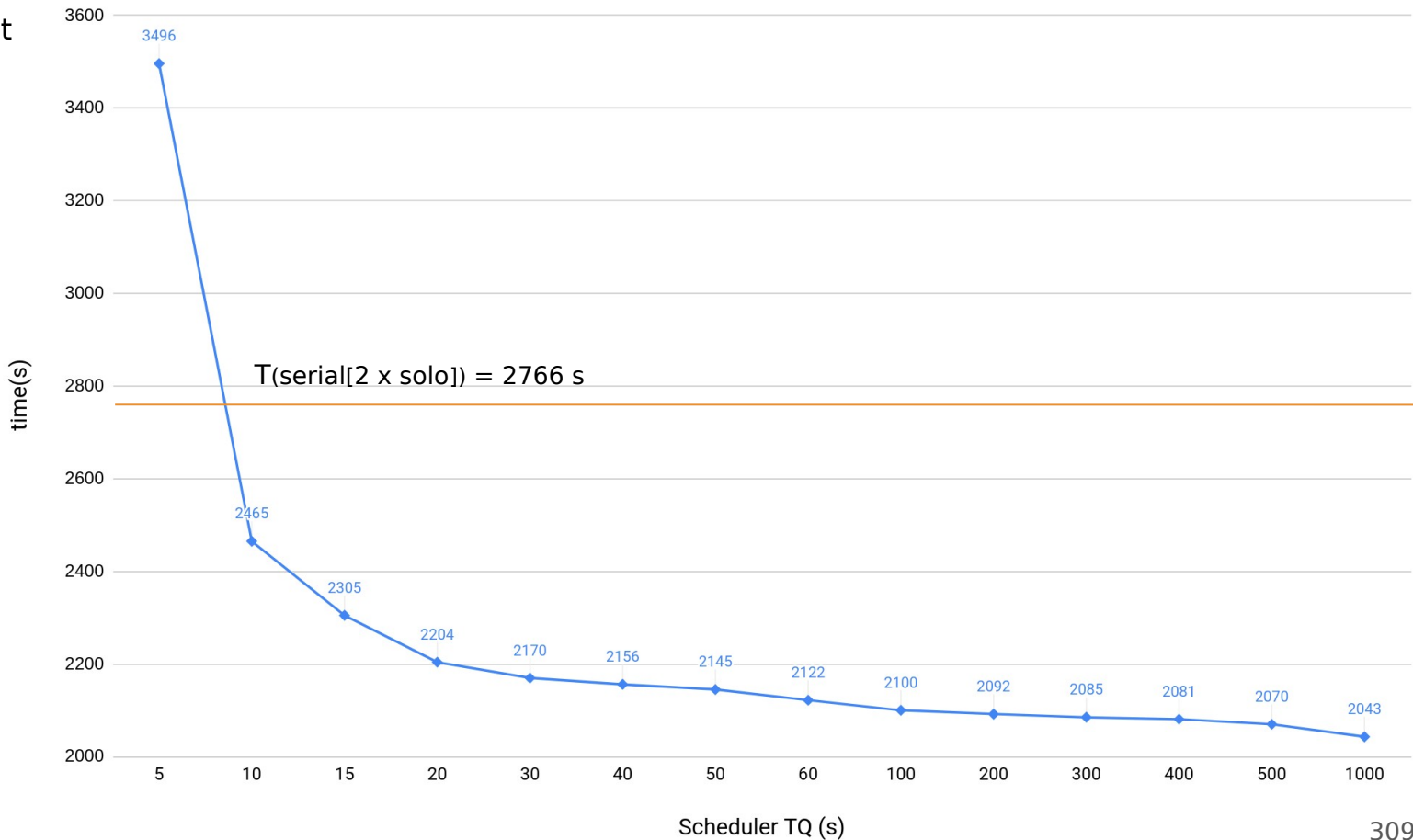


# Results for big\_50

T(libnvshare, no anti-thrash) = 11757

S

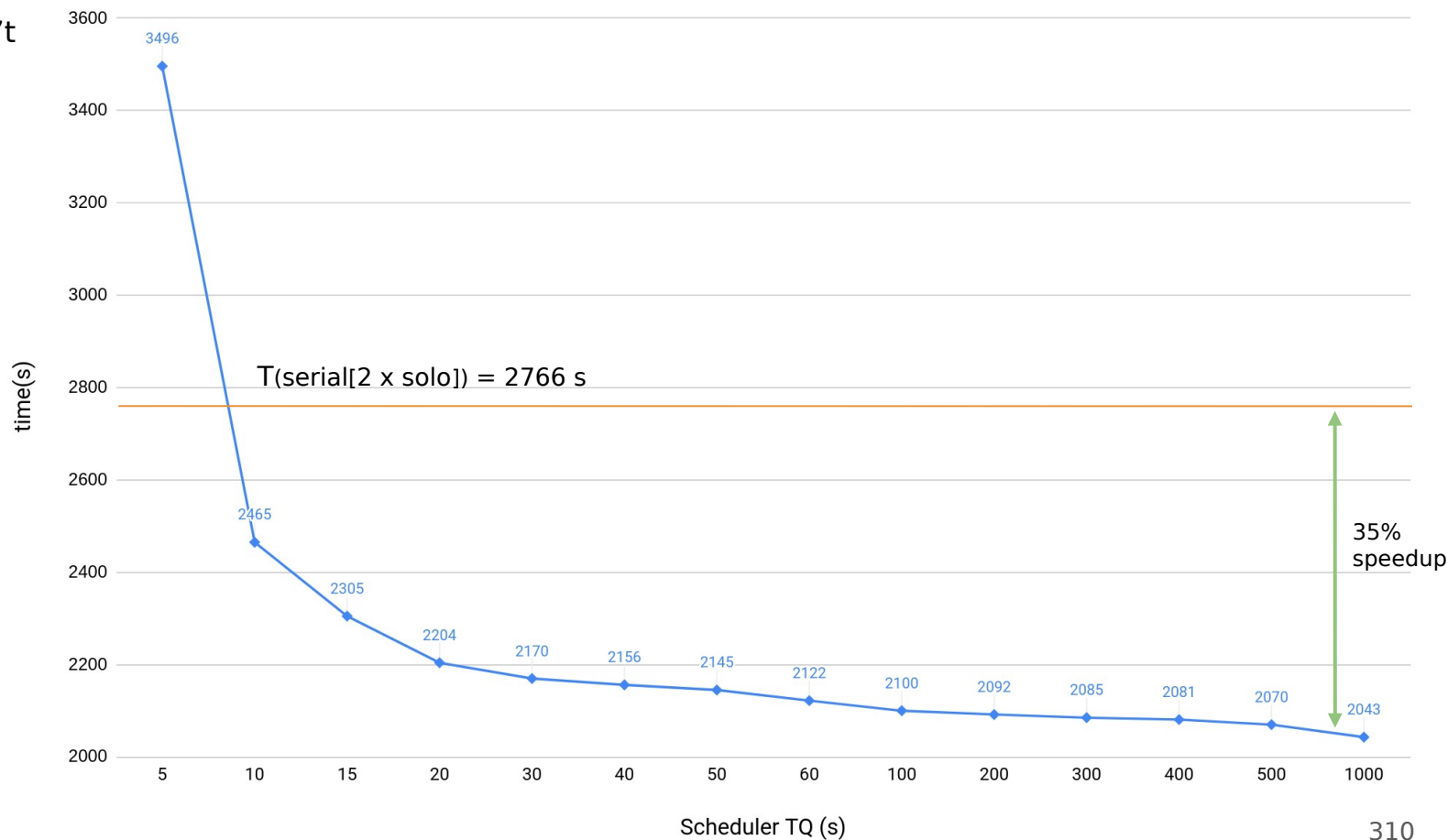
● Kubeshare can't run 2 x big !



# Results for big\_50

T(libnvshare, no anti-thrash) = 11757

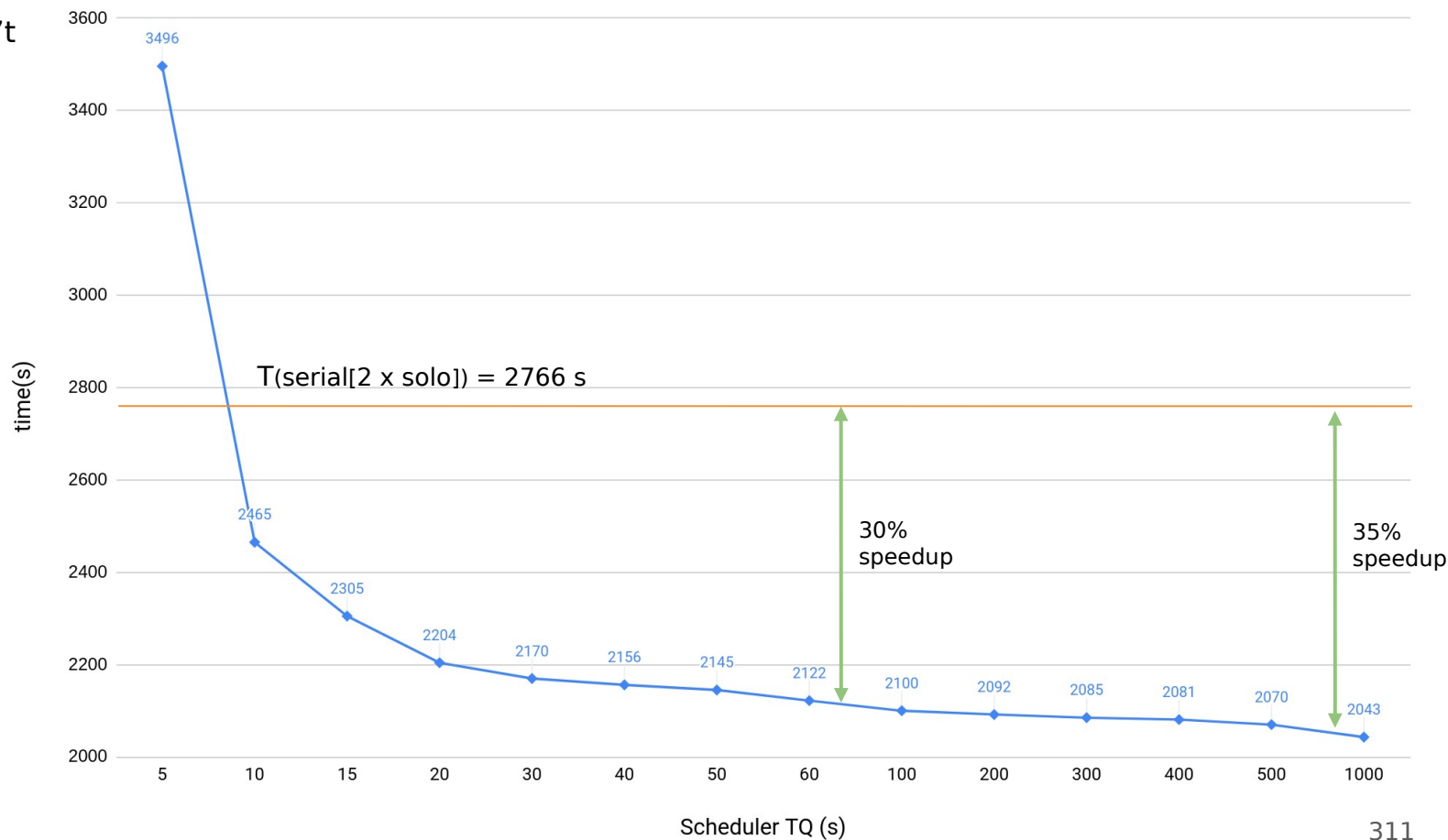
● Kubeshare can't run 2 x big !



# Results for big\_50

T(libnvshare, no anti-thrash) = 11757

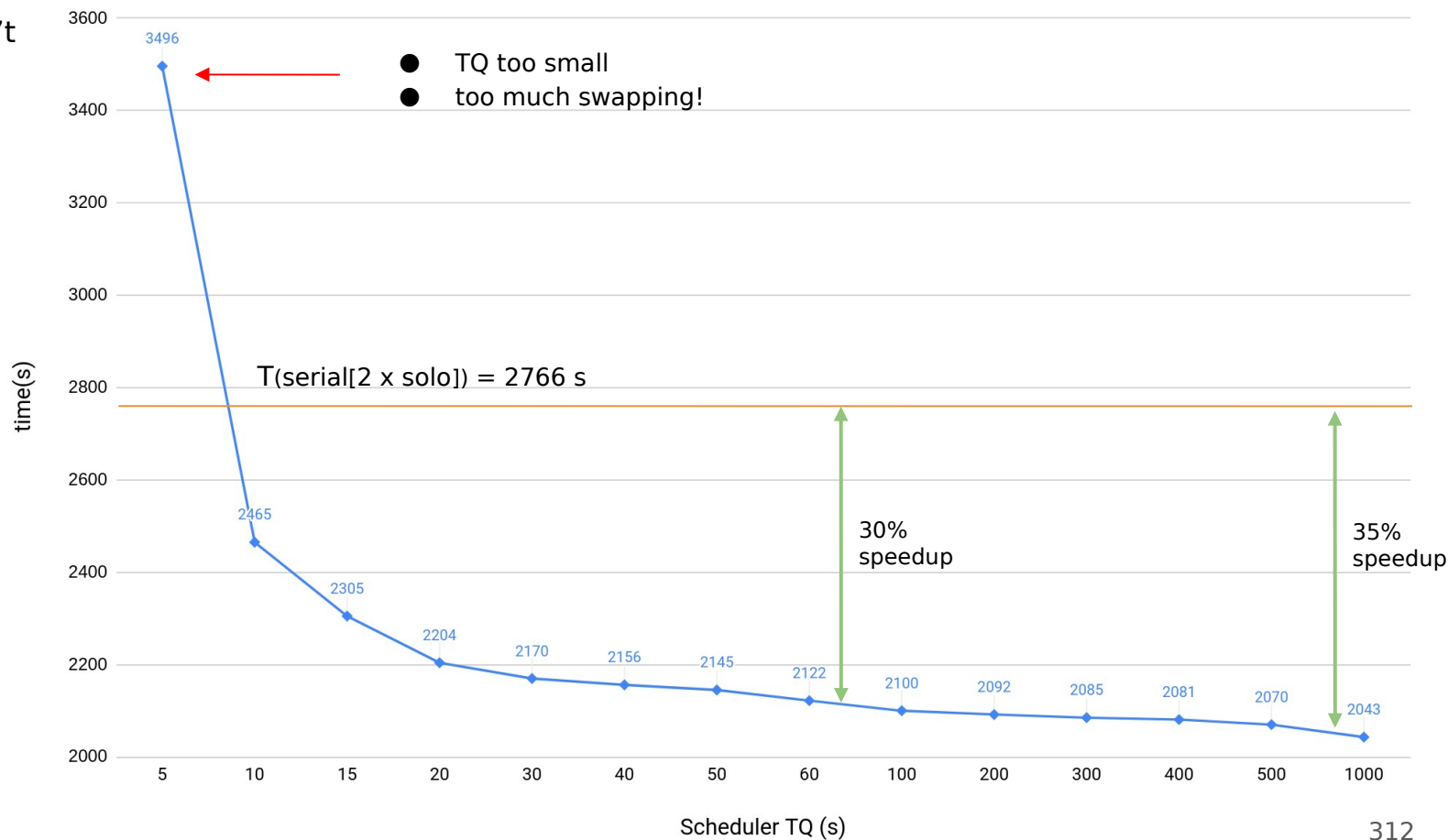
● Kubeshare can't run 2 x big !



# Results for big\_50

T(libnvshare, no anti-thrash) = 11757

● Kubeshare can't run 2 x big !



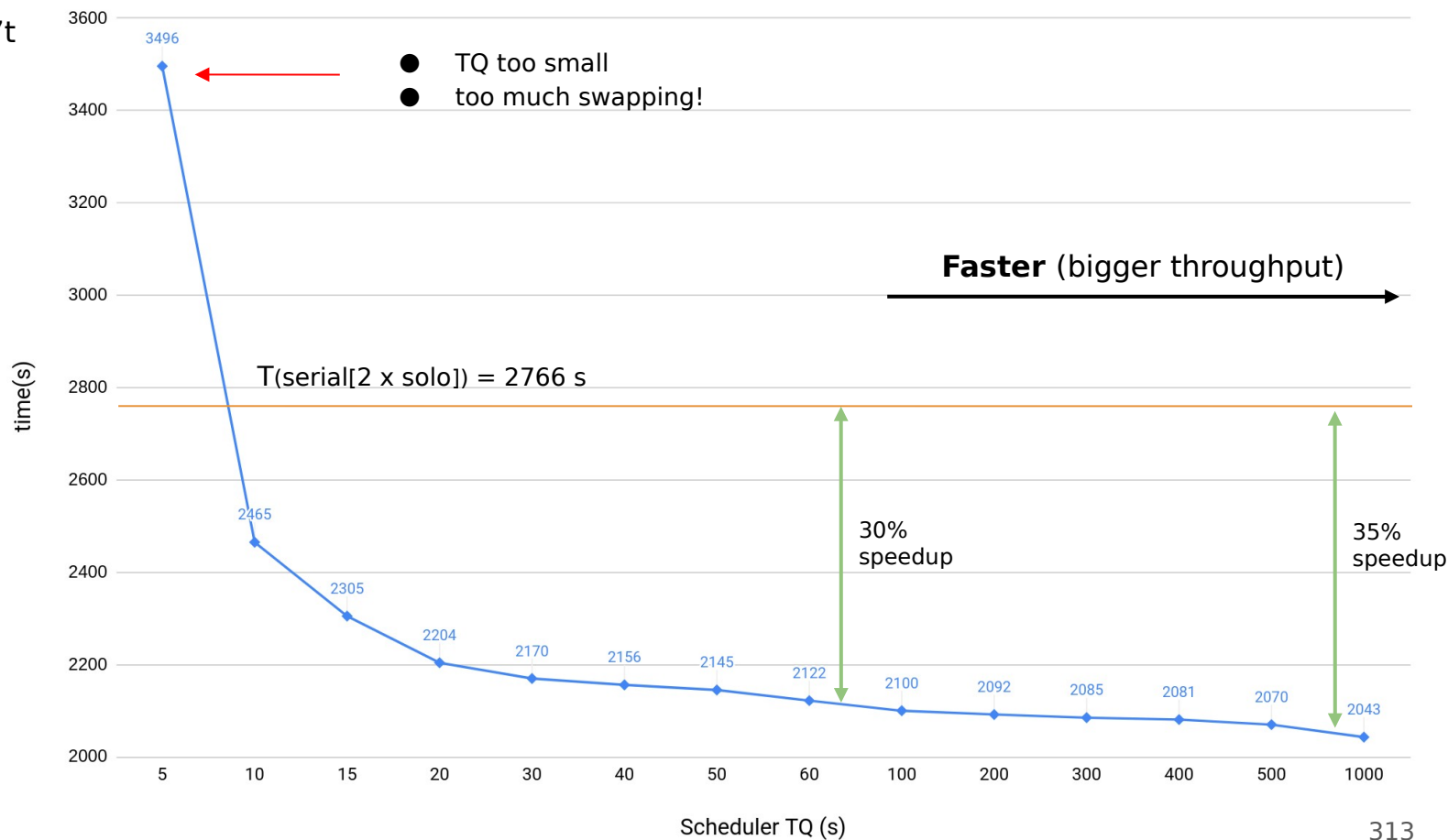
- TQ too small
- too much swapping!



# Results for big\_50

T(libnvshare, no anti-thrash) = 11757

● Kubeshare can't run 2 x big !



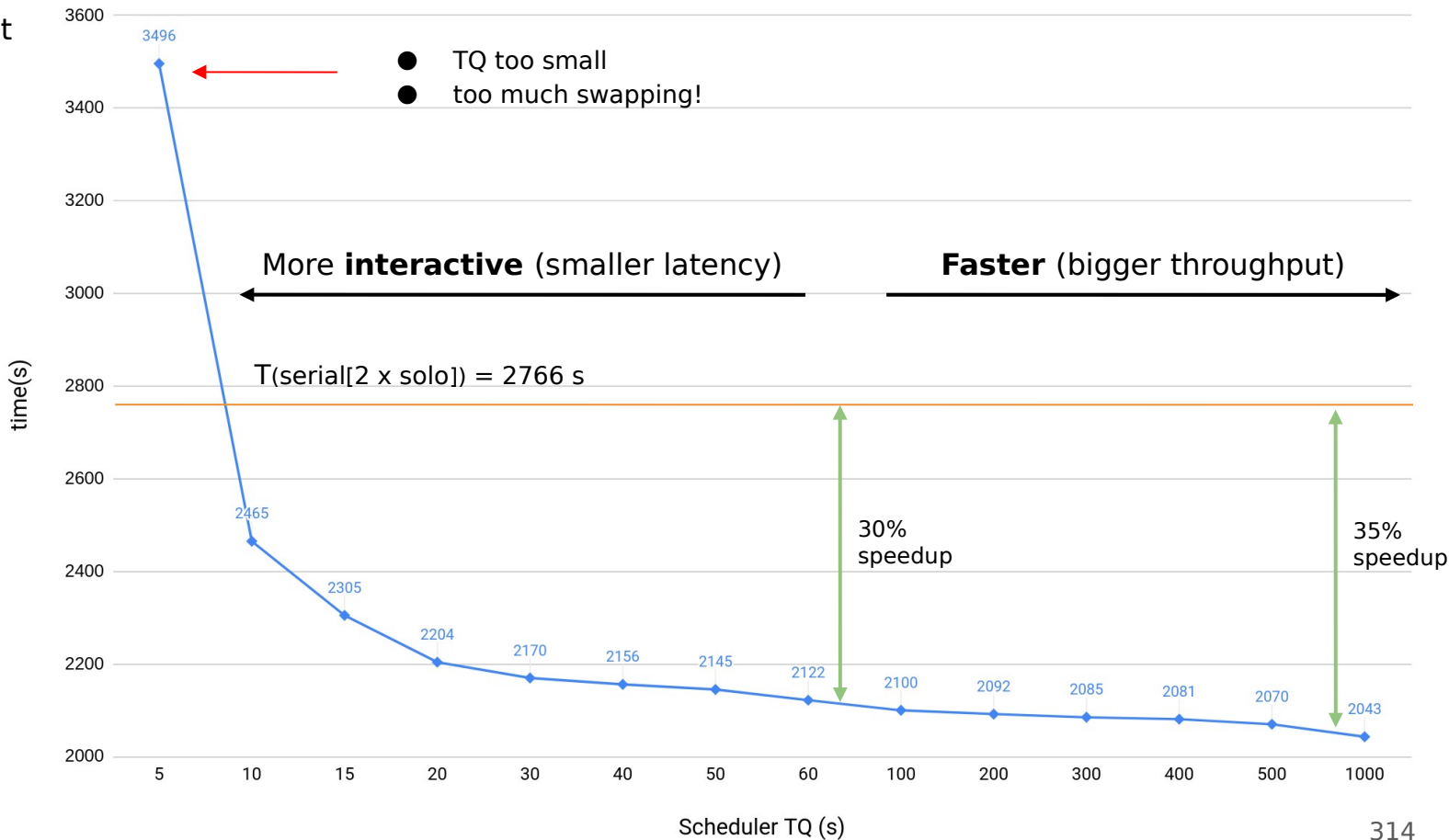
# Results for big\_50

T(libnvshare, no anti-thrash) = 11757

S

● Kubeshare can't run 2 x big !

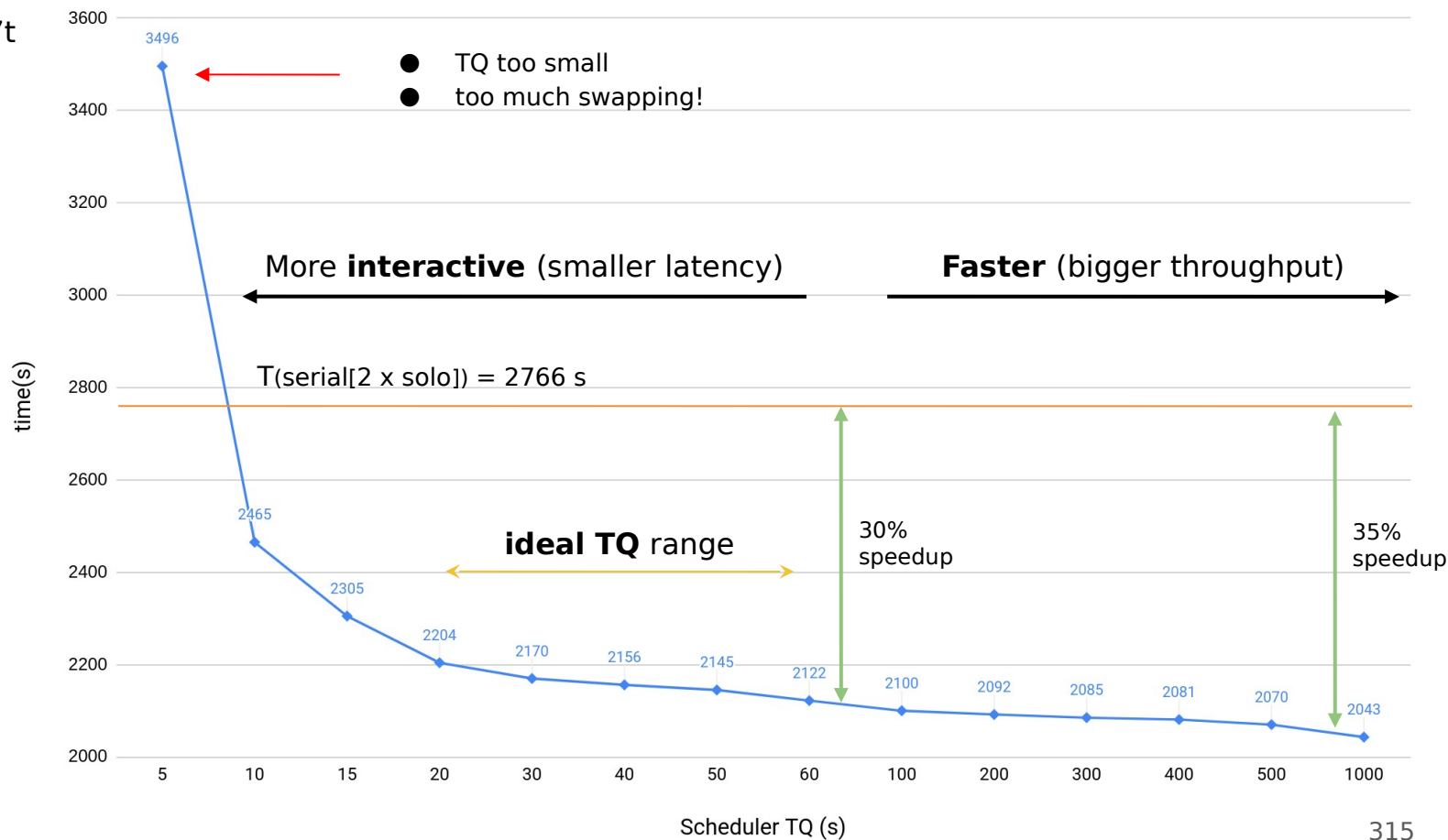
- TQ too small
- too much swapping!



# Results for big\_50

T(libnvshare, no anti-thrash) = 11757

● Kubeshare can't run 2 x big !



A new state of the art

# A new state of the art

- # of co-located processes limited by system RAM

# A new state of the art

- # of co-located processes limited by system RAM
- each process can use the whole GPU memory

# A new state of the art

- # of co-located processes limited by system RAM
- each process can use the whole GPU memory
- GPU utilization increases even for non-interactive jobs

# A new state of the art

- # of co-located processes limited by system RAM
- each process can use the whole GPU memory
- GPU utilization increases even for non-interactive jobs
- No code changes to user programs



# A new state of the art

- # of co-located processes limited by system RAM
- each process can use the whole GPU memory
- GPU utilization increases even for non-interactive jobs
- No code changes to user programs
- Sensible default policies:
  - anti-thrashing scheduler **always on**
  - default TQ = 20 sec

# Future Work

- Deploy to production environments
  - Support multiple GPUs per node
  - Create heuristics (e.g. PCIe traffic) to automatically enable/disable the anti-thrashing mechanism
  - Publish our work
- 
- Later on: Intra-node job migration from one GPU to another
  - Later on: Inter-node GPU job migration (Checkpoint-restart)

Thank you!  
Questions?