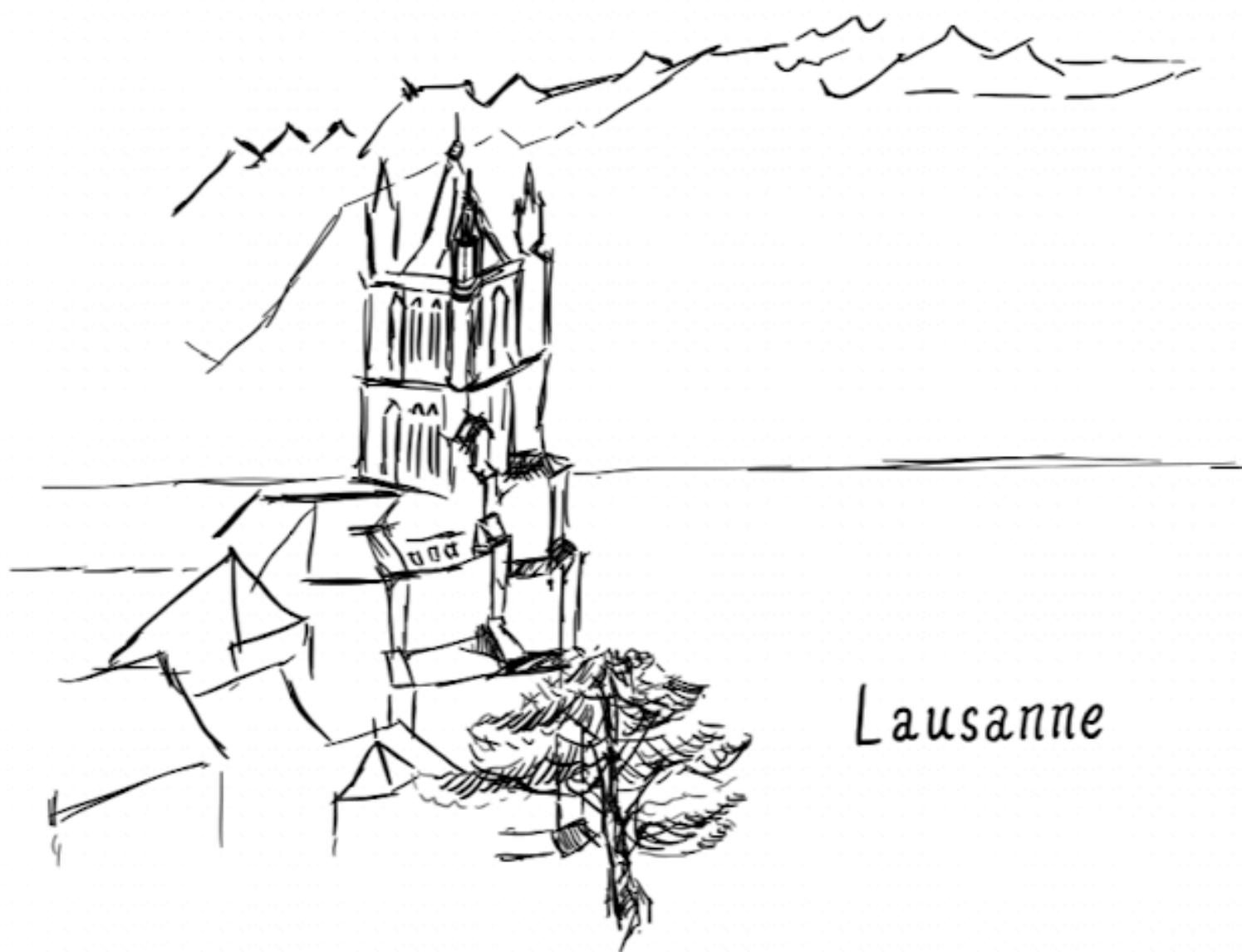


TMO = Transparent Memory Offloading ] in Datacenters  
IDCost = Block IO for Containers



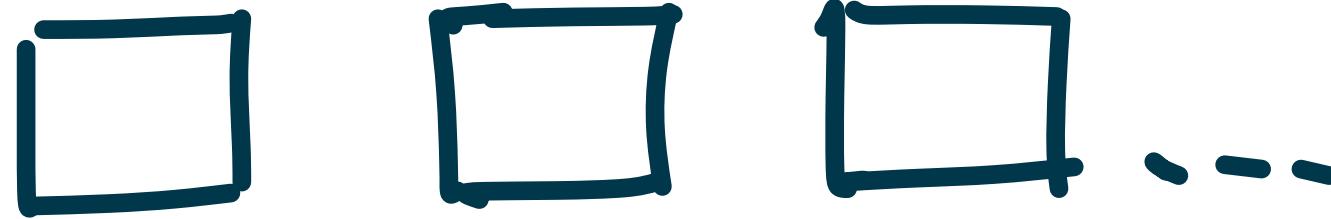
Two papers in ASPLoS 2022

# Memory Offloading

[DRAM]

Apps consume a lot of memory

↓ offload

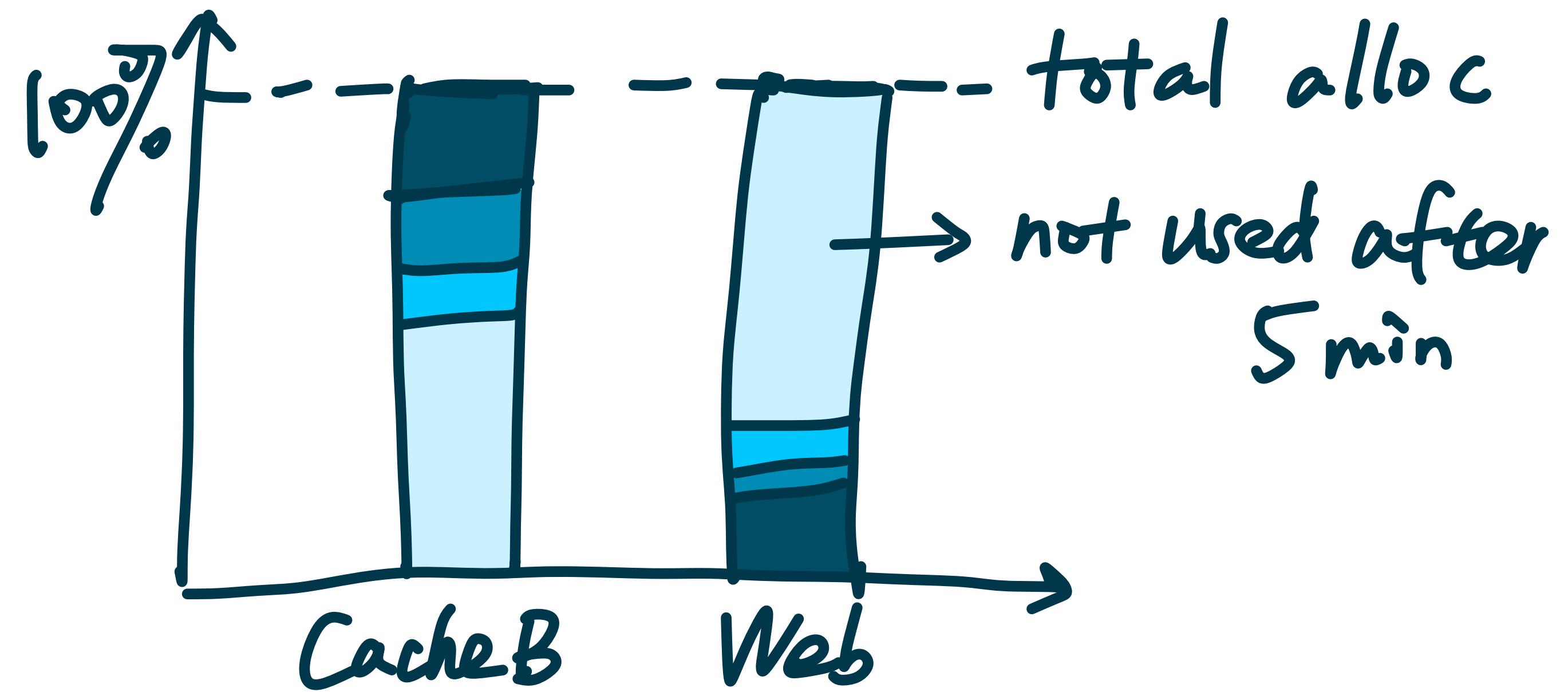


compressed  
mem

NUM  
SSD

NUM

} ⊕ [memory tiering]



# State-of-the-art & Issues

zswap (Google, ASPLOS 19)

- ① single slow memory tier (compressed mem)
- ② offline application profiling  
metric = page-promotion rate

# State-of-the-art & Issues

## zswap (Google, ASPLOS 19)

① single slow memory tier (compressed mem)

⇒ add NVMe SSD → cheaper, power

⇒ Compression not applicable (Sometimes)

② offline application profiling

metric = page-promotion rate

App's sensitivity

⇒ Handle dynamic + diverse apps

differ ⇒ e2e

⇒ Handle device heterogeneity (diff SSDs)

perf

slow device implies low PPR

# TMO

{ Transparent memory offloading  
  { containerized environment => interesting setting

key questions { How much memory to offload?  
                  { What memory to offload?

=> mechanism for better information

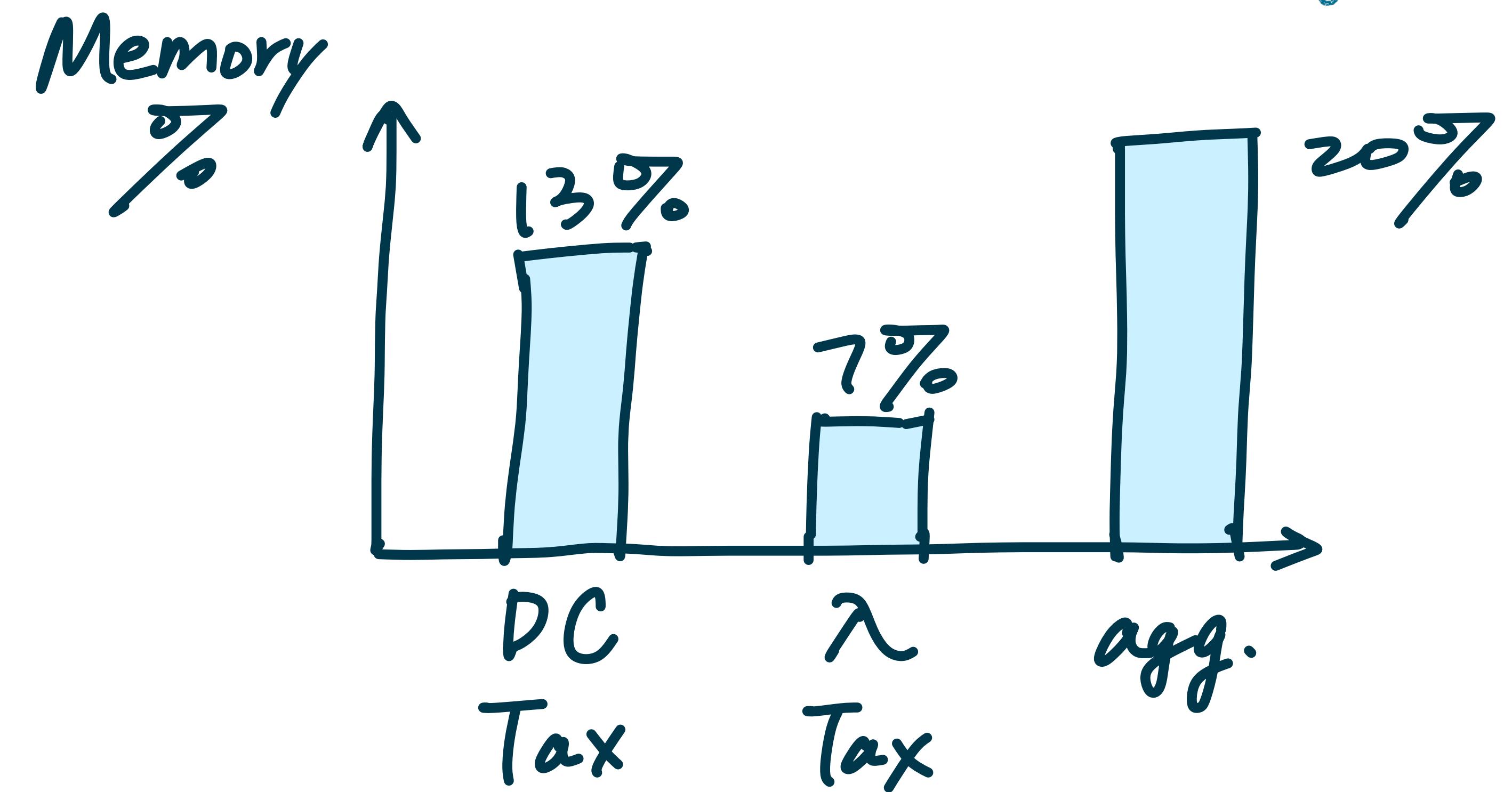
=> perf model captures: app-centric slowdown  $\oplus$  diff devices

=> framework + algorithm to adjust at runtime

# Data center Setting / Deployment $\oplus$ Challenges

Memory Tax: { infrastructure-level functions  
microservices  
routing, proxy...}

Packaging, logging  
profiling - - -

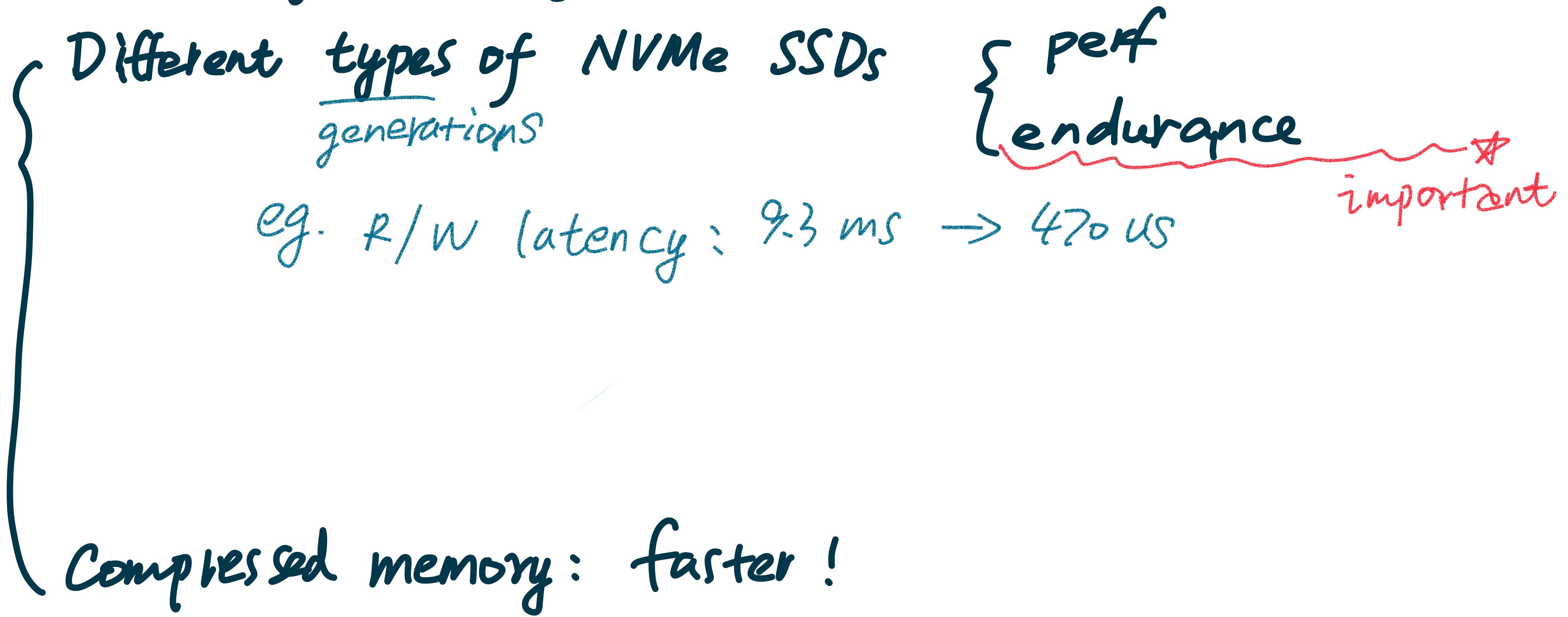


memory tax SLA = relaxed

↓  
primary target of  
offloading (first version)

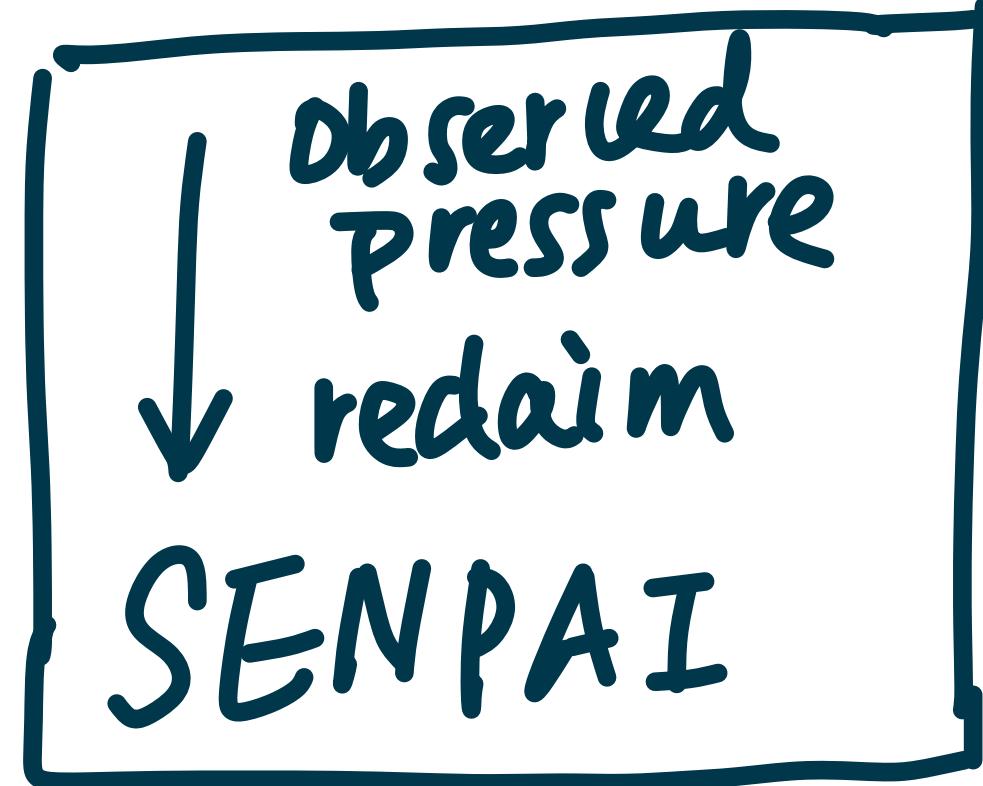
# Data center Setting / Deployment $\oplus$ Challenges

HW heterogeneity of offload backend

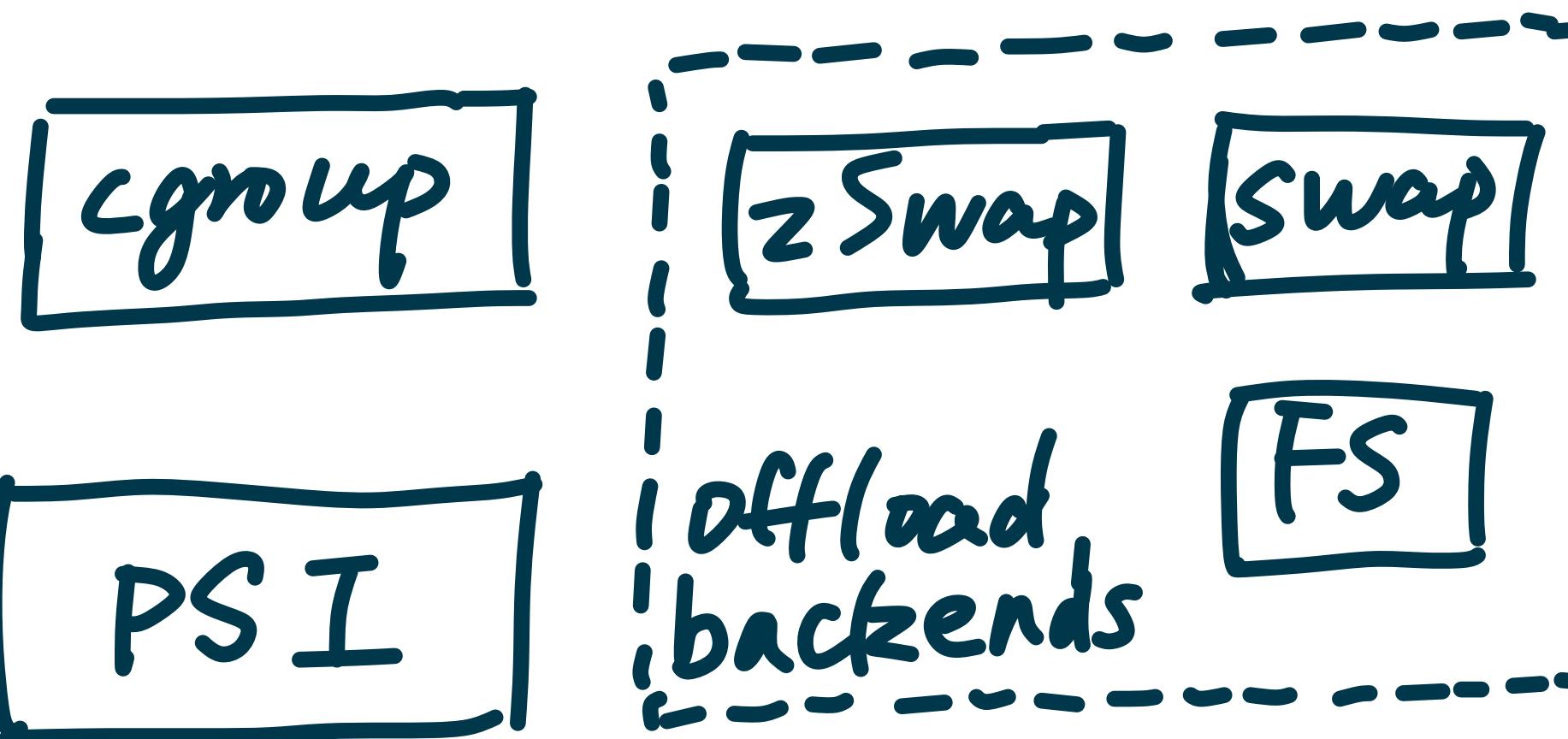
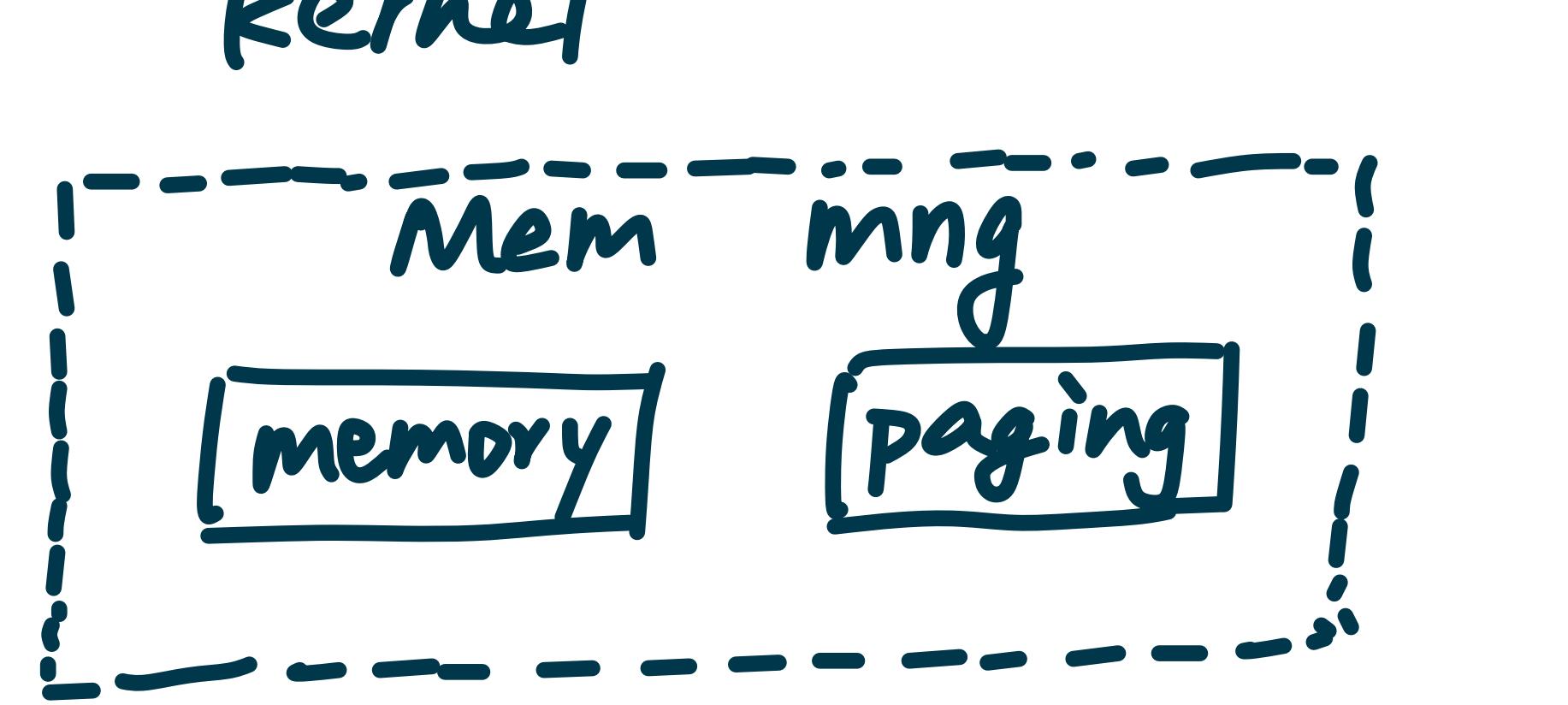


# TMO Design

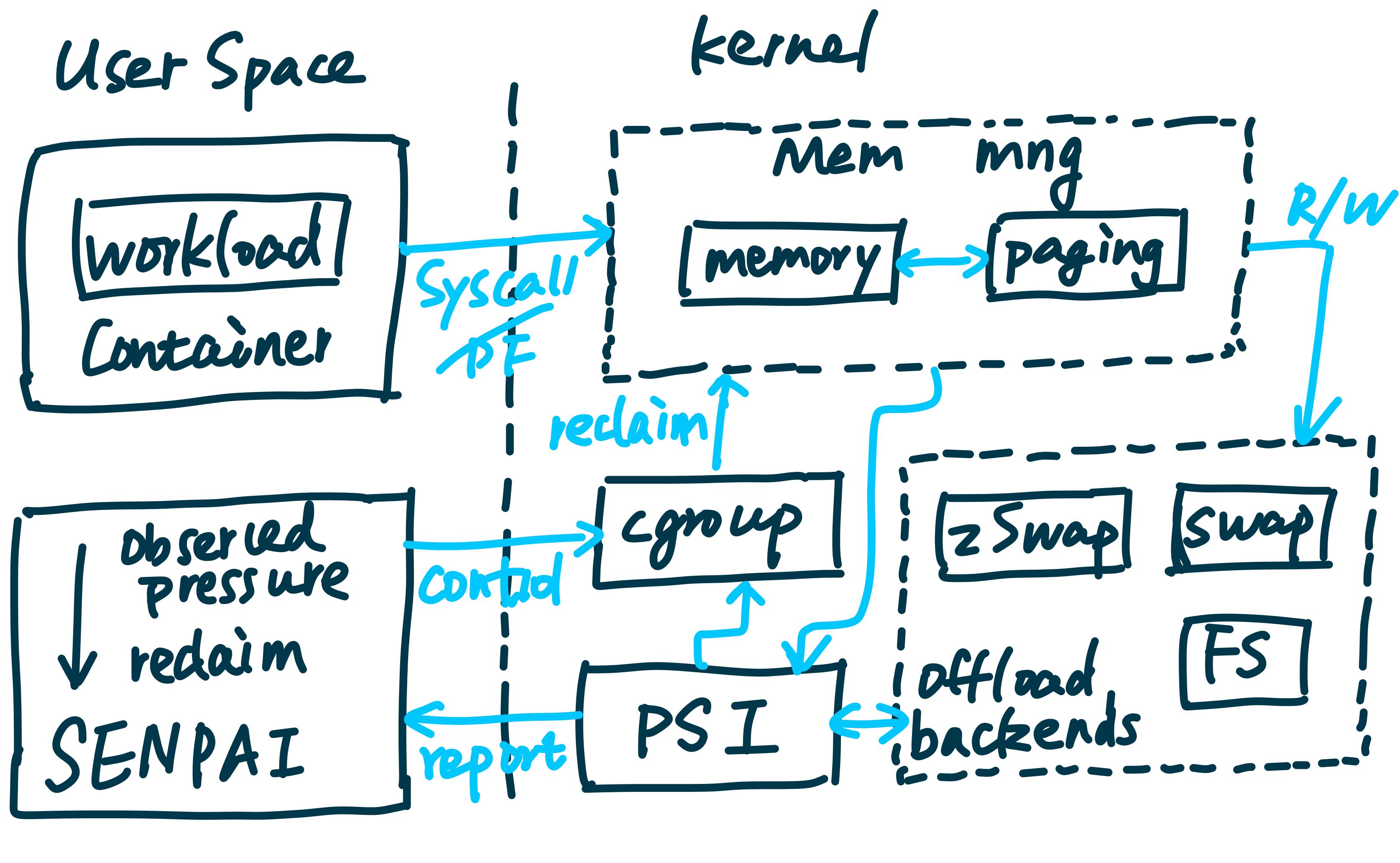
User Space



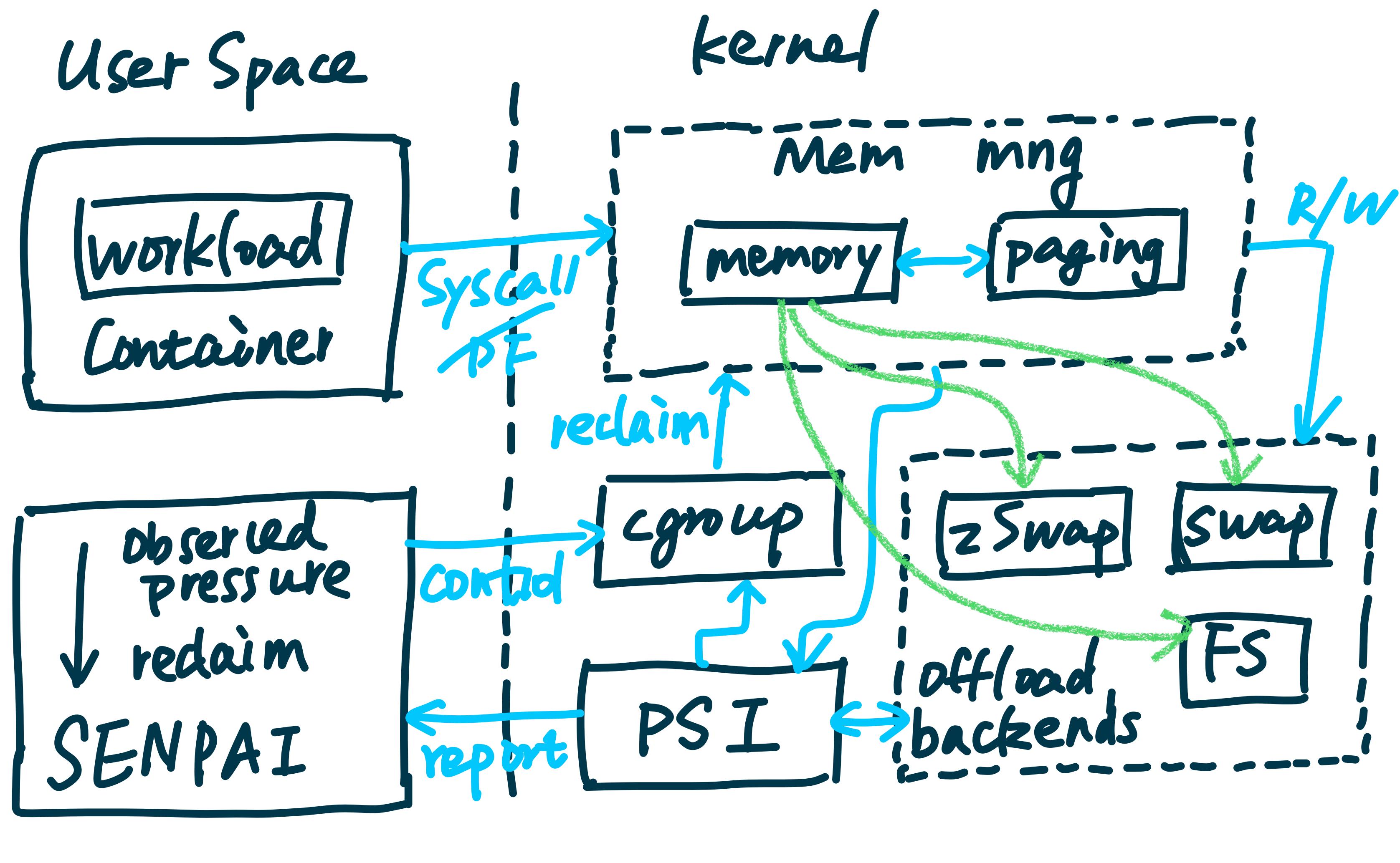
kernel



# TMO Design

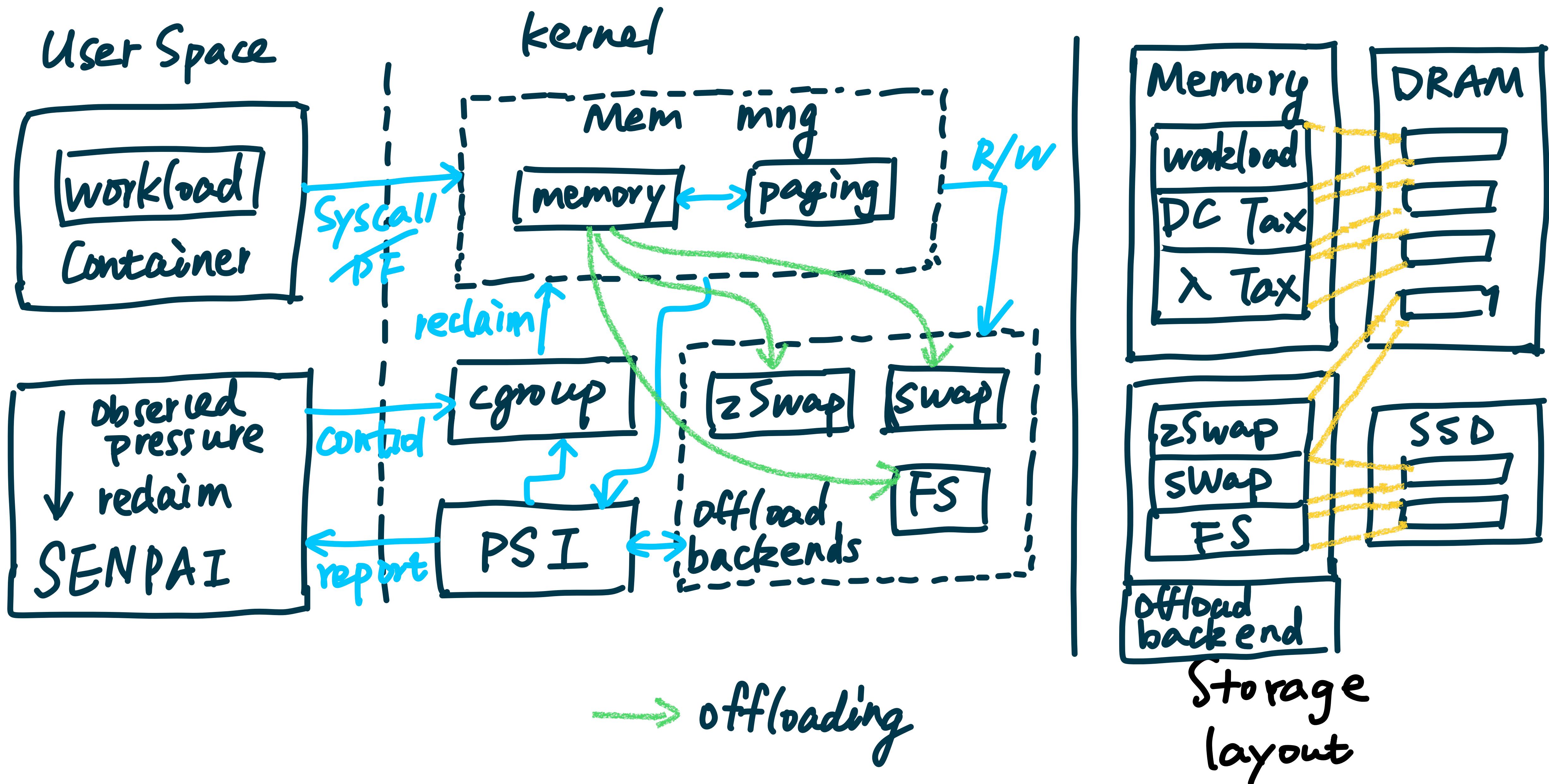


# TMO Design



→ offloading

# TMO Design



# Resource Pressure : Pressure Stall Info (PSI)

How much memory to offload?

metric: predict the resulting performance after offloading

why hard? { kernel's information is indirect  $\Rightarrow$  mechanism support  
single metric is hard

accurately attribute performance penalty to offloading

eg. workload start up

workload working set transition

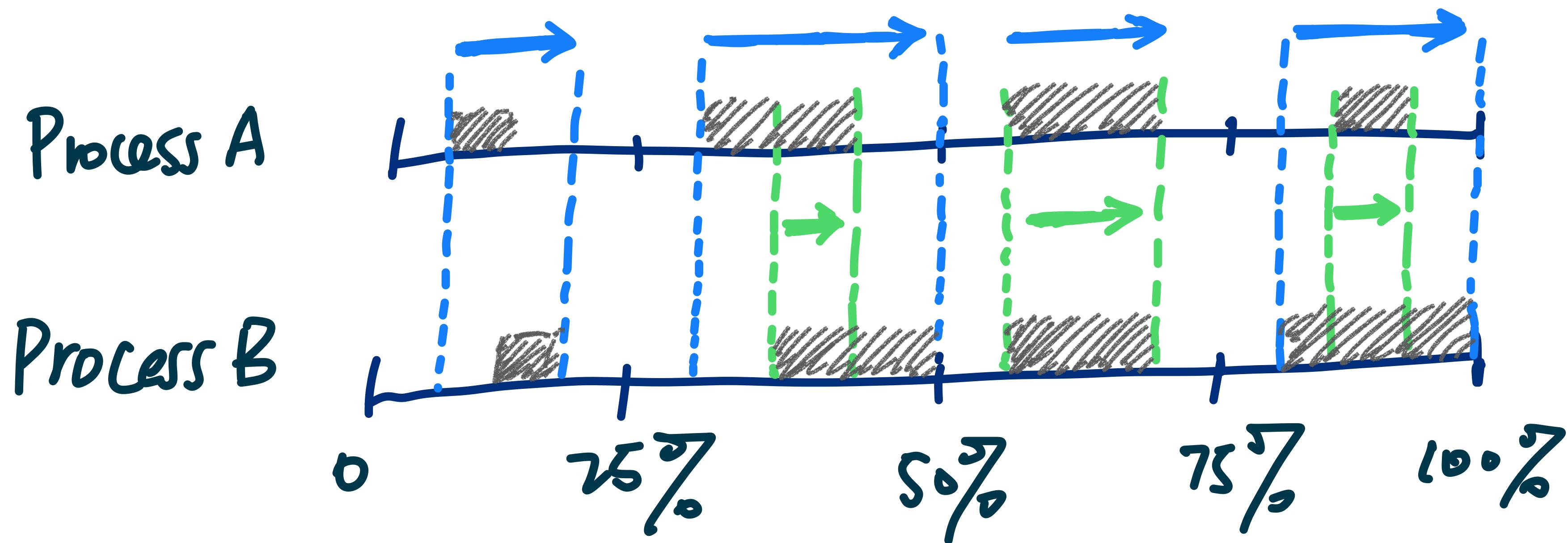
...

mechanism to provide information & accounting to  
container (in kernel)

# Resource Pressure = Pressure Stall Info (PSI)

How much memory to offload?

metric: predict the resulting performance after offloading

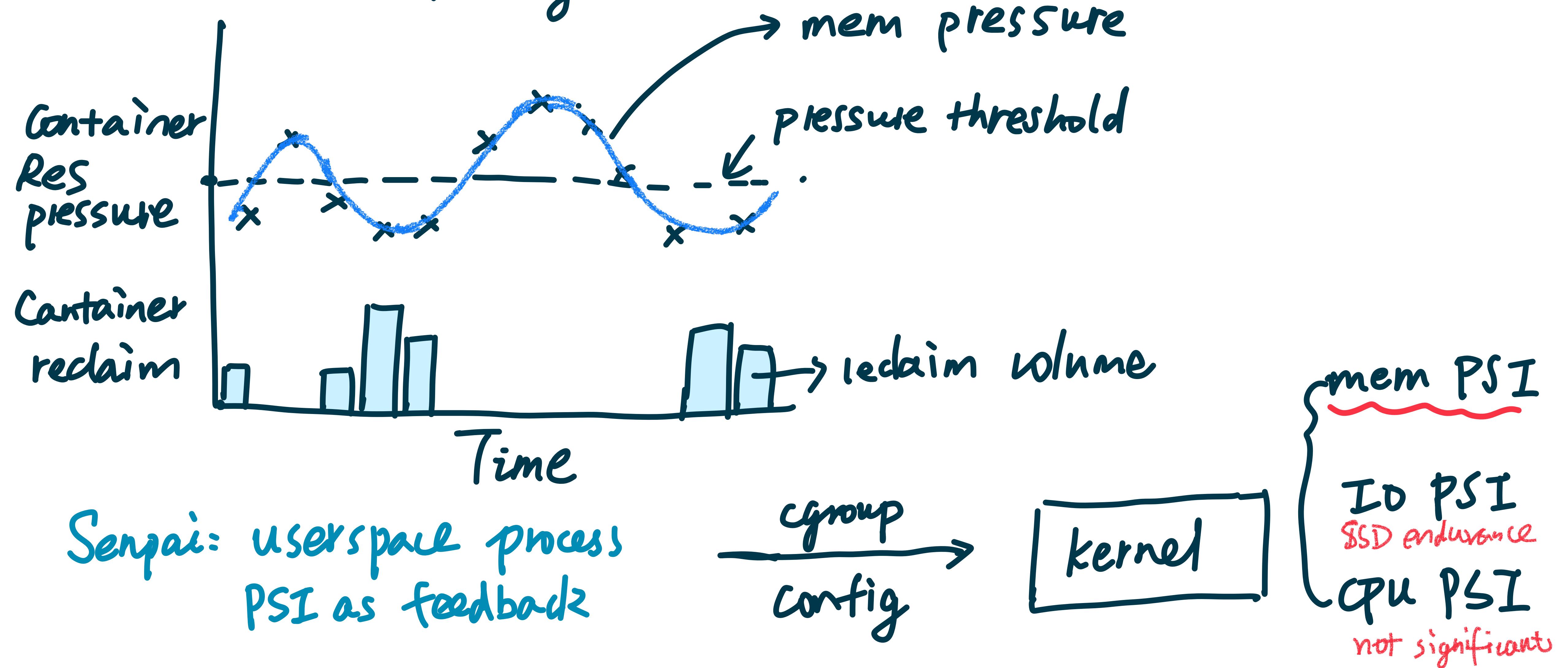


- Some: One process stall for resource
  - Full: All processes stall for resource
- Stall

# Determine Memory Requirement of Containers

How much memory to offload?

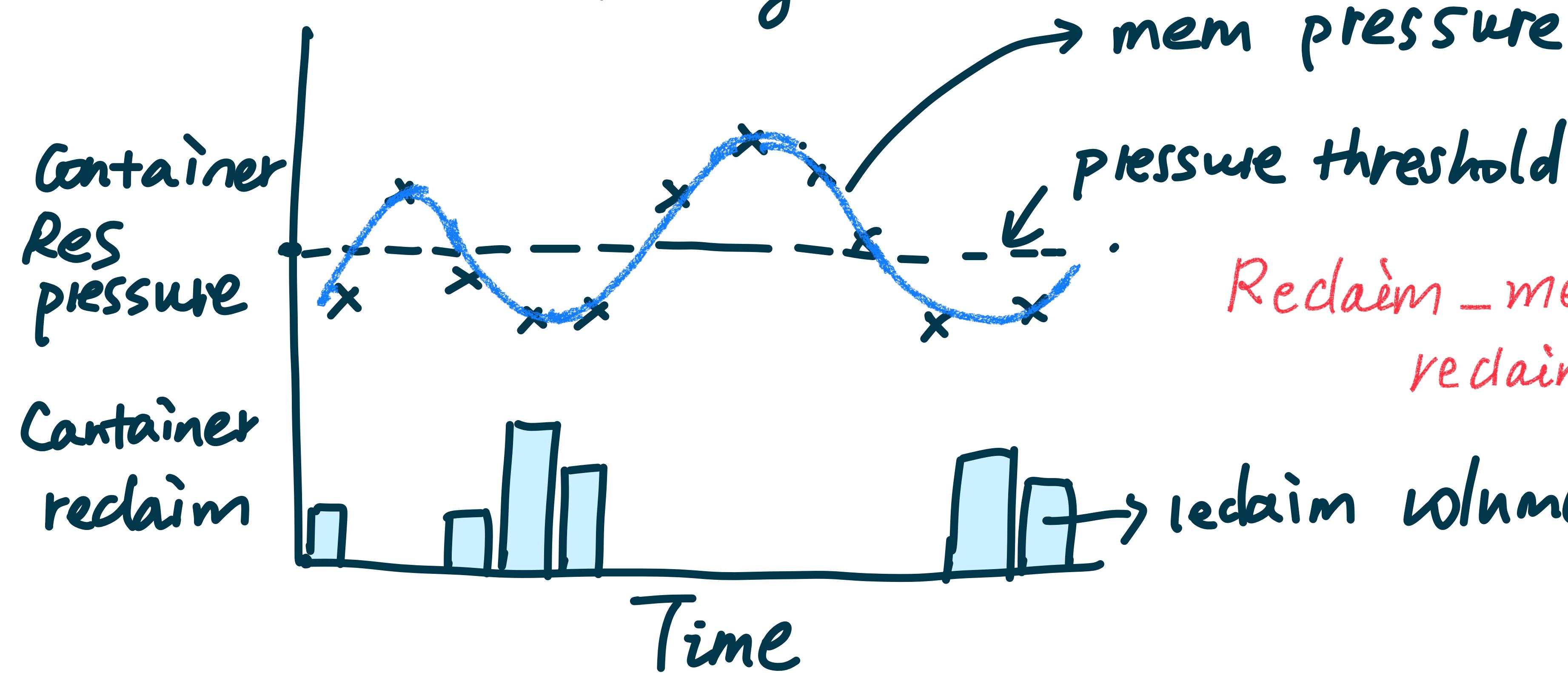
x PSI tracking



# Determine Memory Requirement of Containers

How much memory to offload?

× PSI tracking



$$\text{Reclaim\_mem} = \text{current\_mem} \times \text{reclaim\_ratio} \times \max(0, 1 - \frac{\text{PSI}_{\text{some}}}{\text{PSI}_{\text{thr}}})$$

Senpai: userspace process  
PSI as feedback



# Kernel Optimization for memory offloading

What memory to offload?

mechanism: kernel's reclaim

issue: kernel's swapping is very conservative  
skewed heavily towards file cache

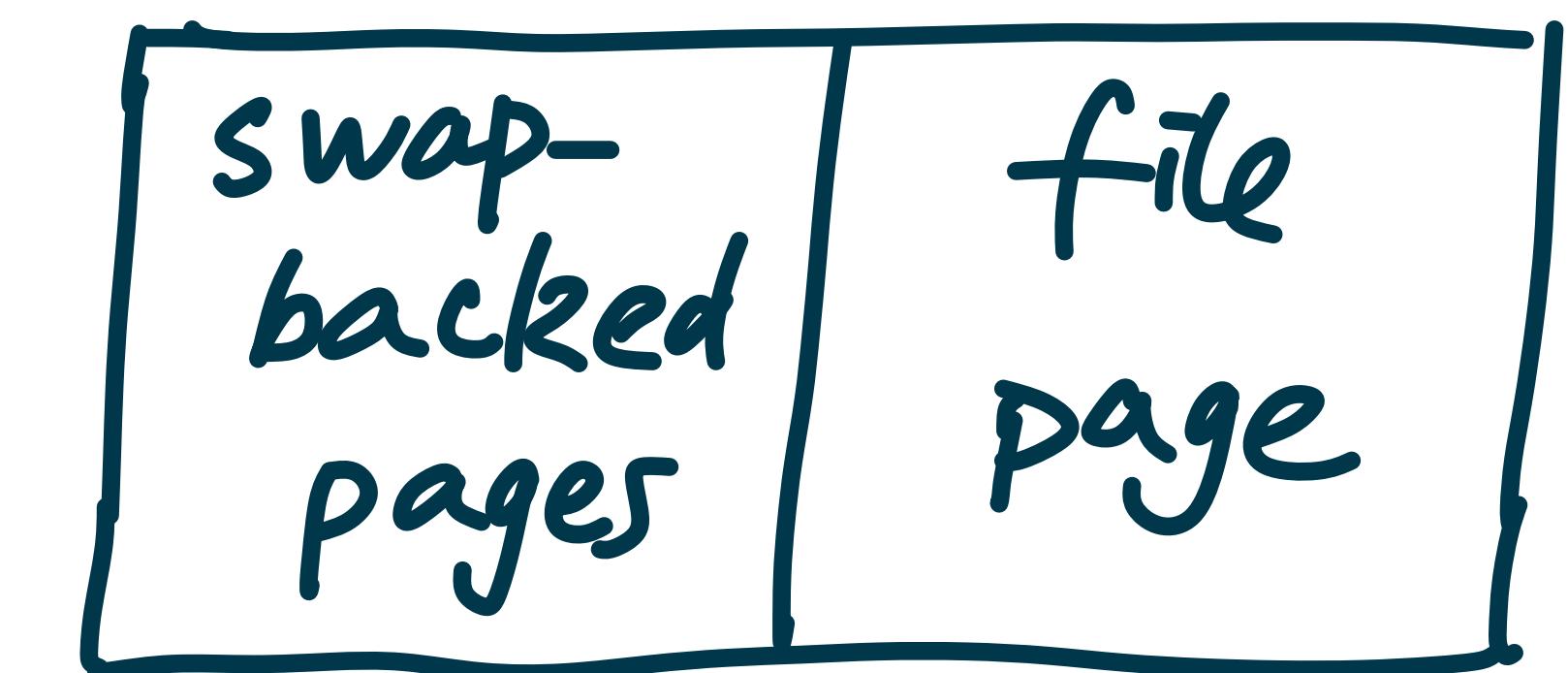
+ refault-detected mechanism

non-resident cache tracking

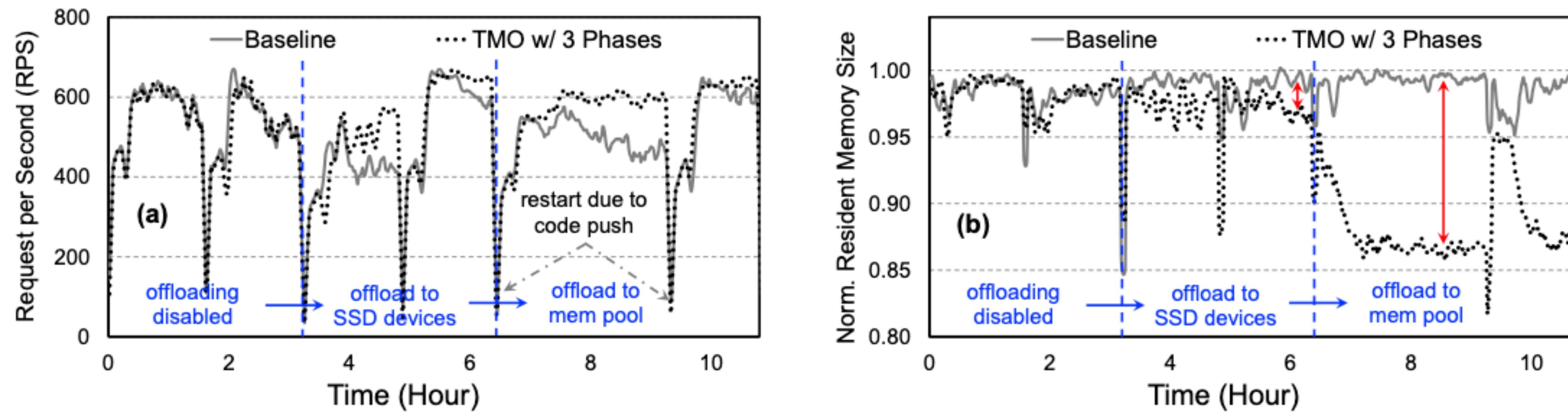
⇒ reuse-distance

⇒ exclude first-time-accessed  
file cache

+ if refault happens { refault rate  
{ swap-in rate } balance }

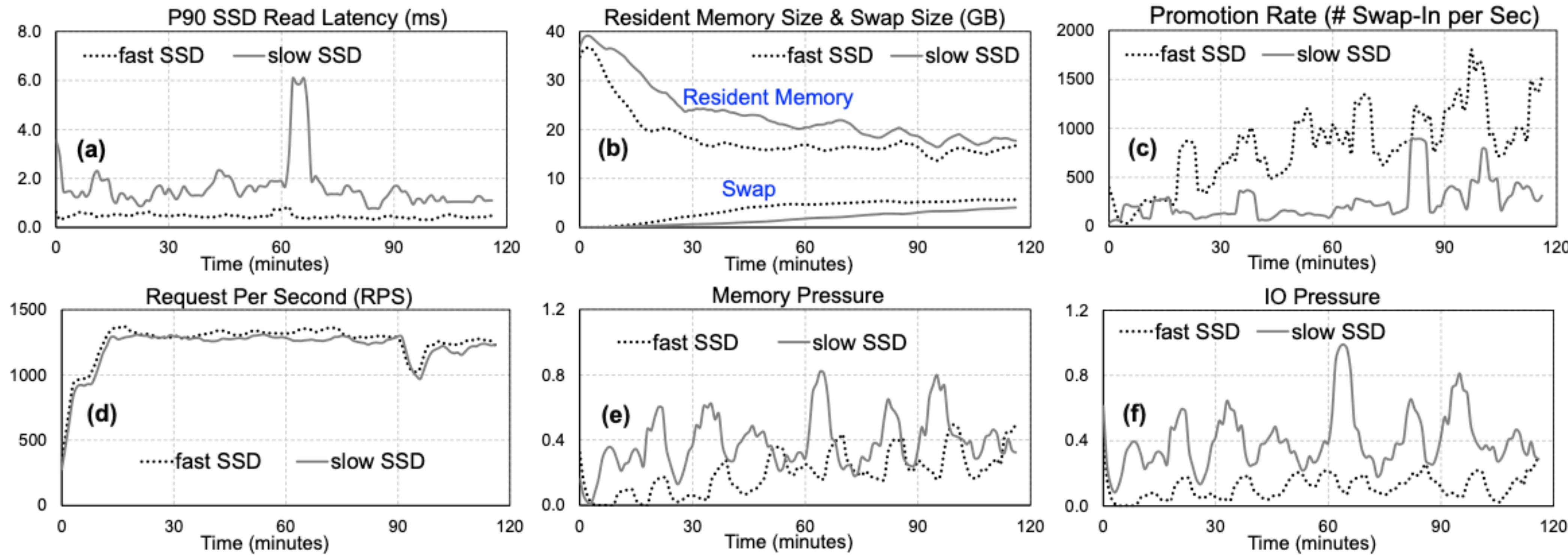


# TMO - Eval: Performance + memory Saving



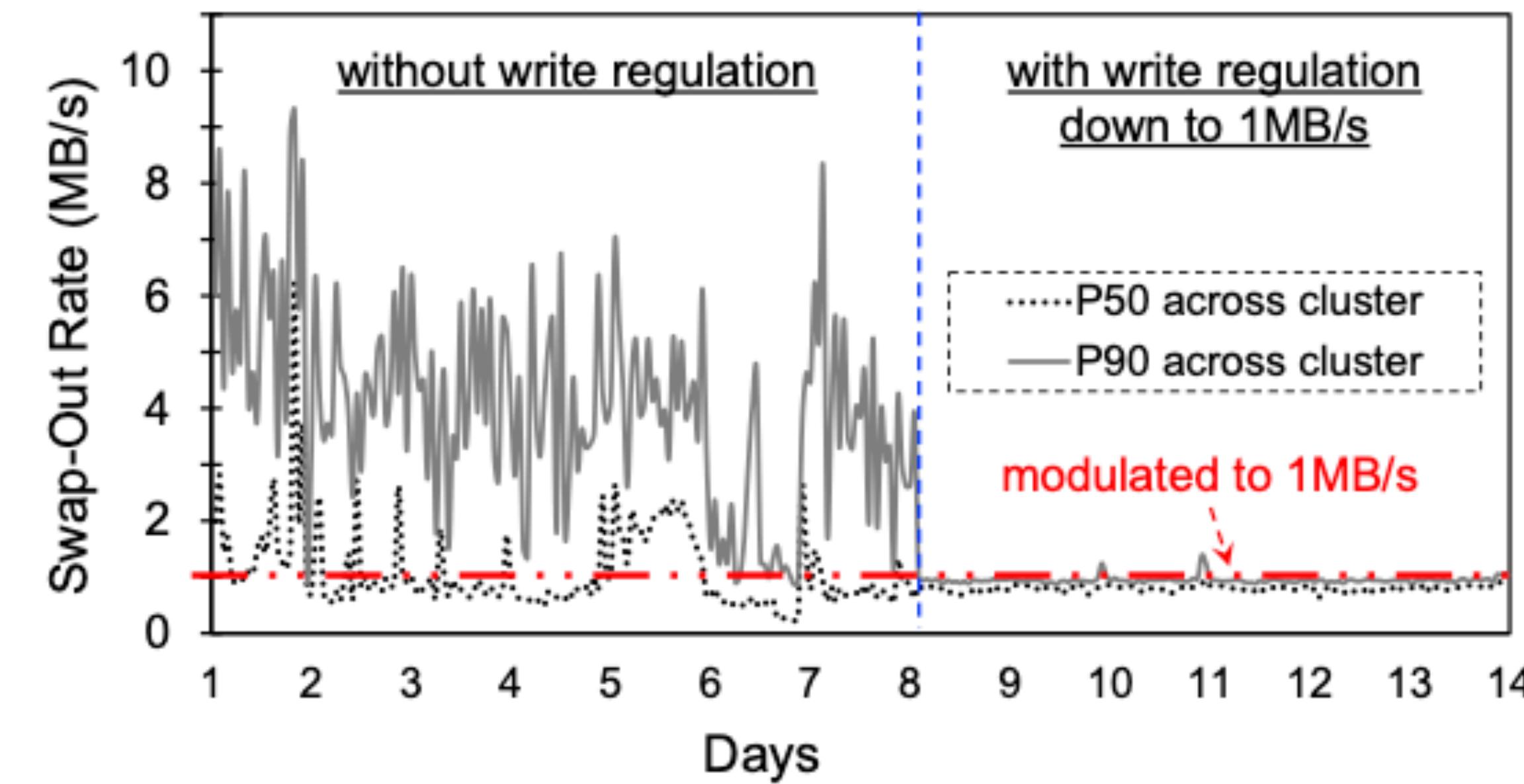
**Figure 11: Performance improvement and memory saving for the Web application on memory-bound hosts. We set up two tiers: started with identical configuration without any swap enabled and later enabled SSD & then compressed memory offloading on the 2nd tier.**

# TMO - Eval: PSI

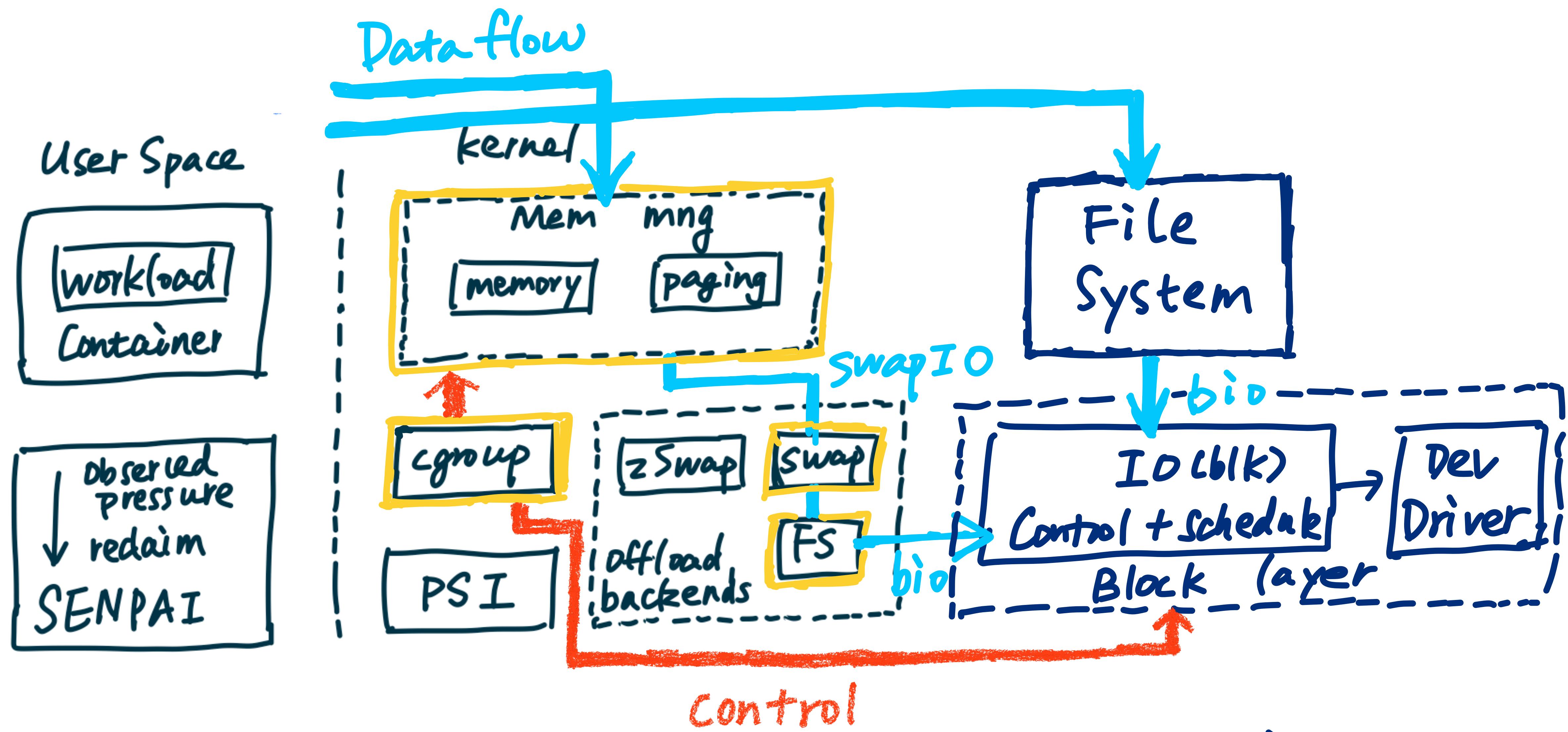


**Figure 12: The Web application under the control of TMO with fast and slow SSDs. The solid-gray configuration has lower RPS/performance but actually lower promotion/swap-in rate; while its PSI metrics are indeed higher.**

# TMO - Eval: Device Endurance



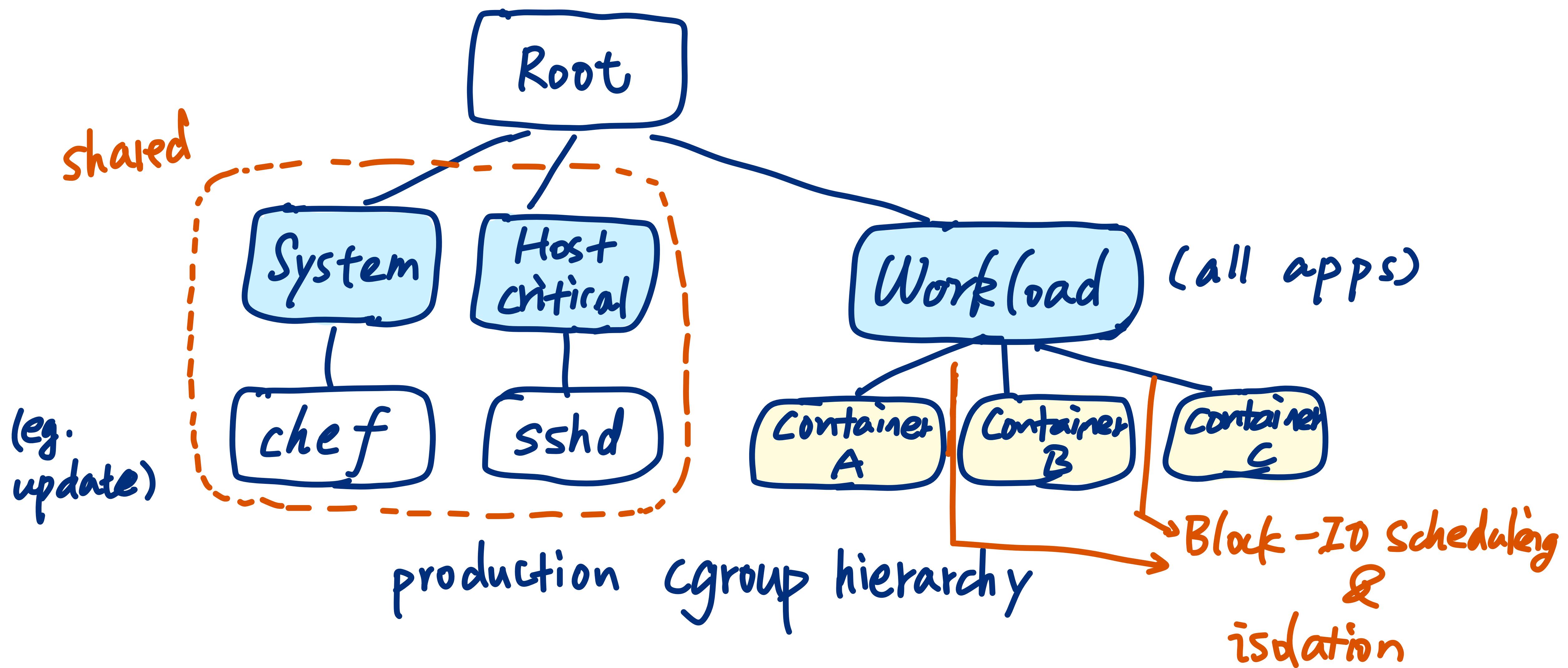
**Figure 14: Swap rate with and without write regulation.**



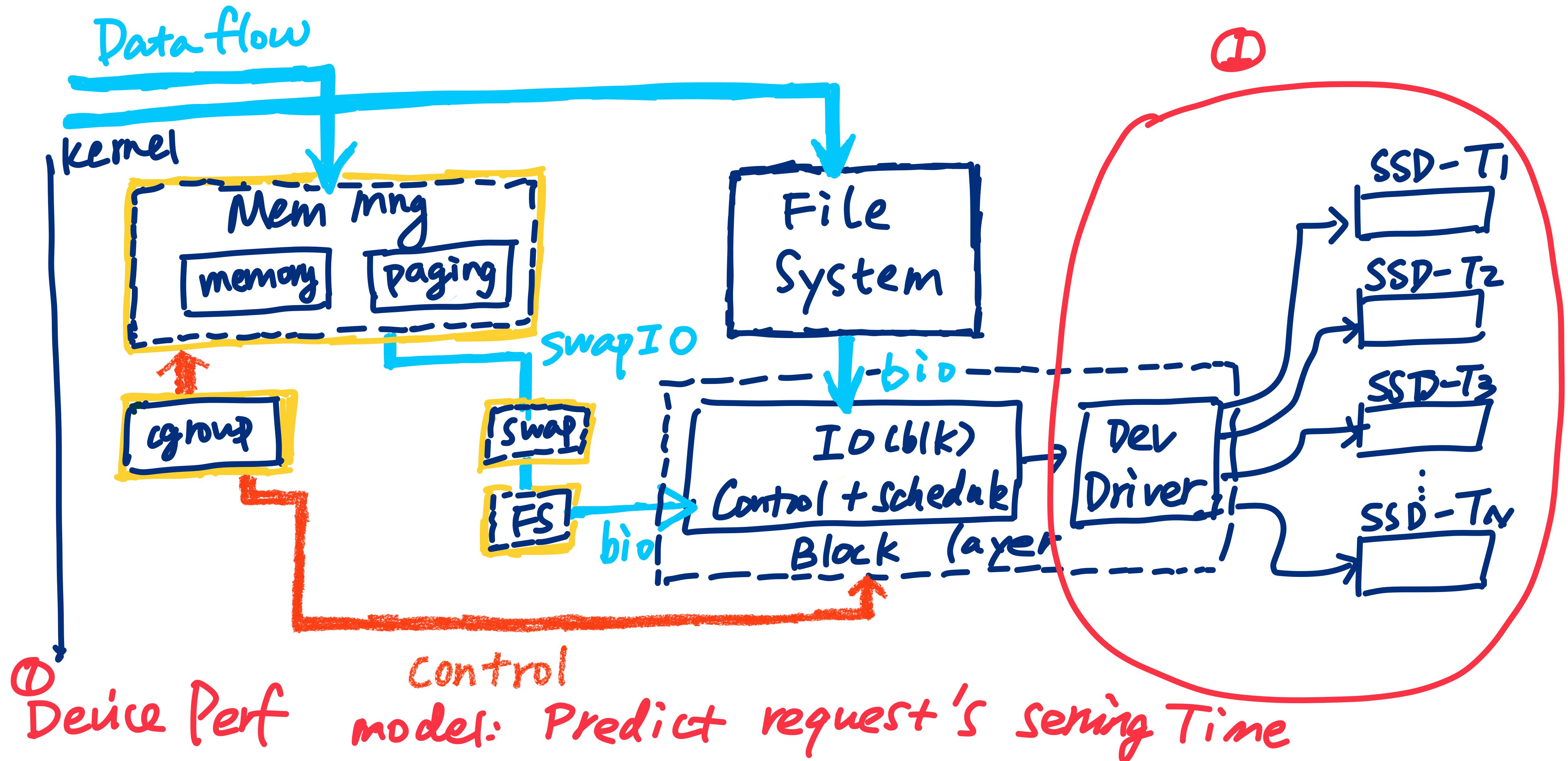
# IO cost: Block IO control

{ PC containers  
heterogeneity in Devices

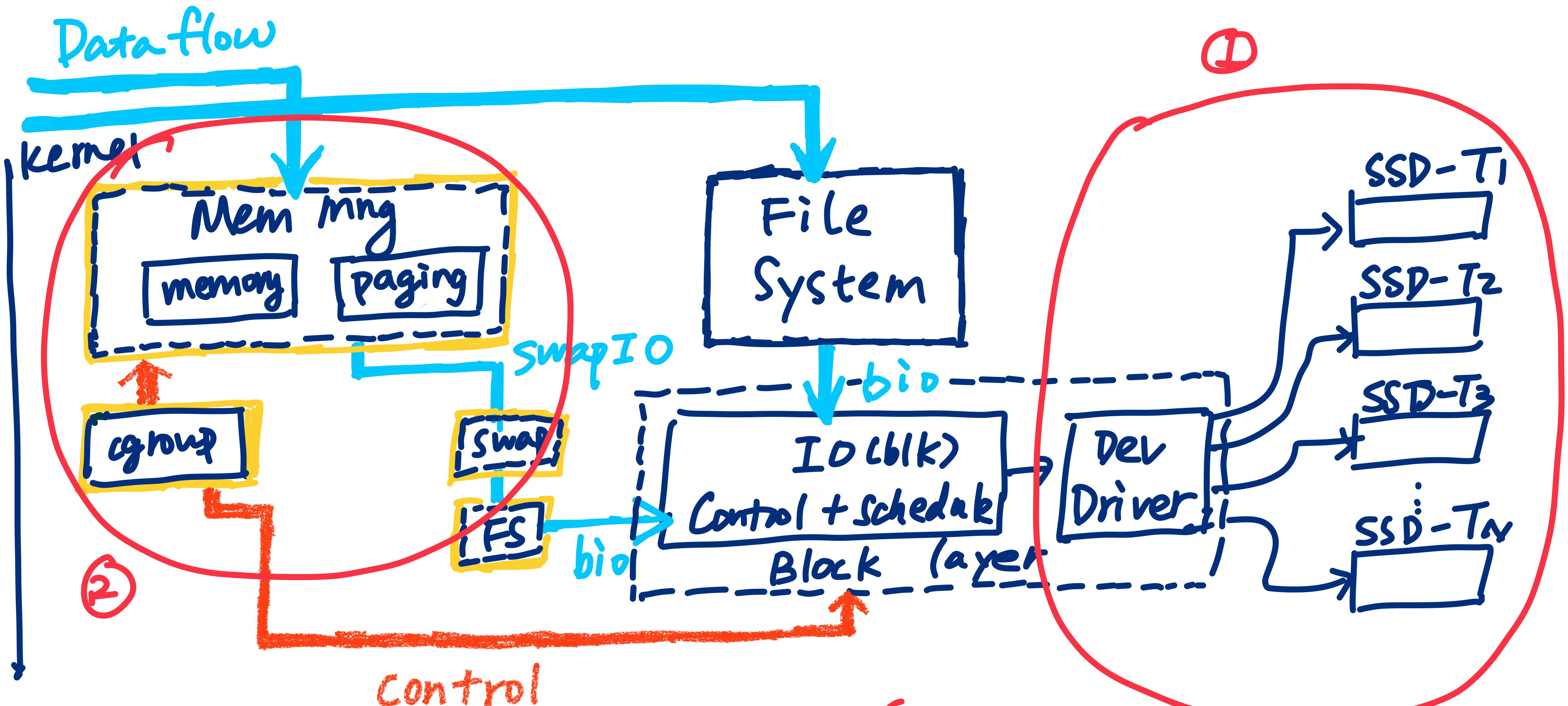
# Data center Setting / Deployment $\oplus$ Challenges



# Data center Setting / Deployment $\oplus$ Challenges

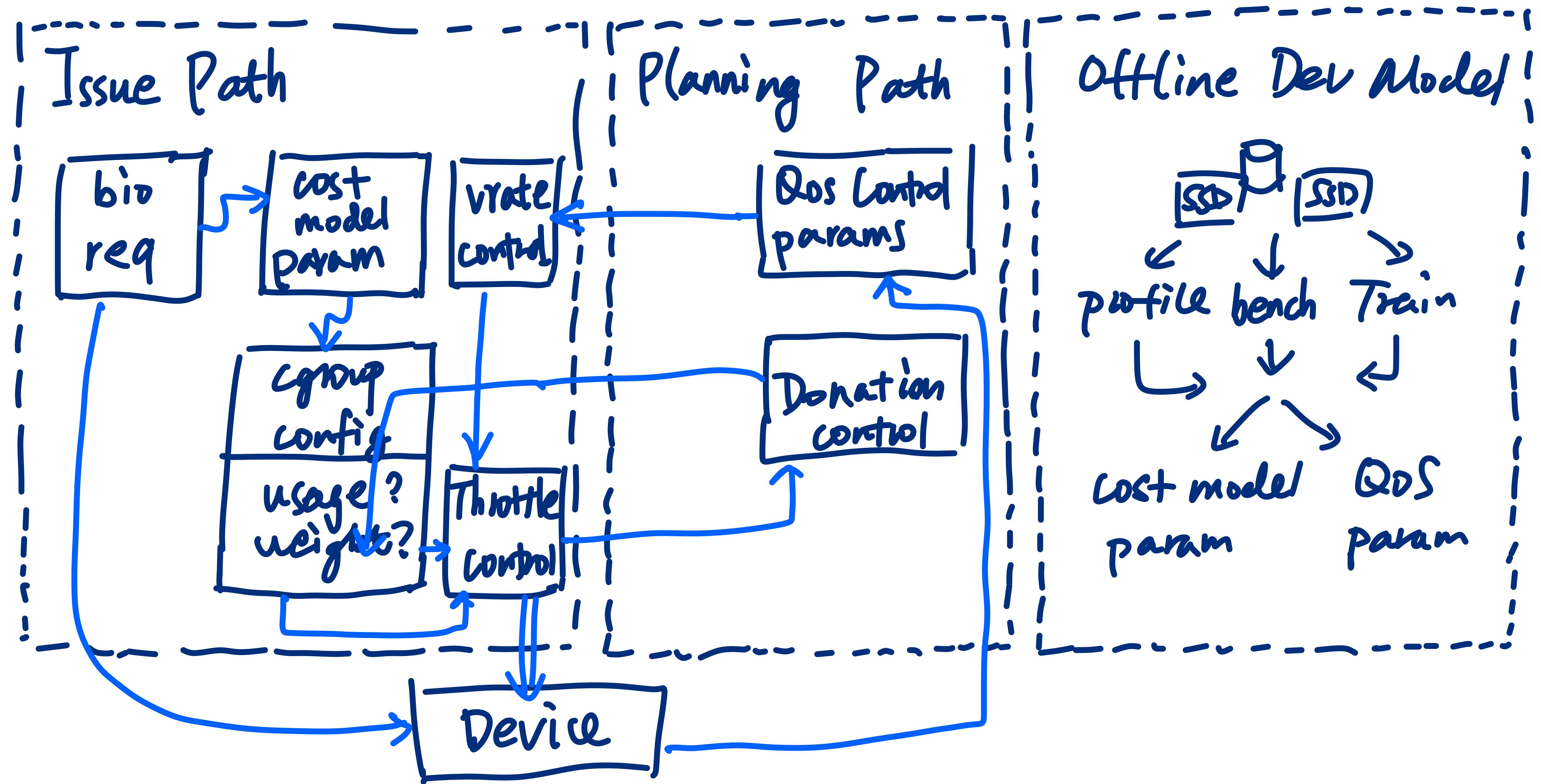


# Data center Setting / Deployment $\oplus$ Challenges



② Interaction with memory subsystem (mem + I/o control make coherent decision)

# IO Cost Overview



## Summary of IO Cost

- ⇒ Offline performance model of device (heavyweight model)
- ⊕ runtime block IO control under dynamic workload (live model) ← ⊕
- { fast path : io operation / policy enforcement
- slow path : policy decision

Summary: TMO + IO cost

- ⇒ Mechanism support for more control in kernel  
+ cross component coordination
- ⇒ userspace process for policy level decision
- ⇒ Modeling the device from different aspects  
+ online approximate + query
- ⇒ more direct app-performance metrics after certain resource control:  
 $\Delta \text{res} \Rightarrow \Delta \text{perf}$   
|-----|  
levels, a lot of components, shared infra  
e.g. sidecar container