



# Kata Containers: the speed of containers, security of VMs - even in a nested environment!

Eric Ernst, Intel  
K. Y Srinivasan, Microsoft  
Shiny Sebastian, Intel

# Agenda

- Overview of Kata
- Nested use case
- KVM on Hyper-V
- A look at Kata nested

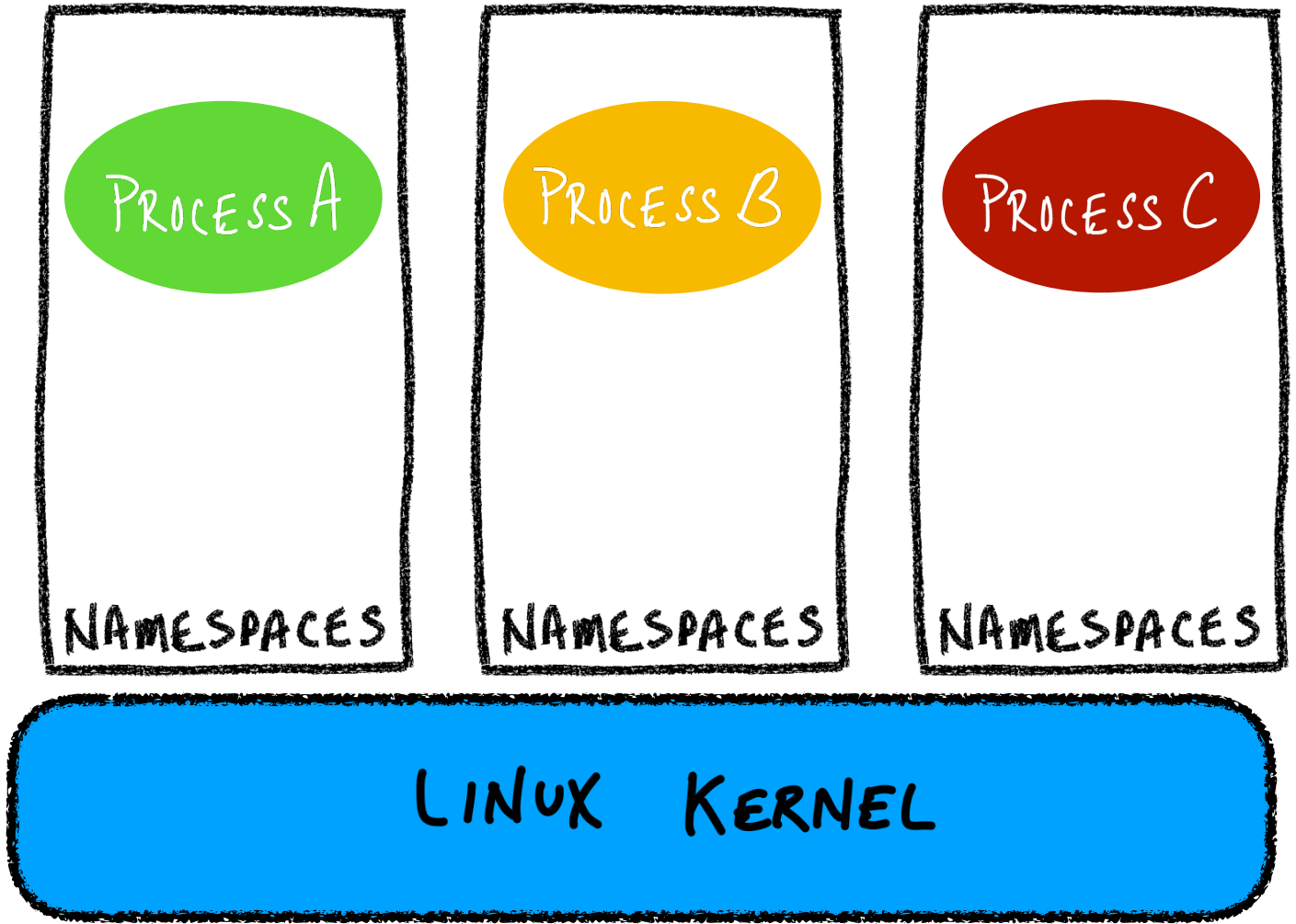
PROCESS A

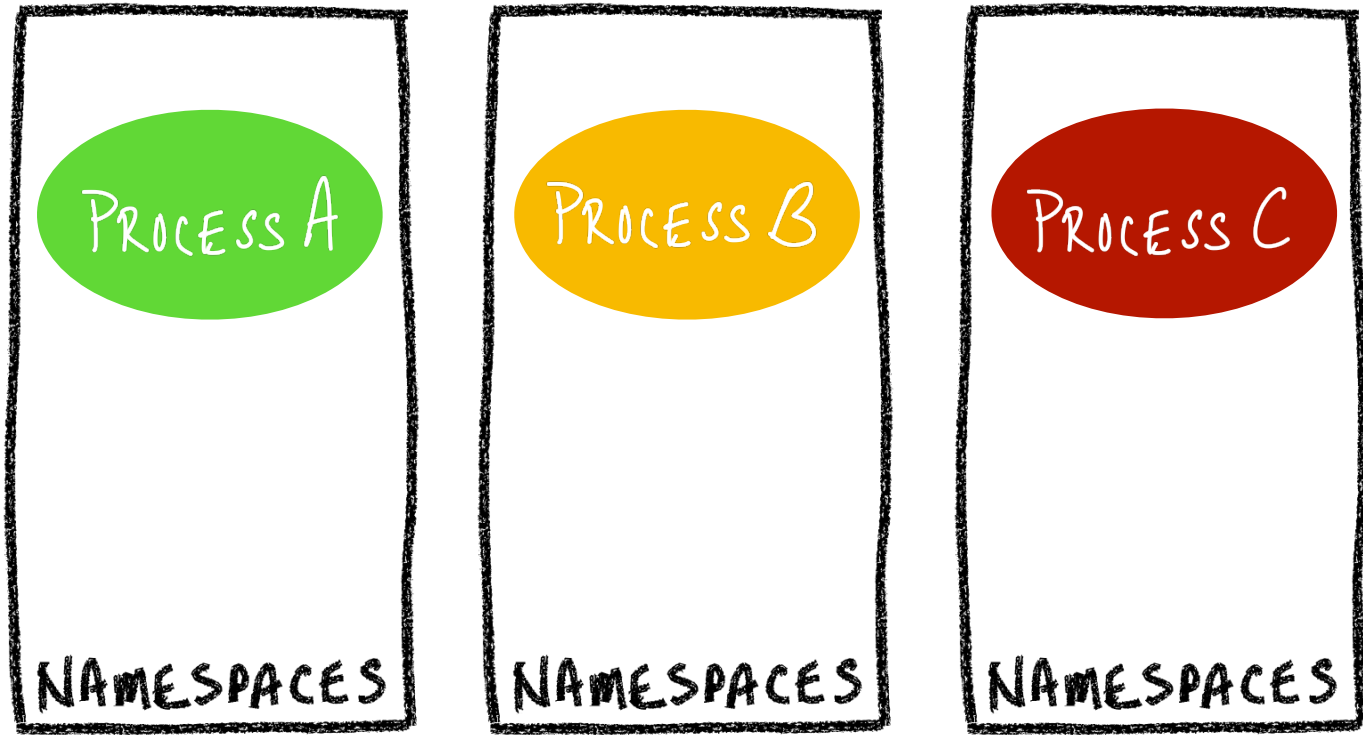
PROCESS B

PROCESS C

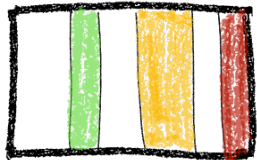
LINUX KERNEL



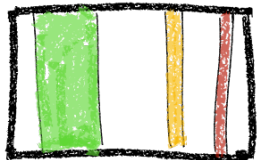




LINUX KERNEL



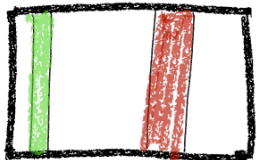
CPU



MEMORY

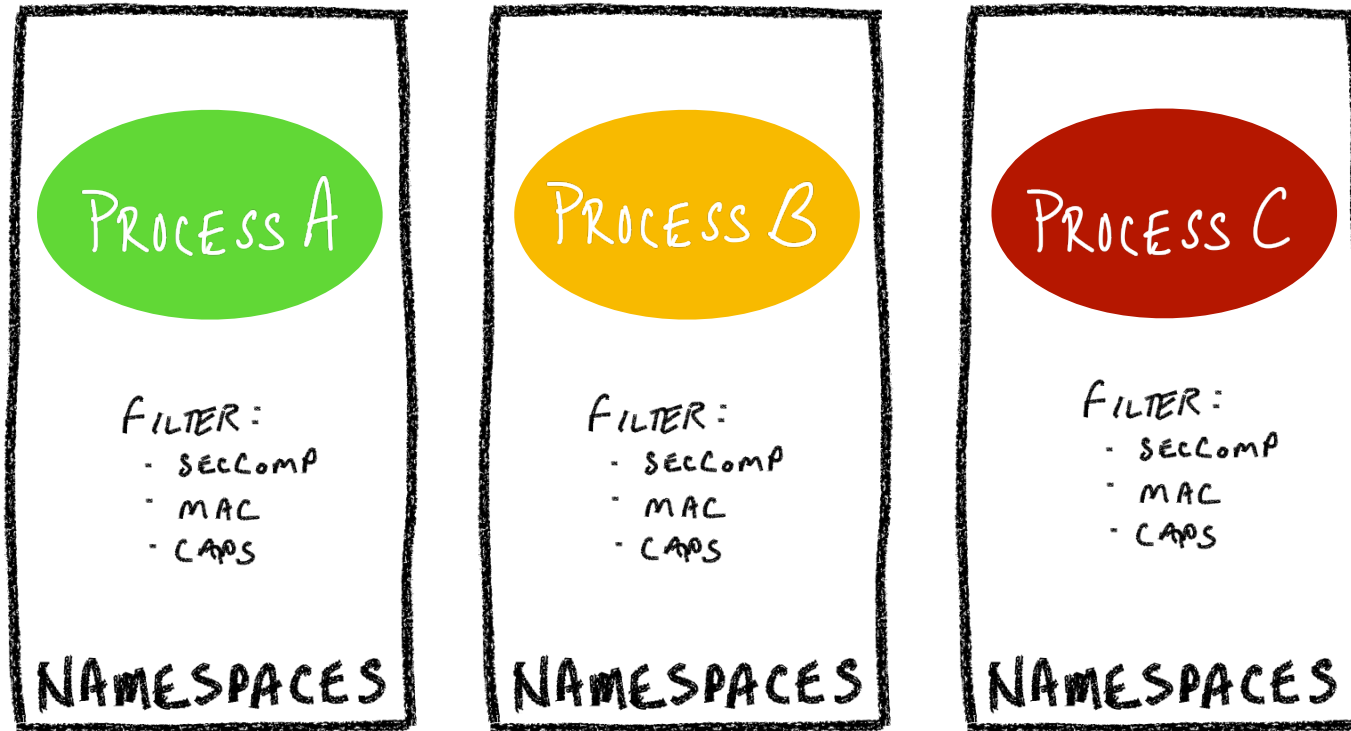


NETWORK



STORAGE





PROCESS A

- FILTER:
- SECCOMP
  - MAC
  - CAPS

NAMESPACES

PROCESS B

- FILTER:
- SECCOMP
  - MAC
  - CAPS

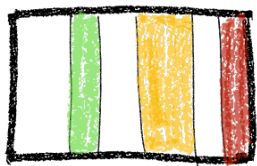
NAMESPACES

PROCESS C

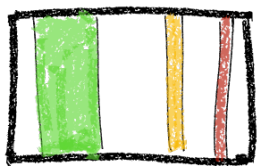
- FILTER:
- SECCOMP
  - MAC
  - CAPS

NAMESPACES

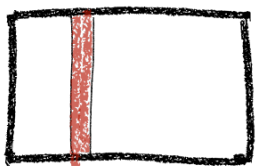
LINUX KERNEL



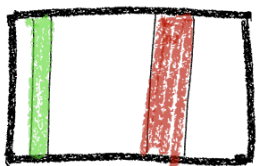
CPU



MEMORY

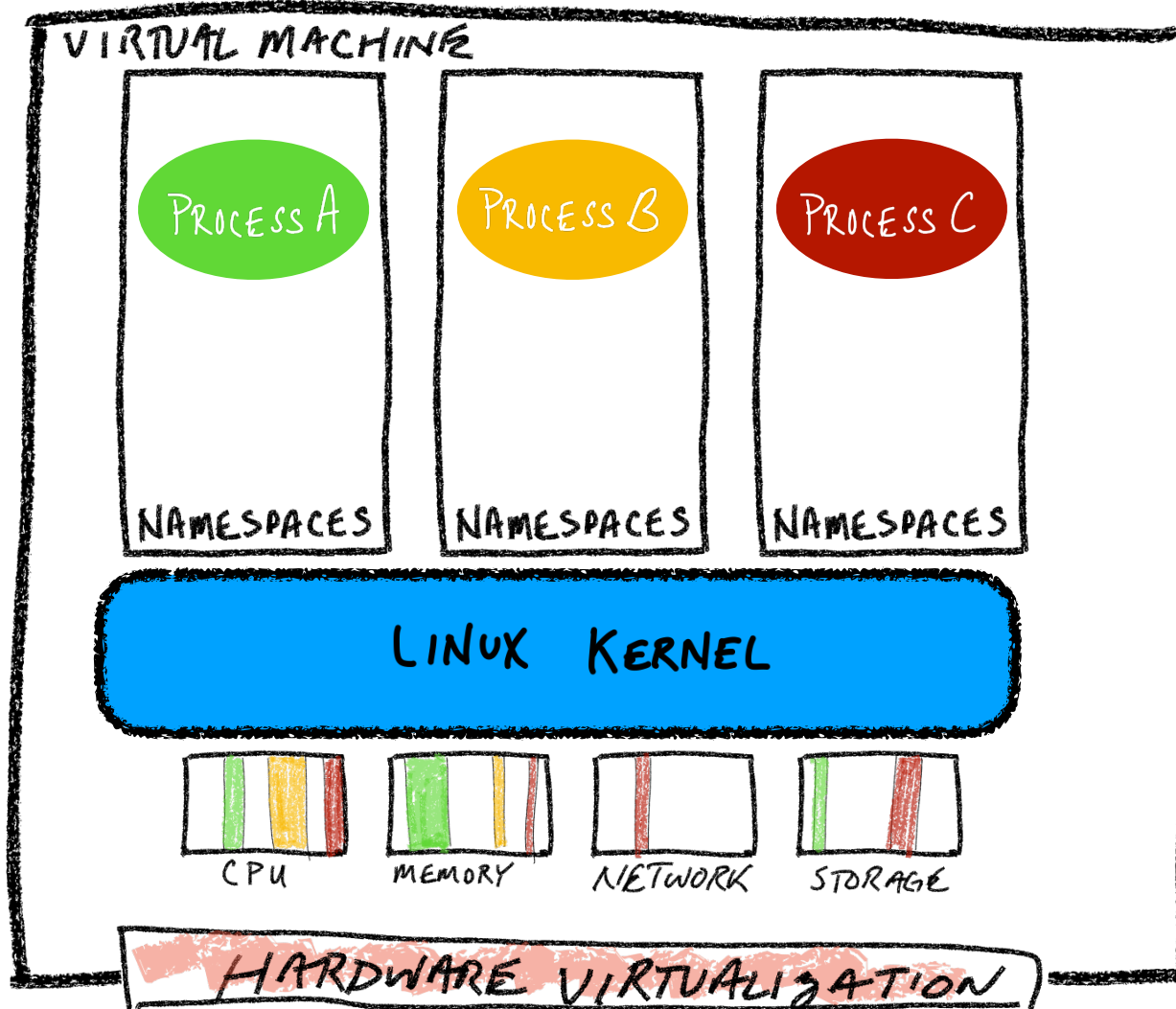


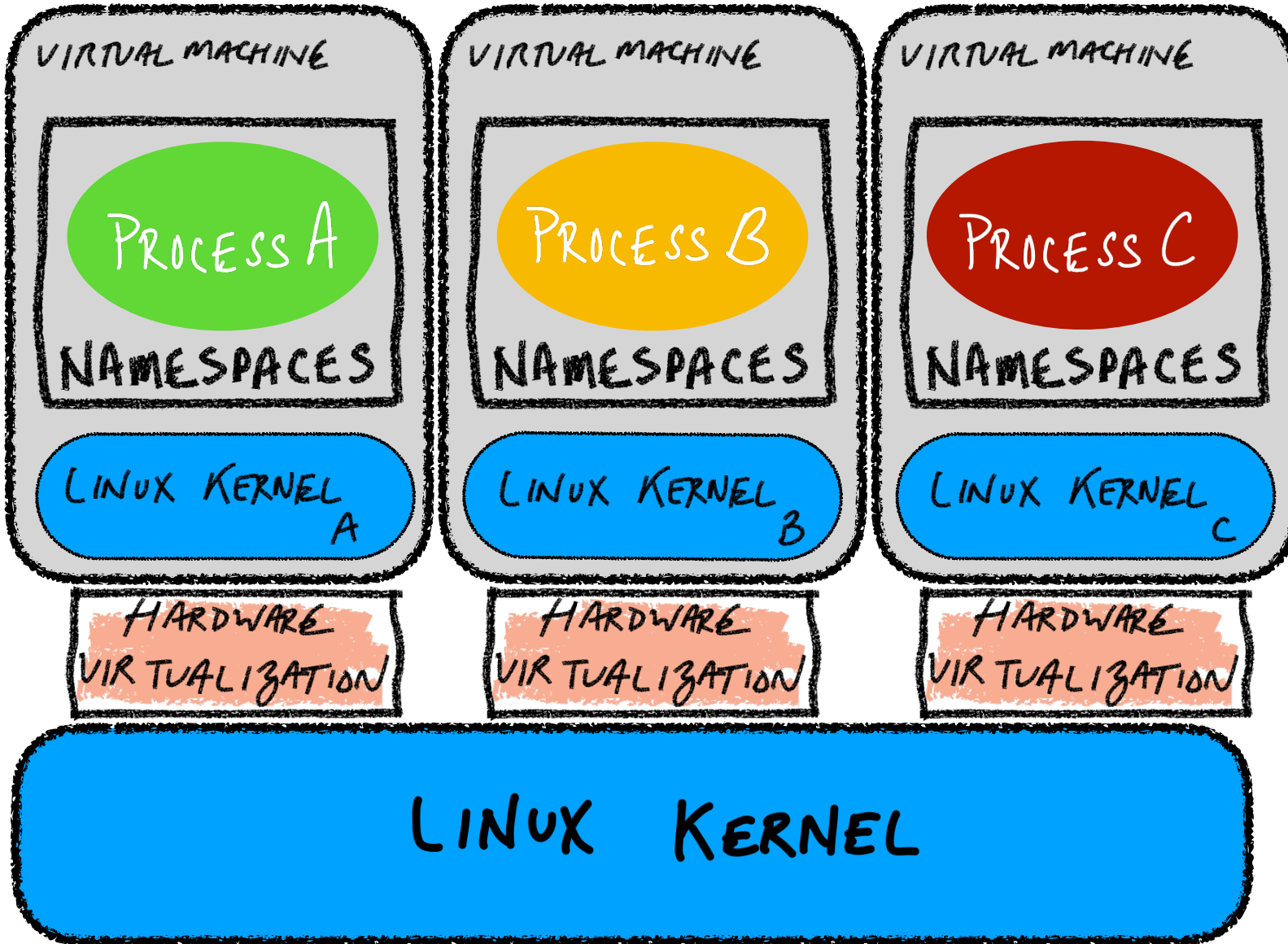
NETWORK



STORAGE





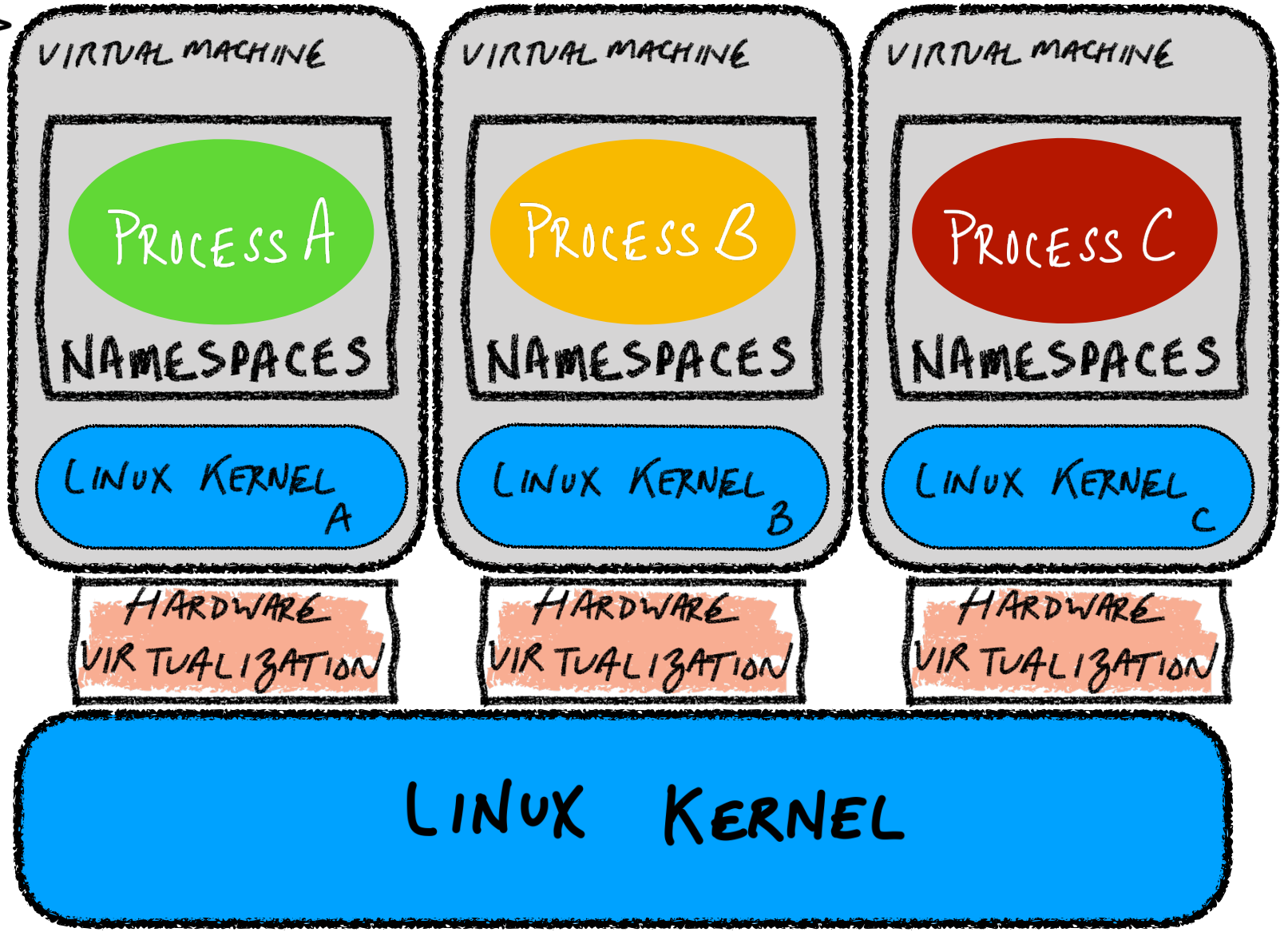


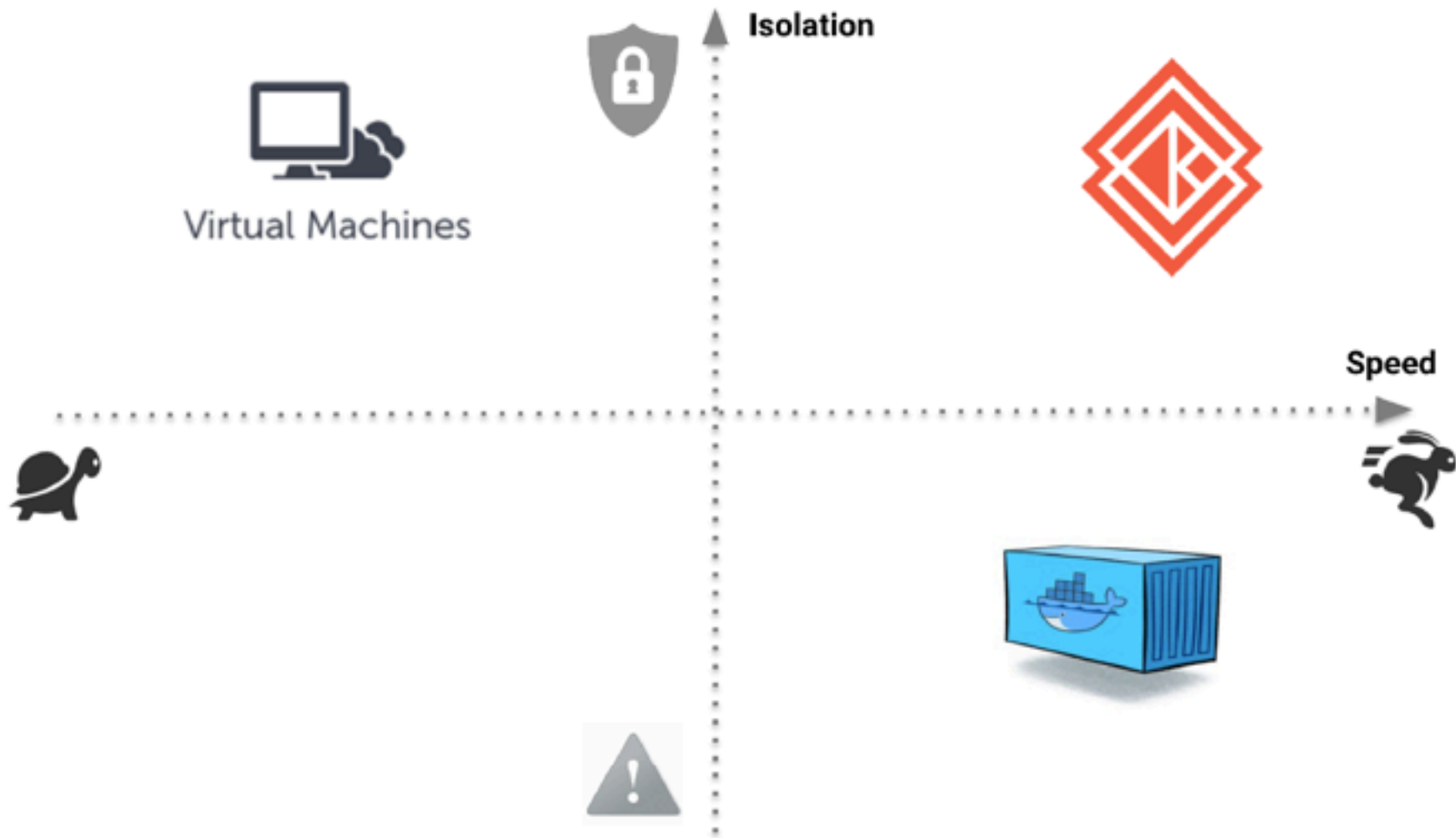


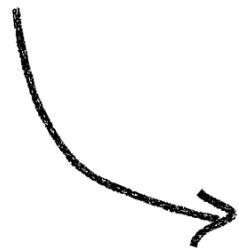
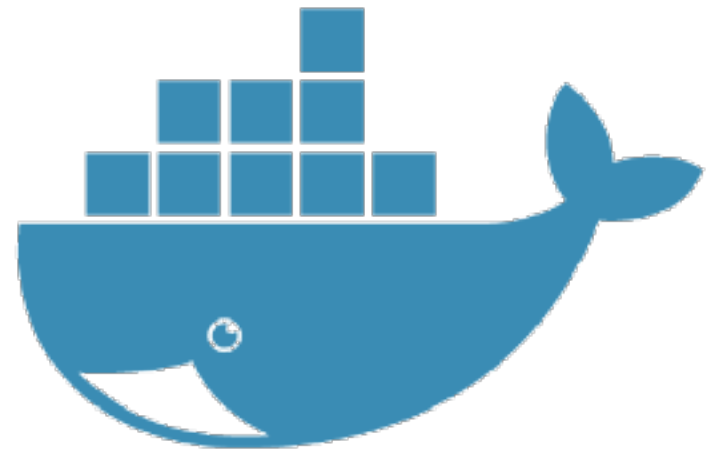
\* VIRTUAL MACHINE CREATED WITH KVM/QEMU

\* ACTIVELY WORKING ON NEMLU  
NO EMULATION

- ✓ New enlightened machine type
- ✓ Based on QEMU
- ✓ See talk @ KVM-Forum on Thursday

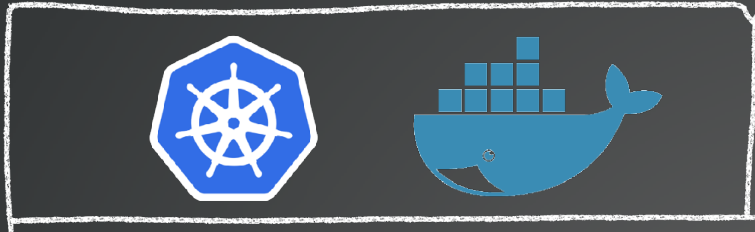






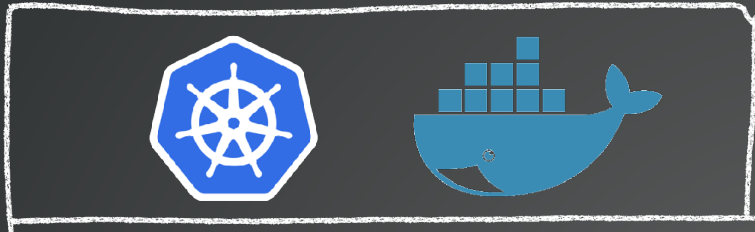
OCI  
COMPLIANT  
RUNTIME





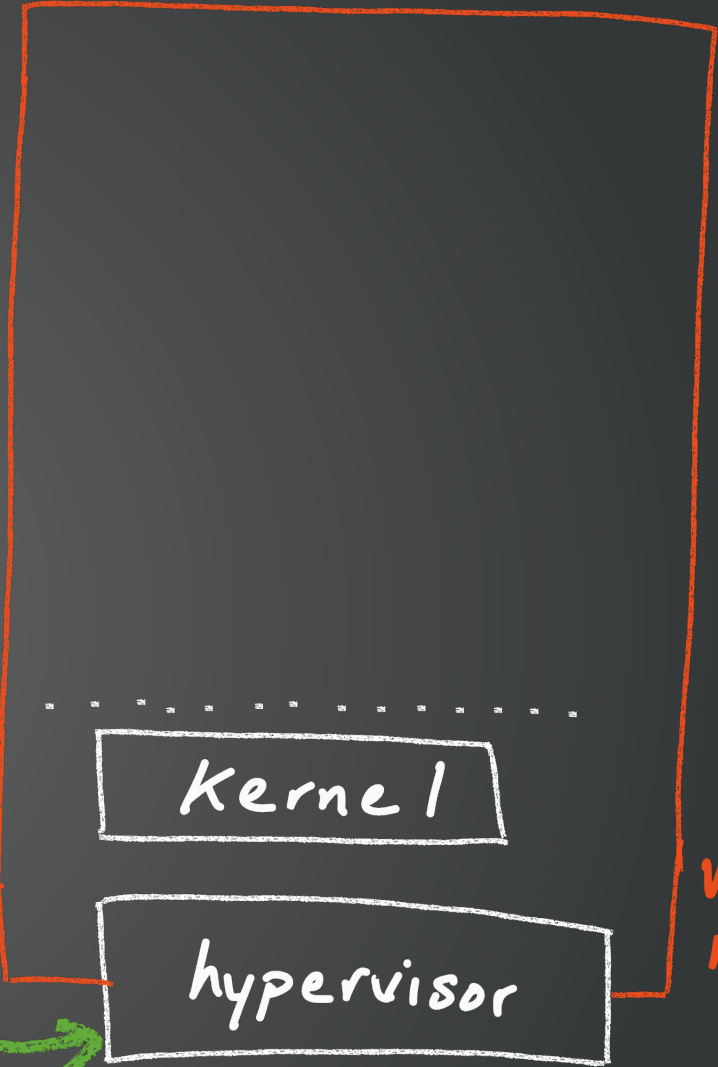
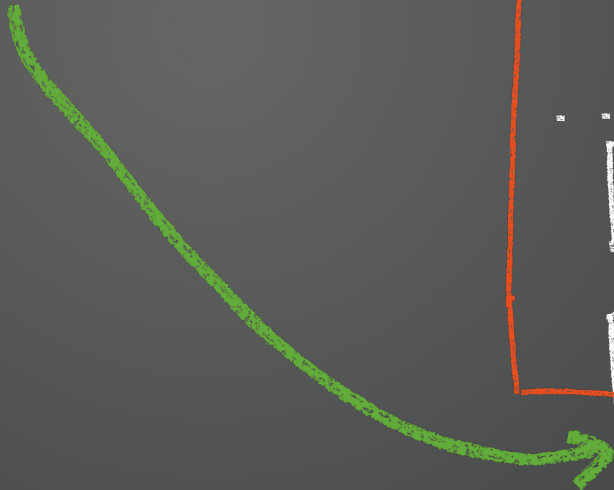
↓ (oci)

Kata-runtime

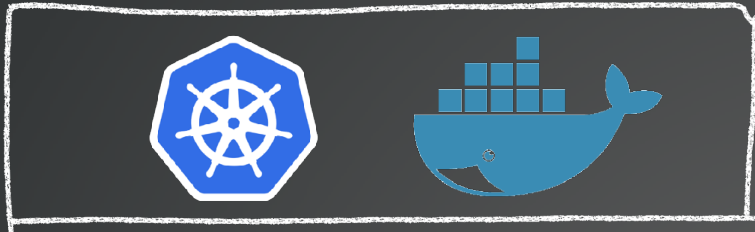


(oci)

Kata-runtime

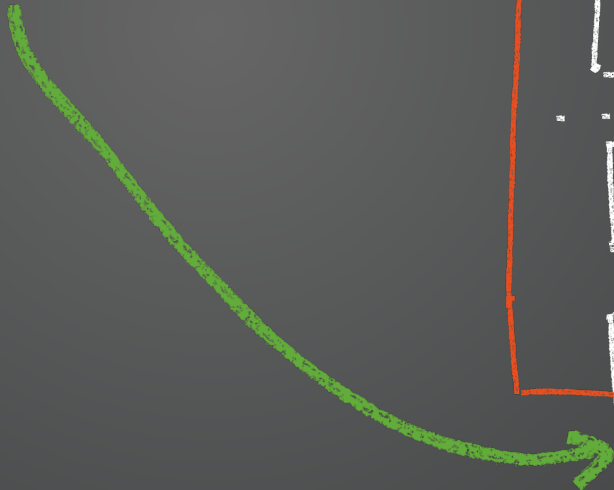


Virtual Machine



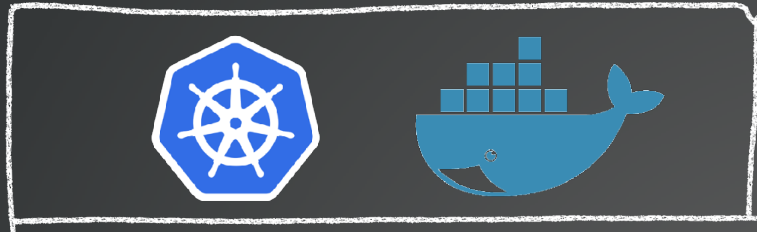
(oci)  
↓

Kata-runtime



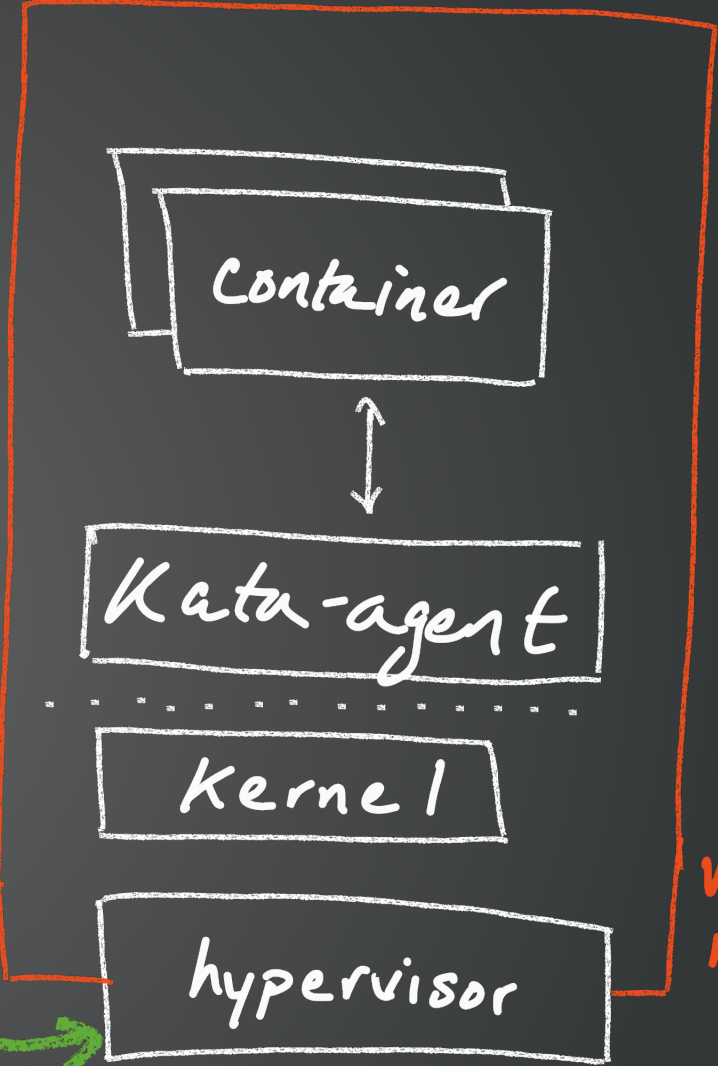
Virtual Machine

vsock

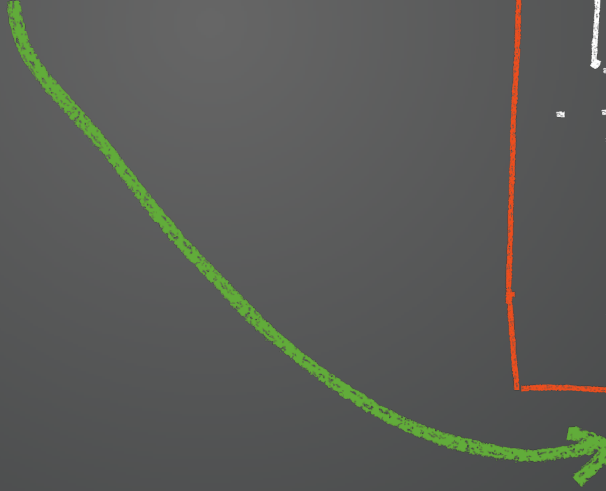


(oci)

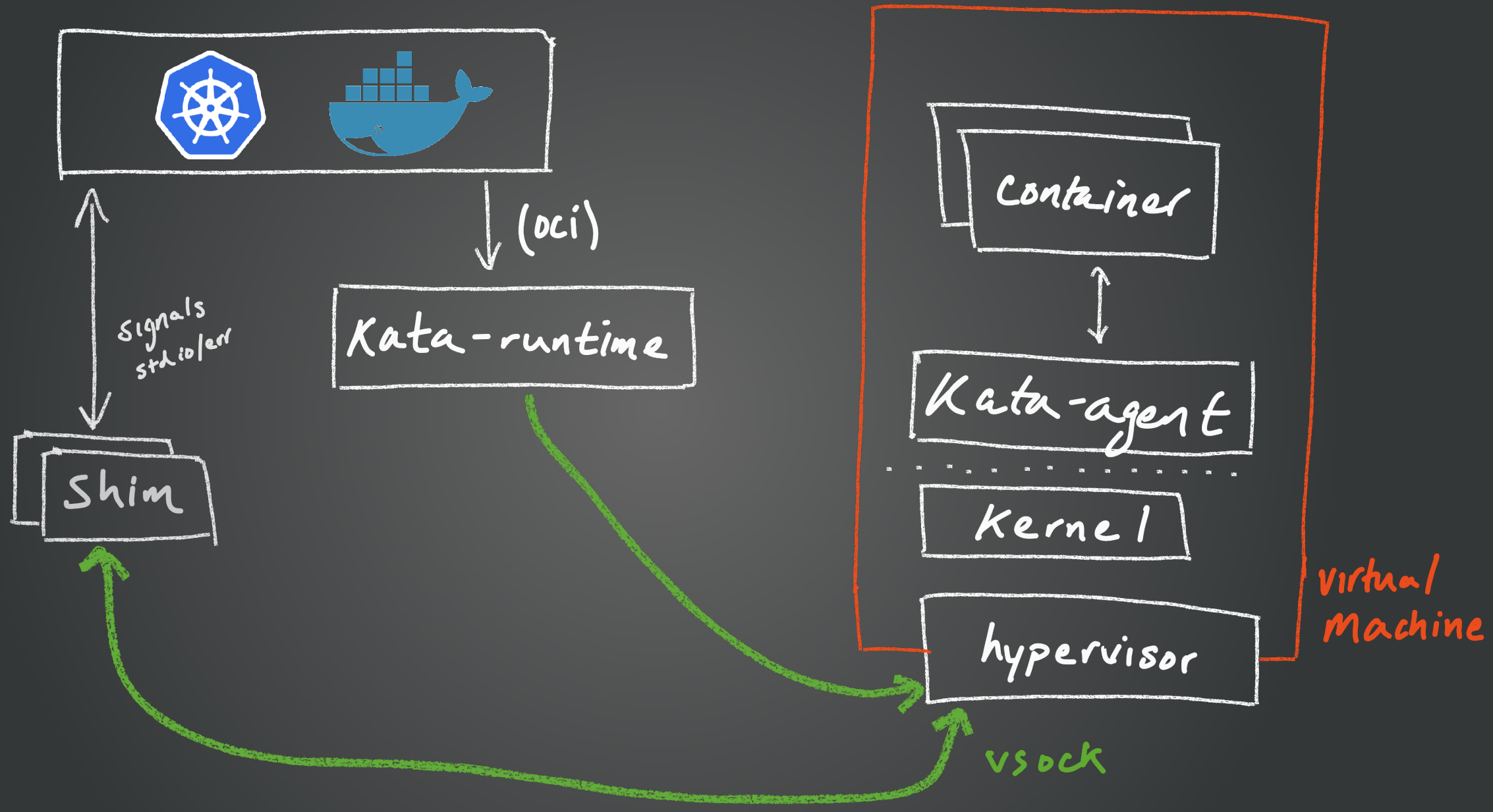
Kata-runtime



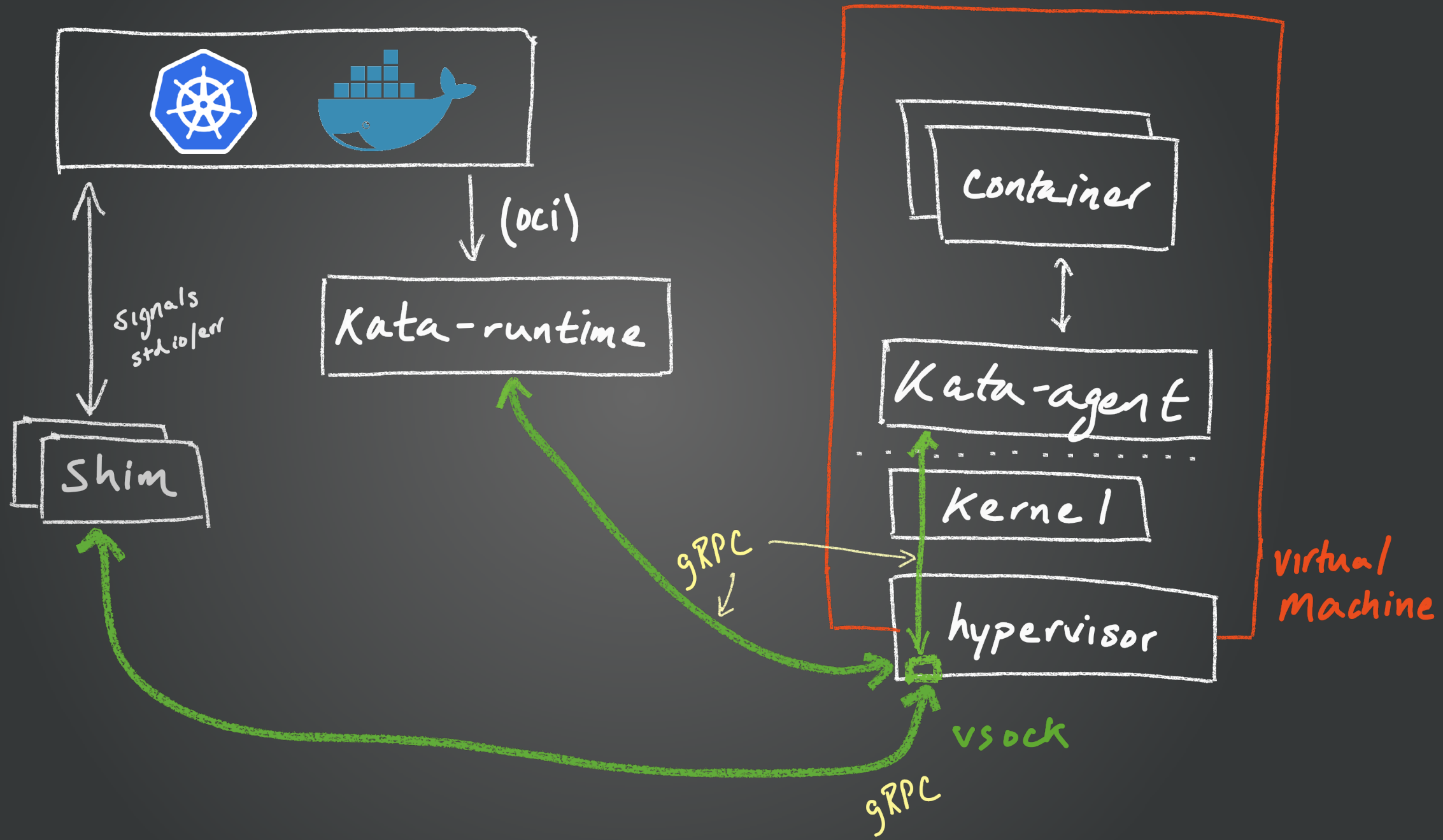
Virtual Machine



vsock





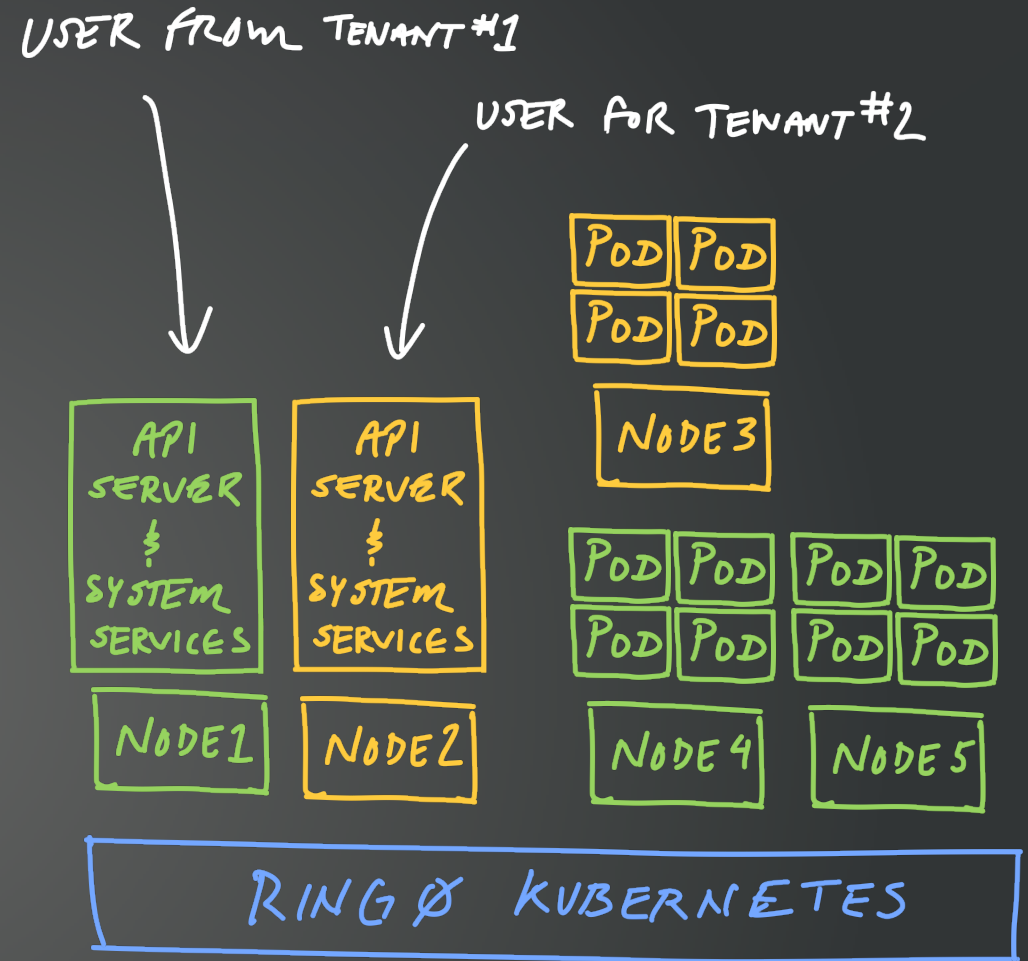


A decorative graphic on the left side of the slide, consisting of several overlapping circles in various shades of green. The background of the slide is a dark grey gradient.

# Nested Use Case

# Nested Kata use case

- How many?
  - What's an appropriate pool size for tenant 1, or a particular workload?
  - Does each untrusted workload need its own VM?
- Infrastructure work:
  - Need to manage virtual machine creation, and tag each per workload/tenant.
  - May need to size virtual machines conservatively based on what a workload \*could\* need
  - SDN, SDS and fabric overheads with spinning up VMs

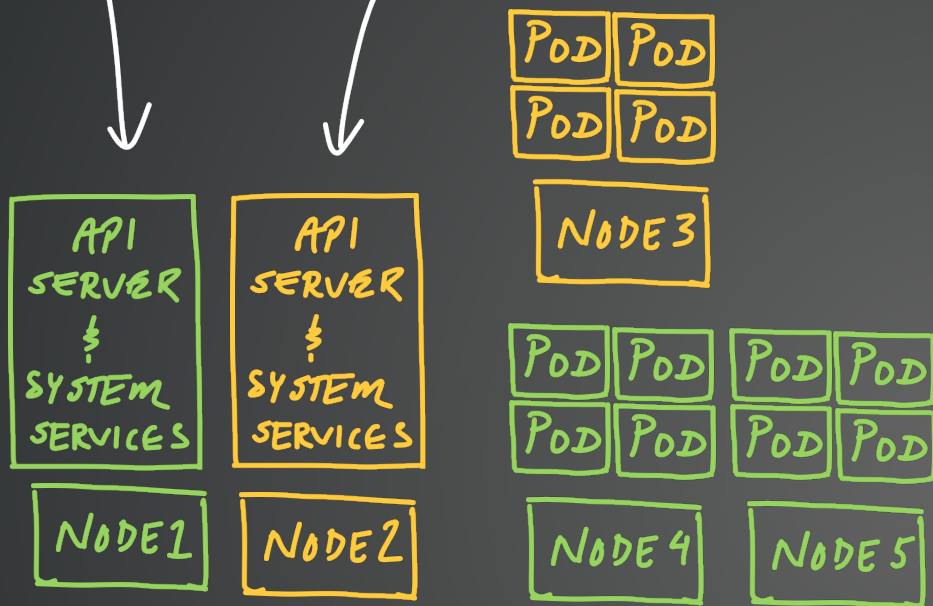


(NAMESPACE 1) (NAMESPACE 2)

# Nested Kata use case

USER FROM TENANT #1

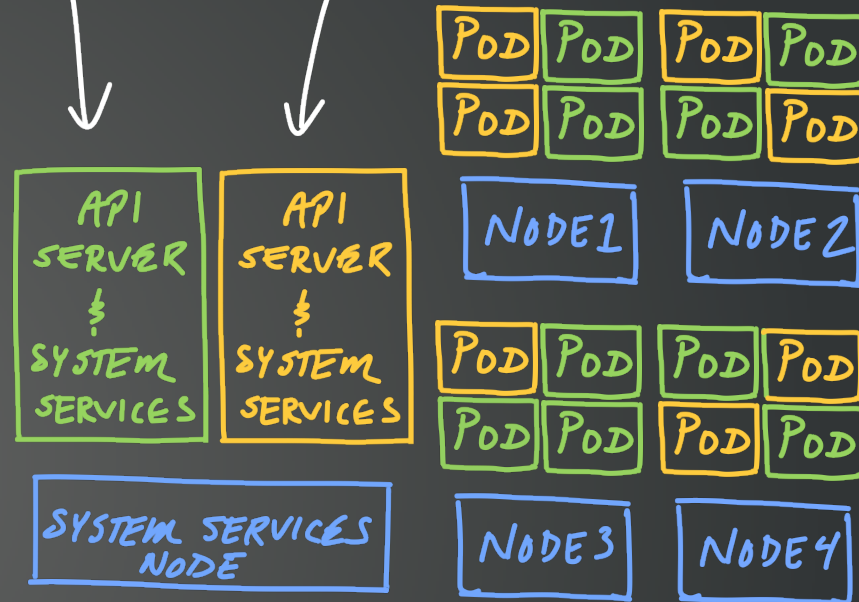
USER FOR TENANT #2



RING OF KUBERNETES

USER FROM TENANT #1

USER FOR TENANT #2



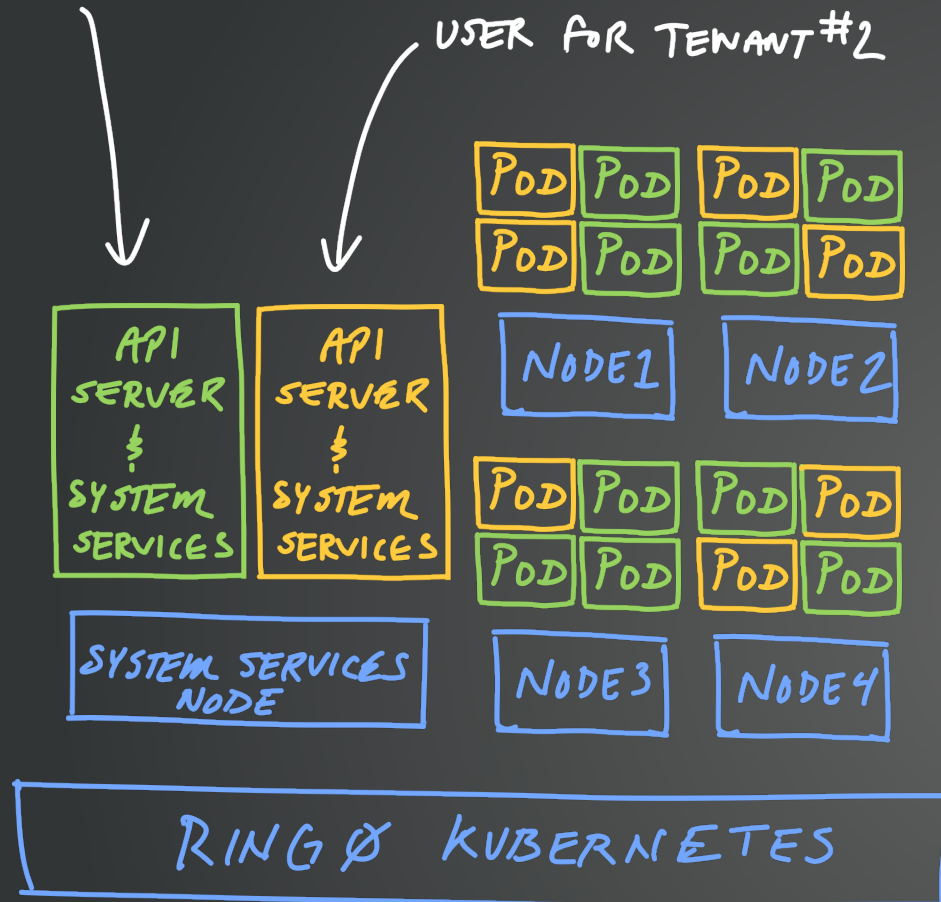
RING OF KUBERNETES

(NAMESPACE 1) (NAMESPACE 2)

# Nested Kata use case

USER FROM TENANT #1

USER FOR TENANT #2



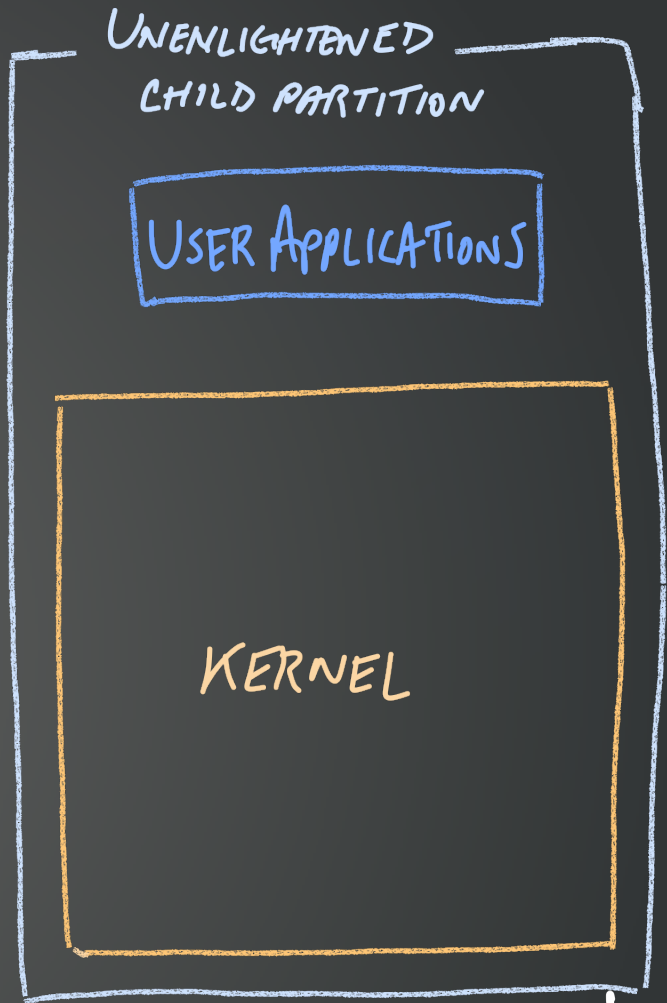
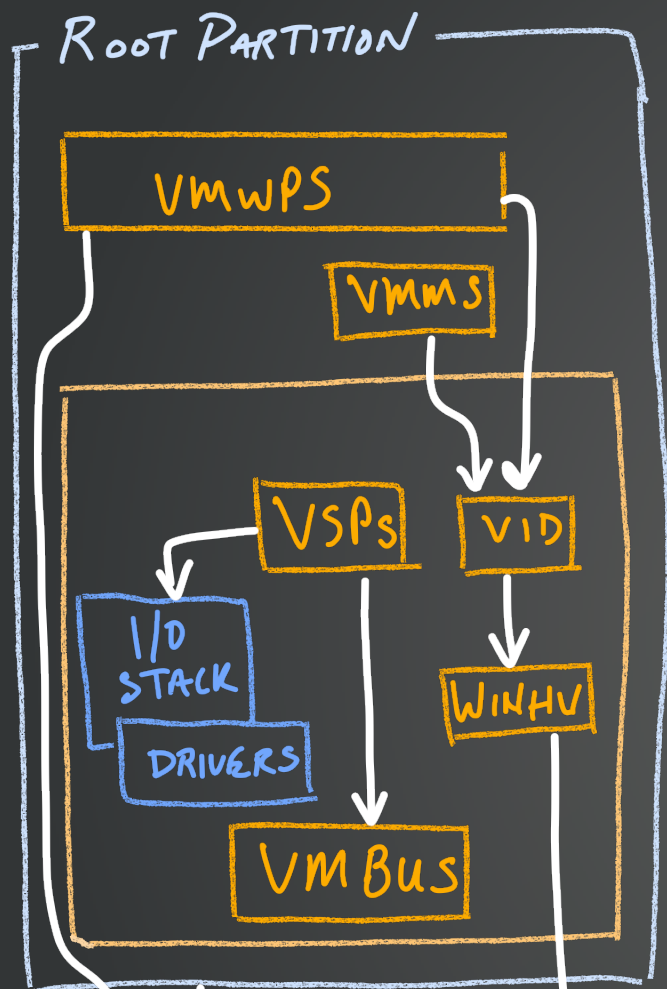
- Better utilization of resources:

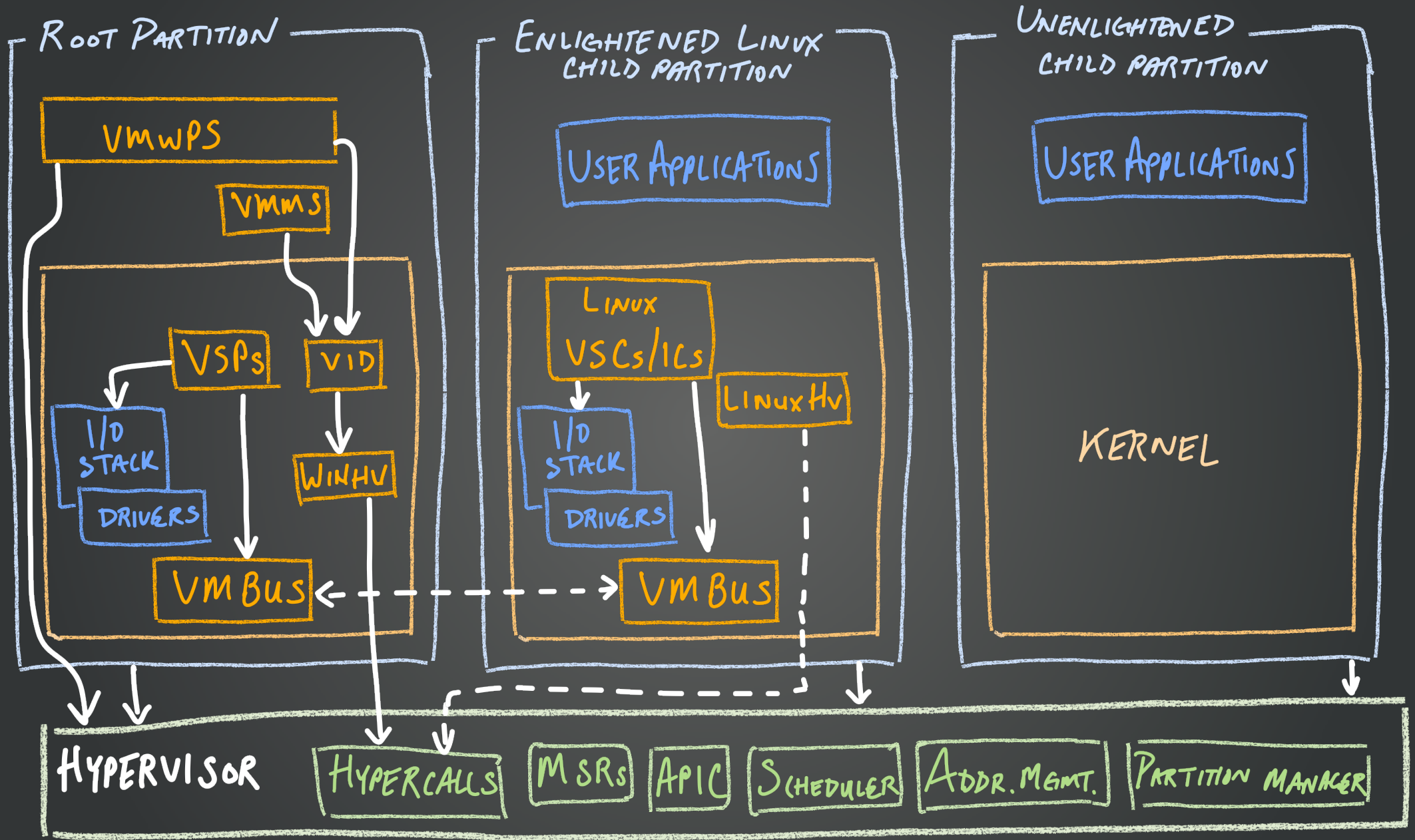
- No need to size VM based on workload needs
- No need to size node pool based on potential tenants potential need
- Pool is shared at a finer granularity

(NAMESPACE 1) (NAMESPACE 2)

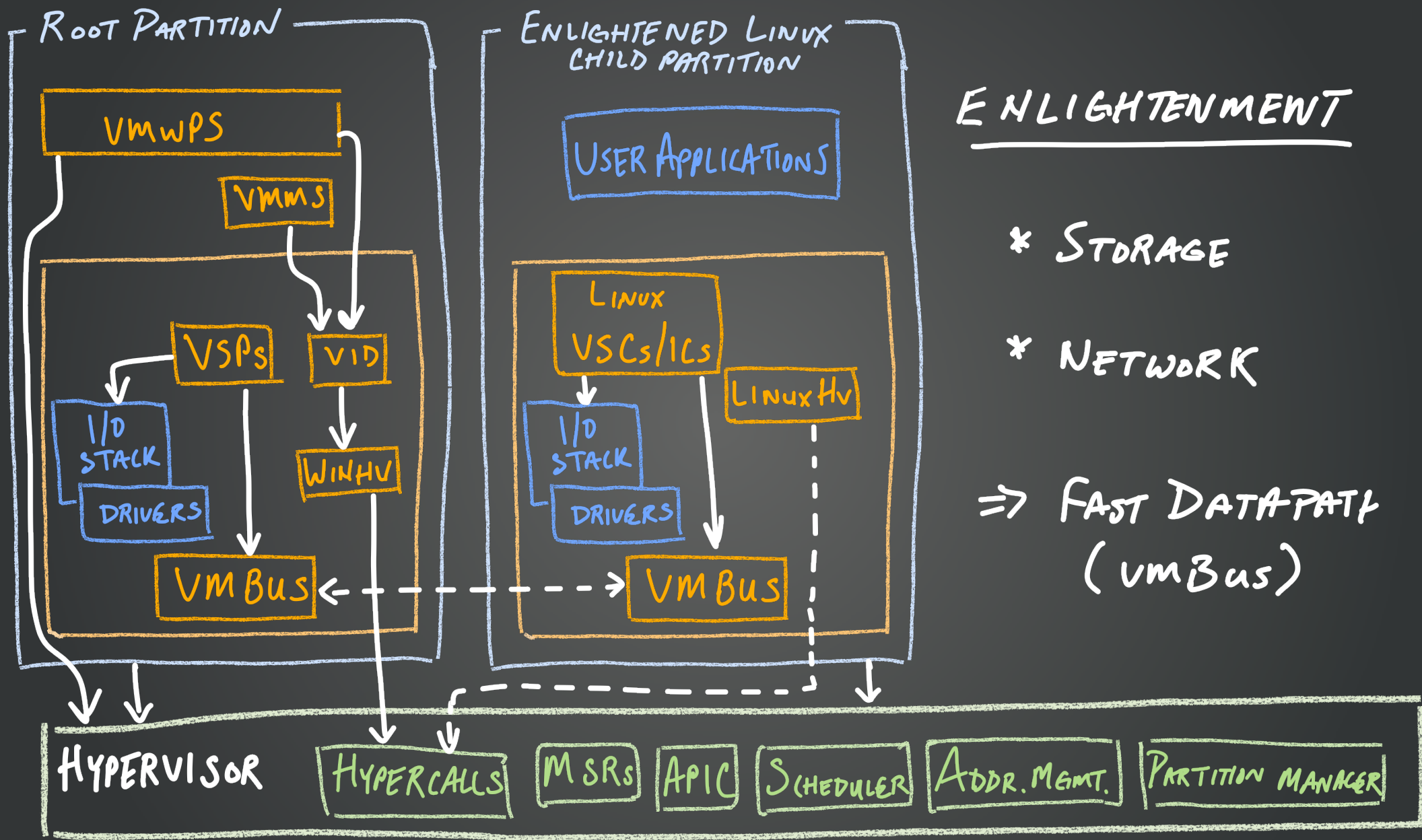
A decorative graphic on the left side of the slide, consisting of several overlapping circles in various shades of green. The background of the slide is a dark grey gradient.

# **KVM on Hyper-V**





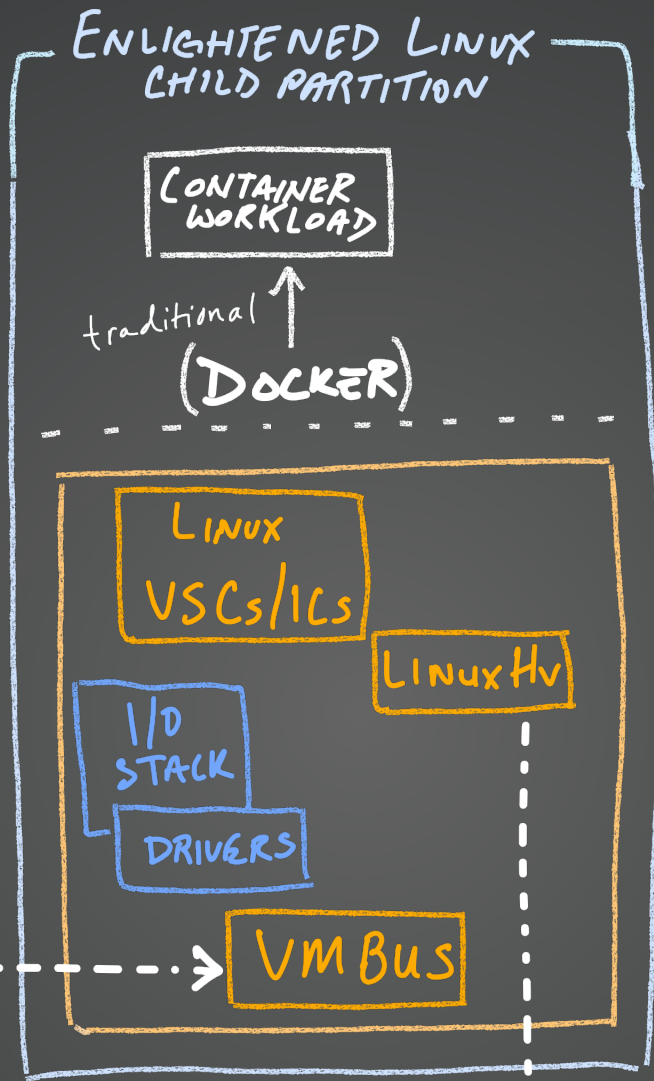
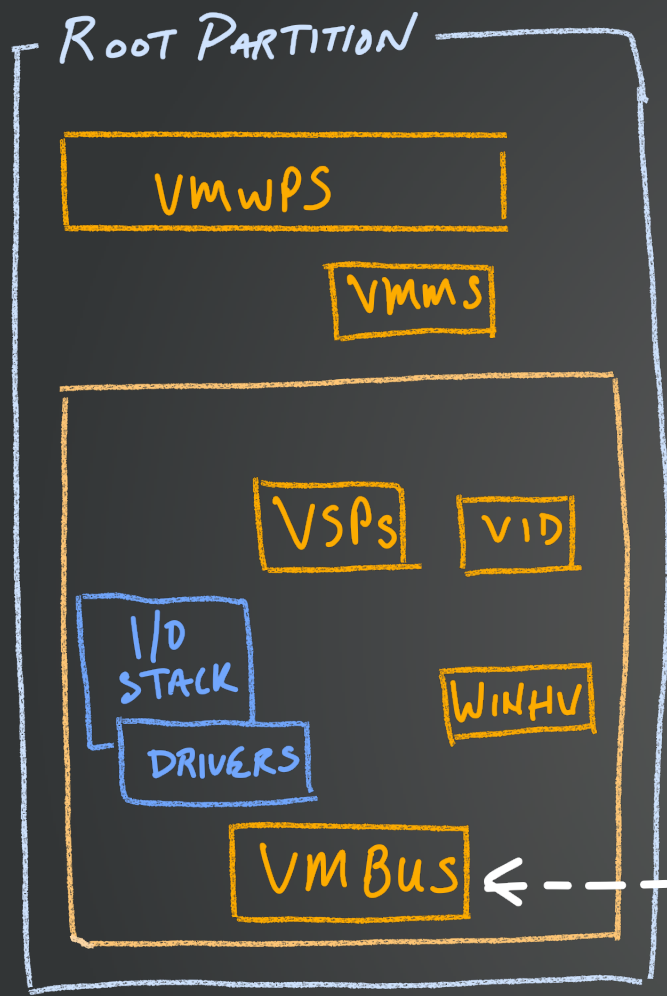


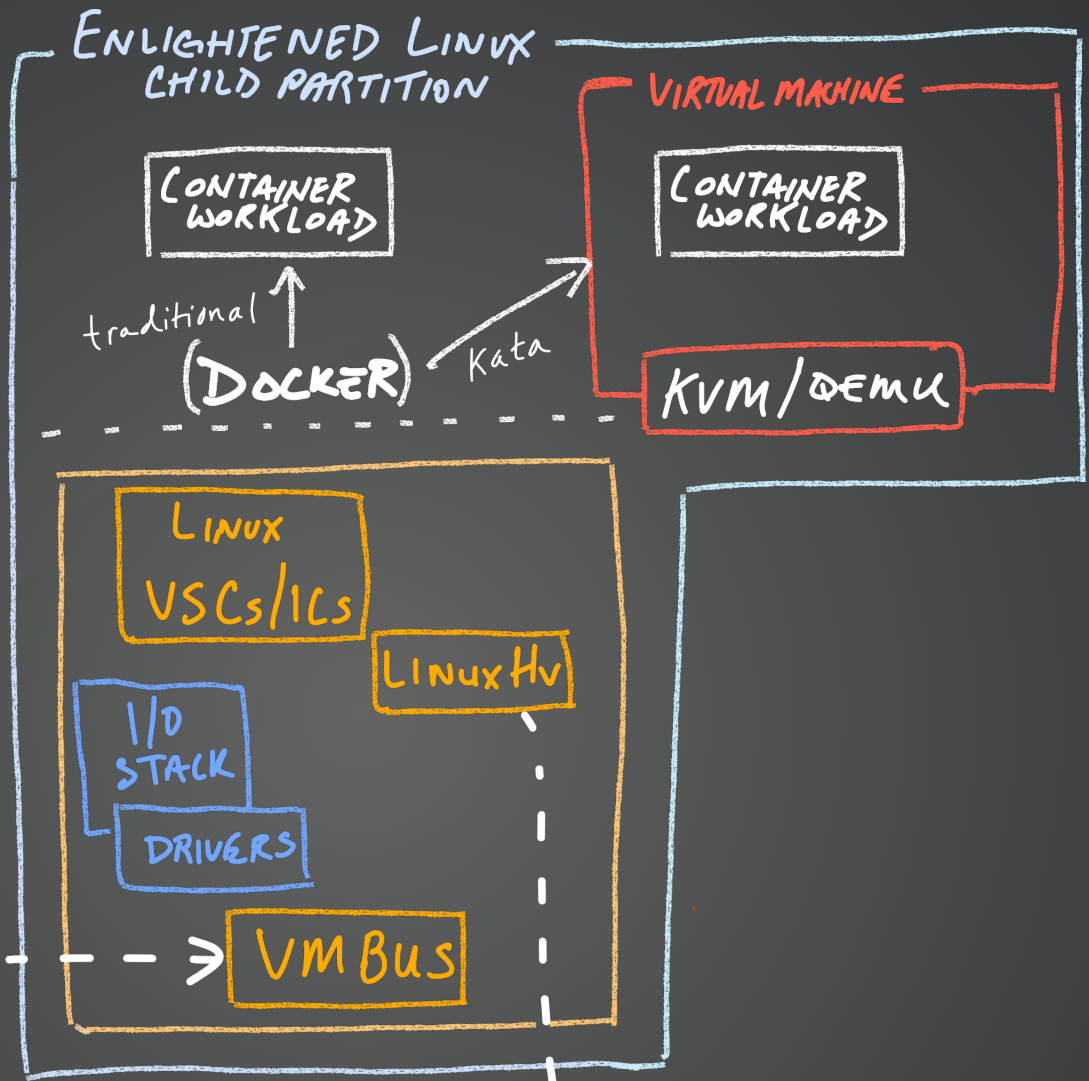
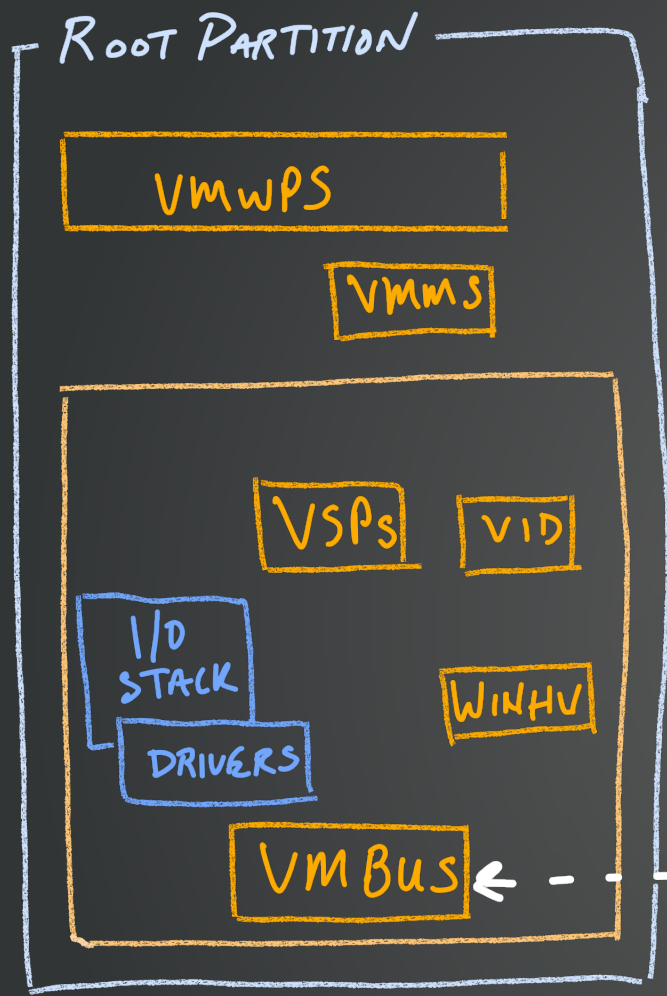


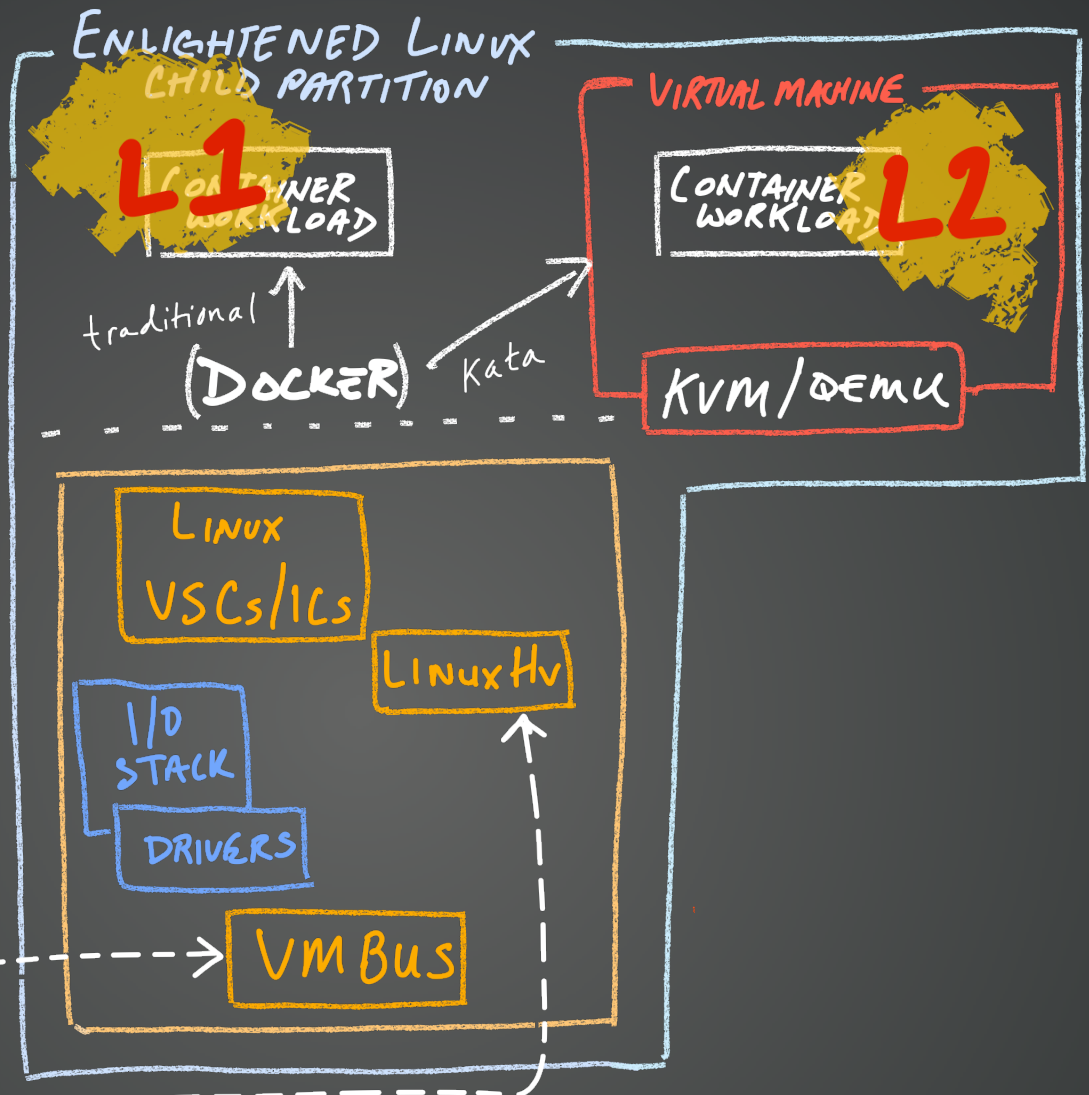
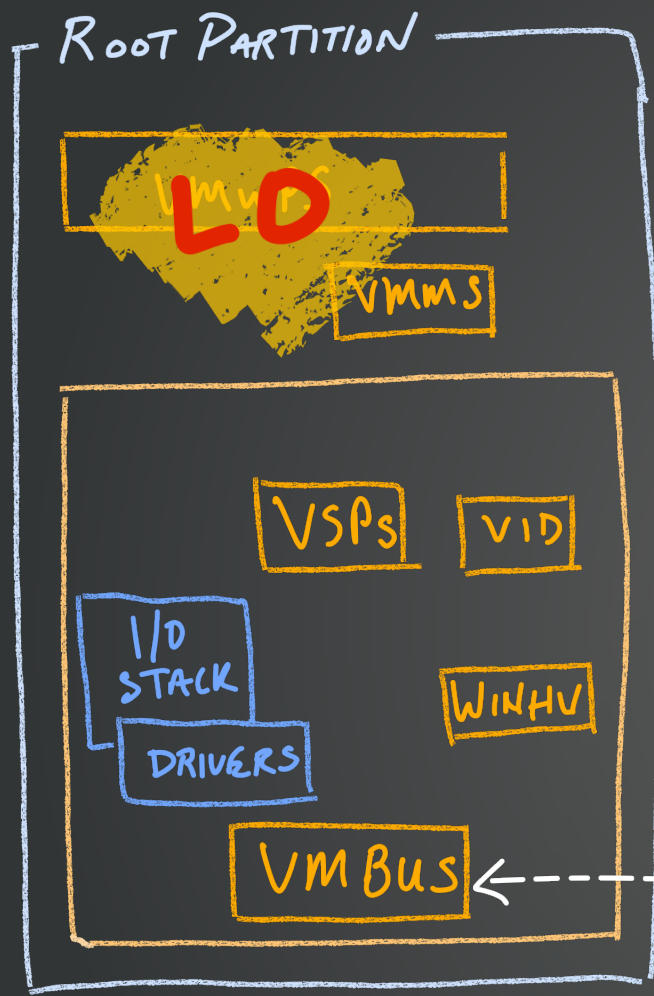
## ENLIGHTENMENT

- \* STORAGE
- \* NETWORK

⇒ FAST DATAPATH (vmbus)



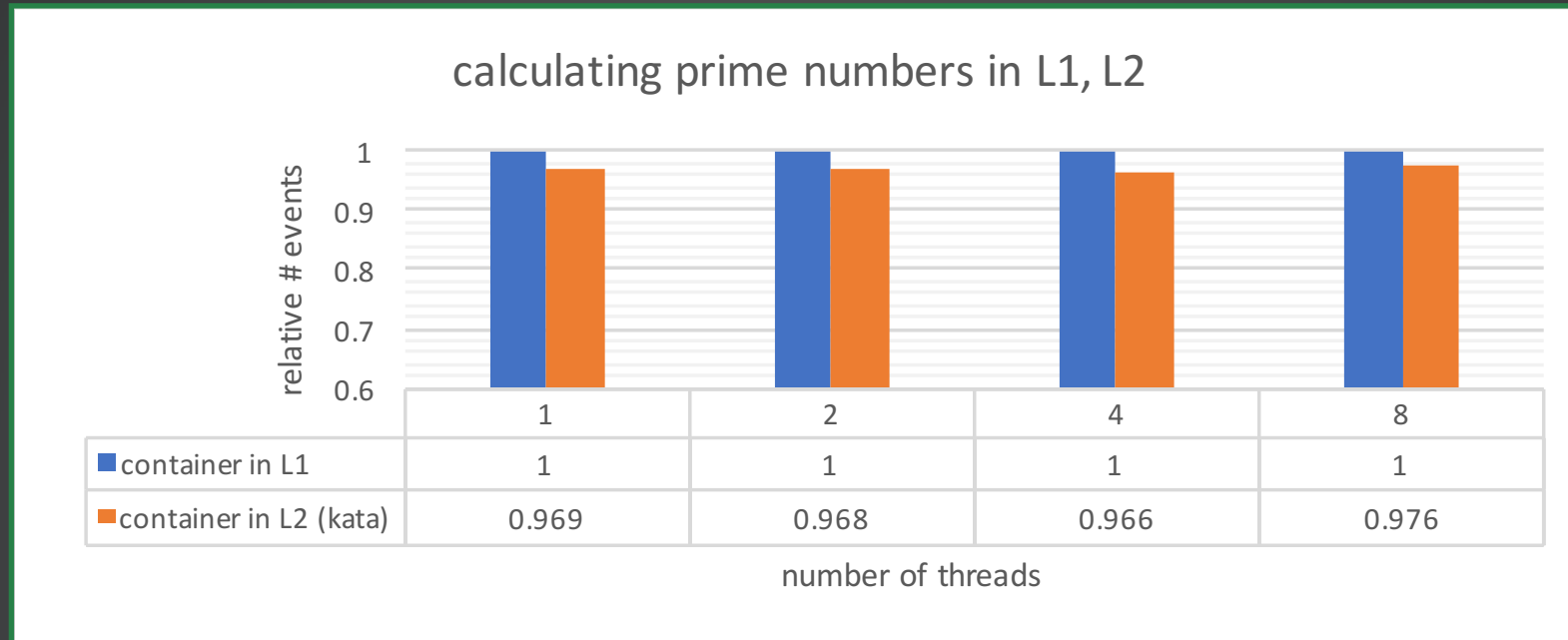






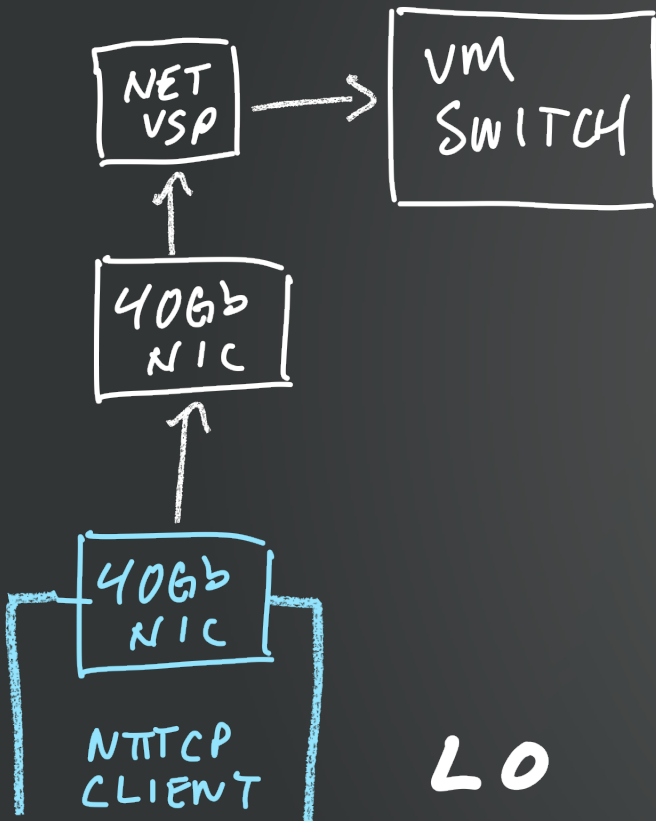
# Nested Kata

# Nested Kata: CPU Measurements

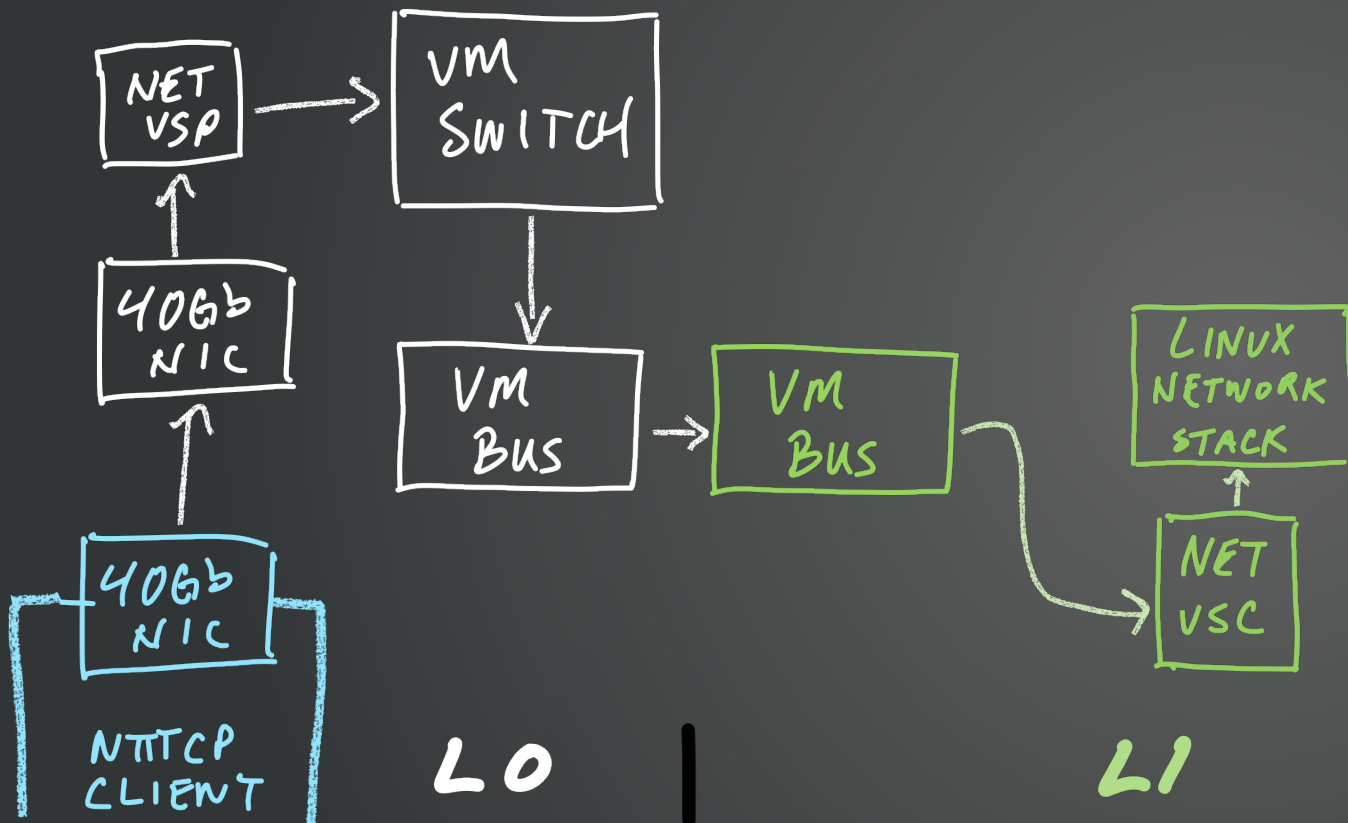


Approximately 3% degradation seen when running with varying Number of threads on prime number calculation workload

# Nested Kata: Network I/O - setup

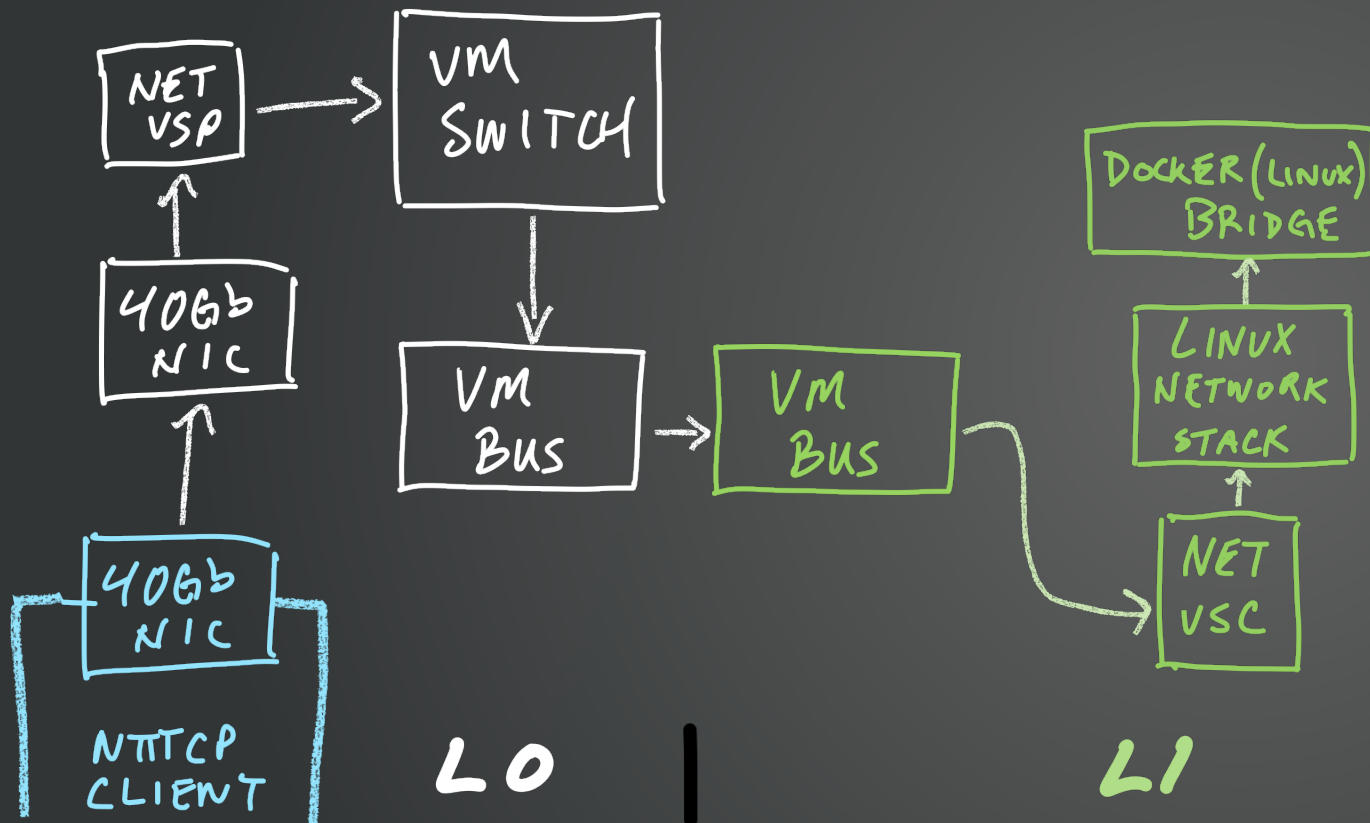


# Nested Kata: Network I/O - setup

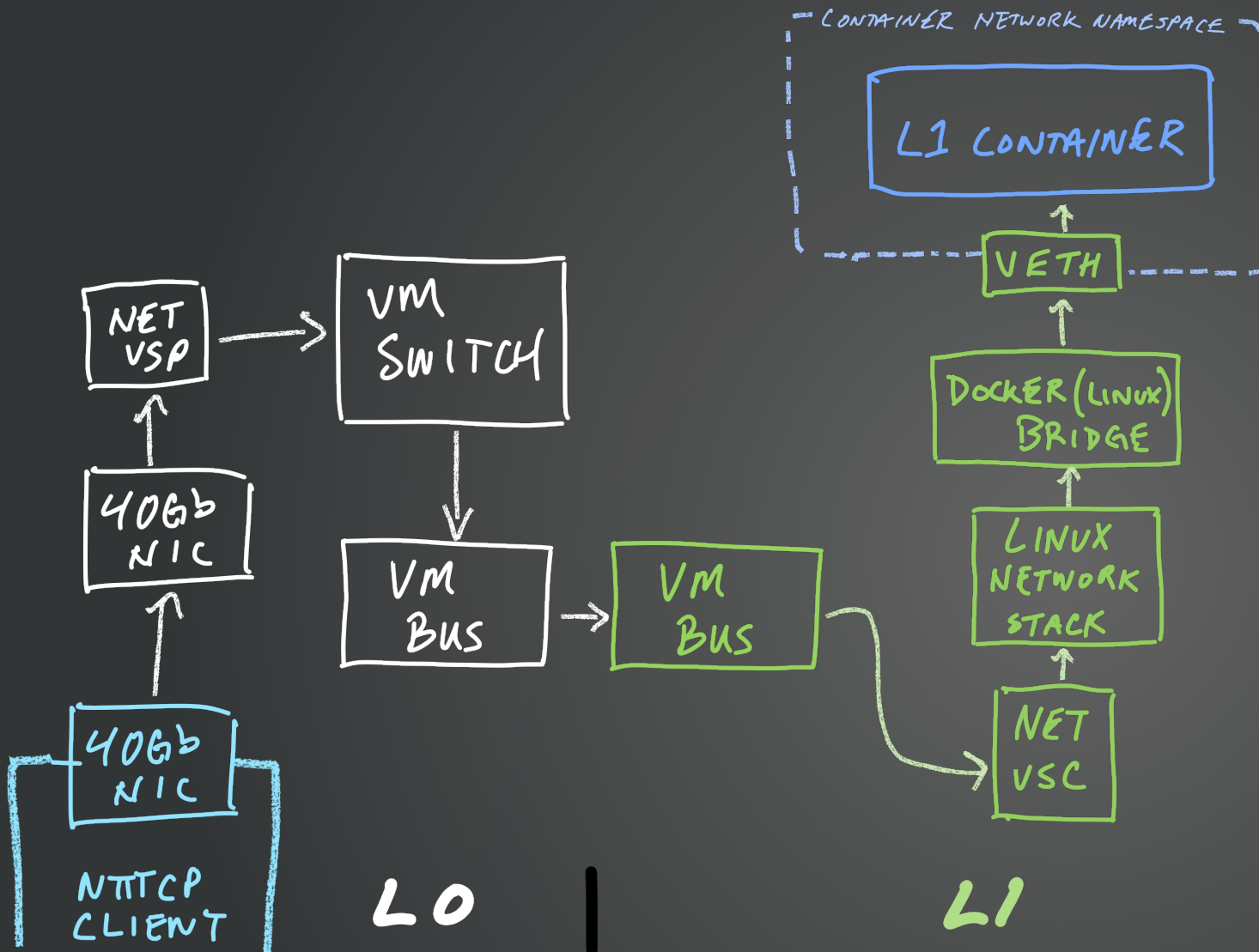




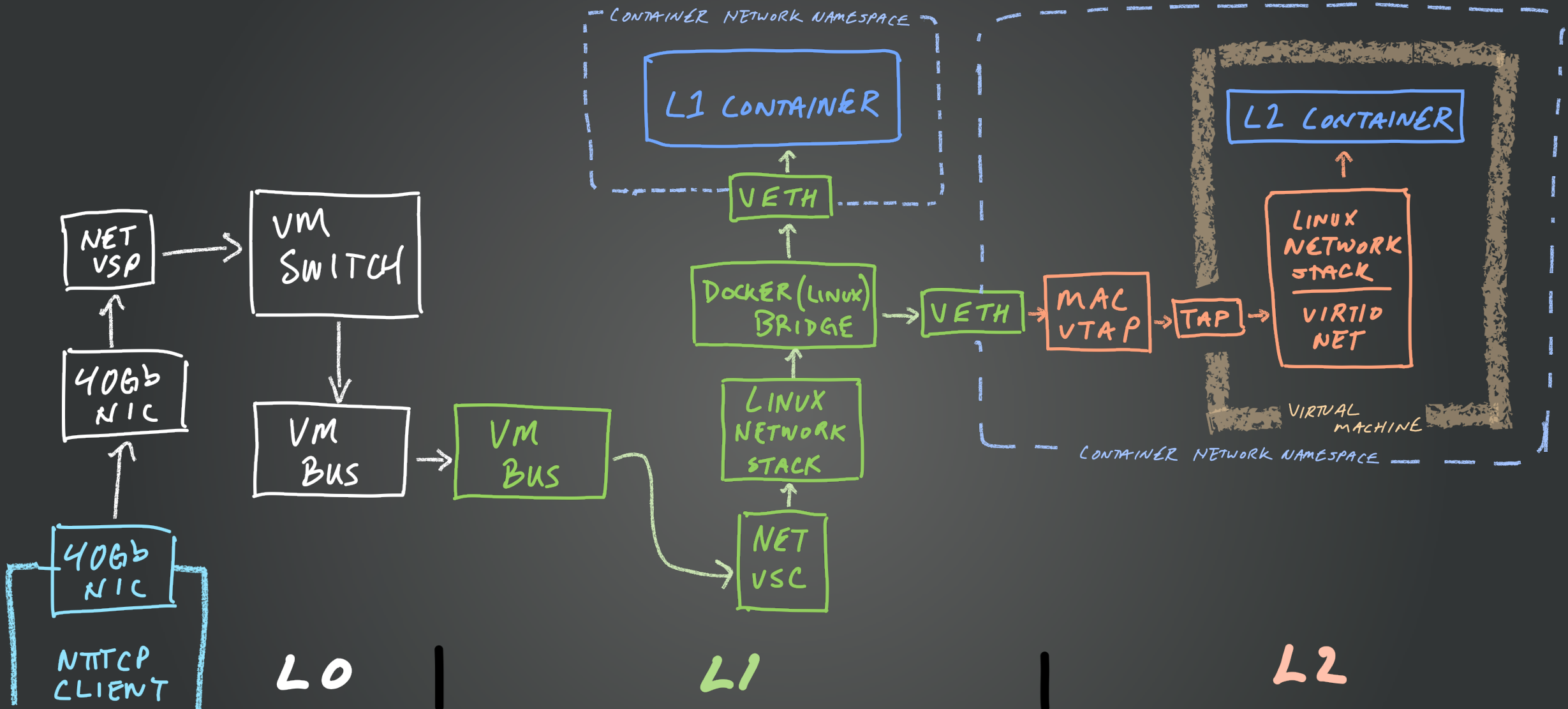
# Nested Kata: Network I/O - setup



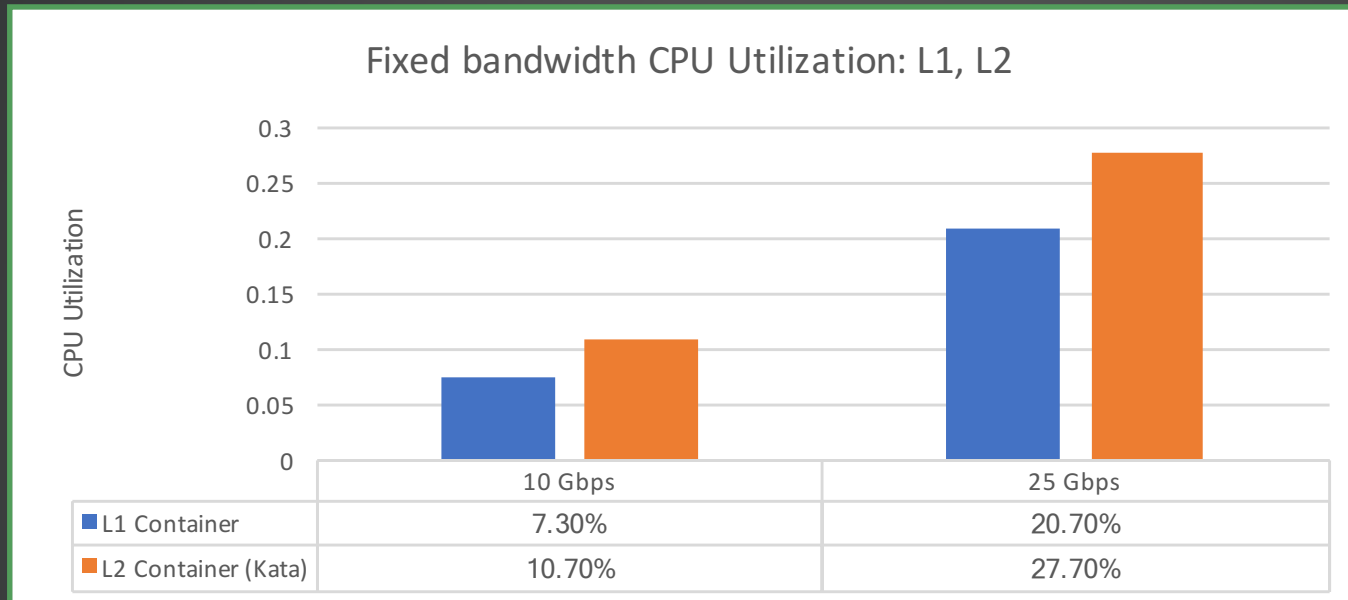
# Nested Kata: Network I/O - setup



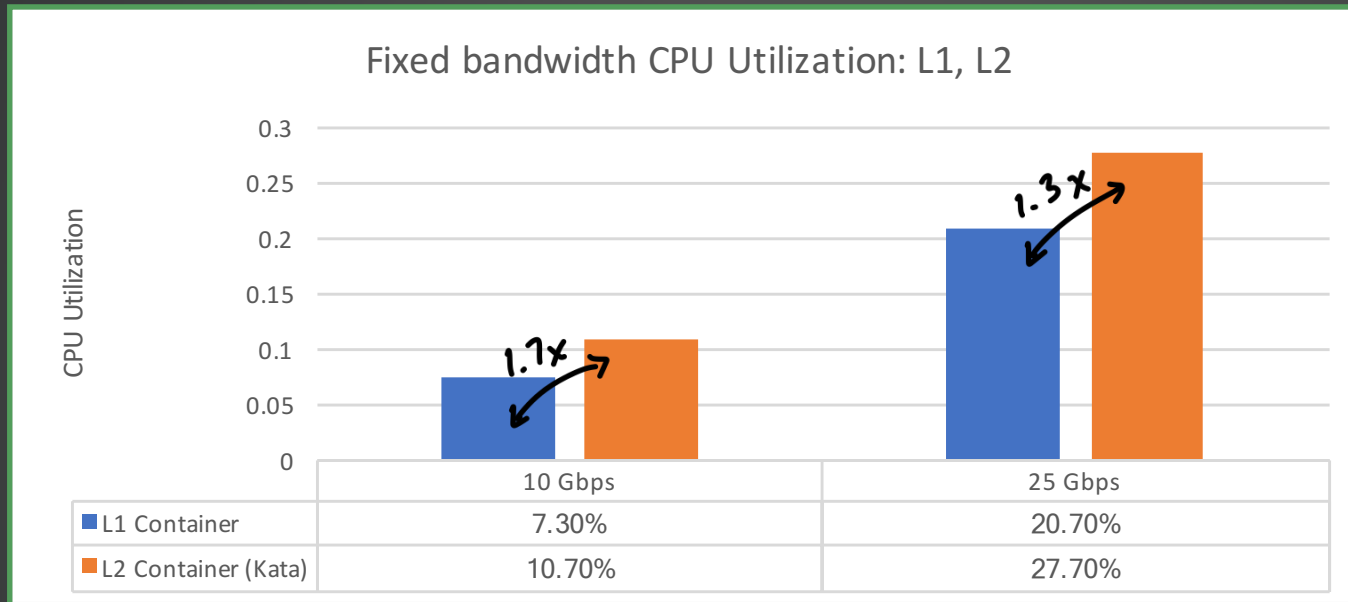
# Nested Kata: Network I/O - setup



# Nested Kata Network I/O



# Nested Kata Network I/O

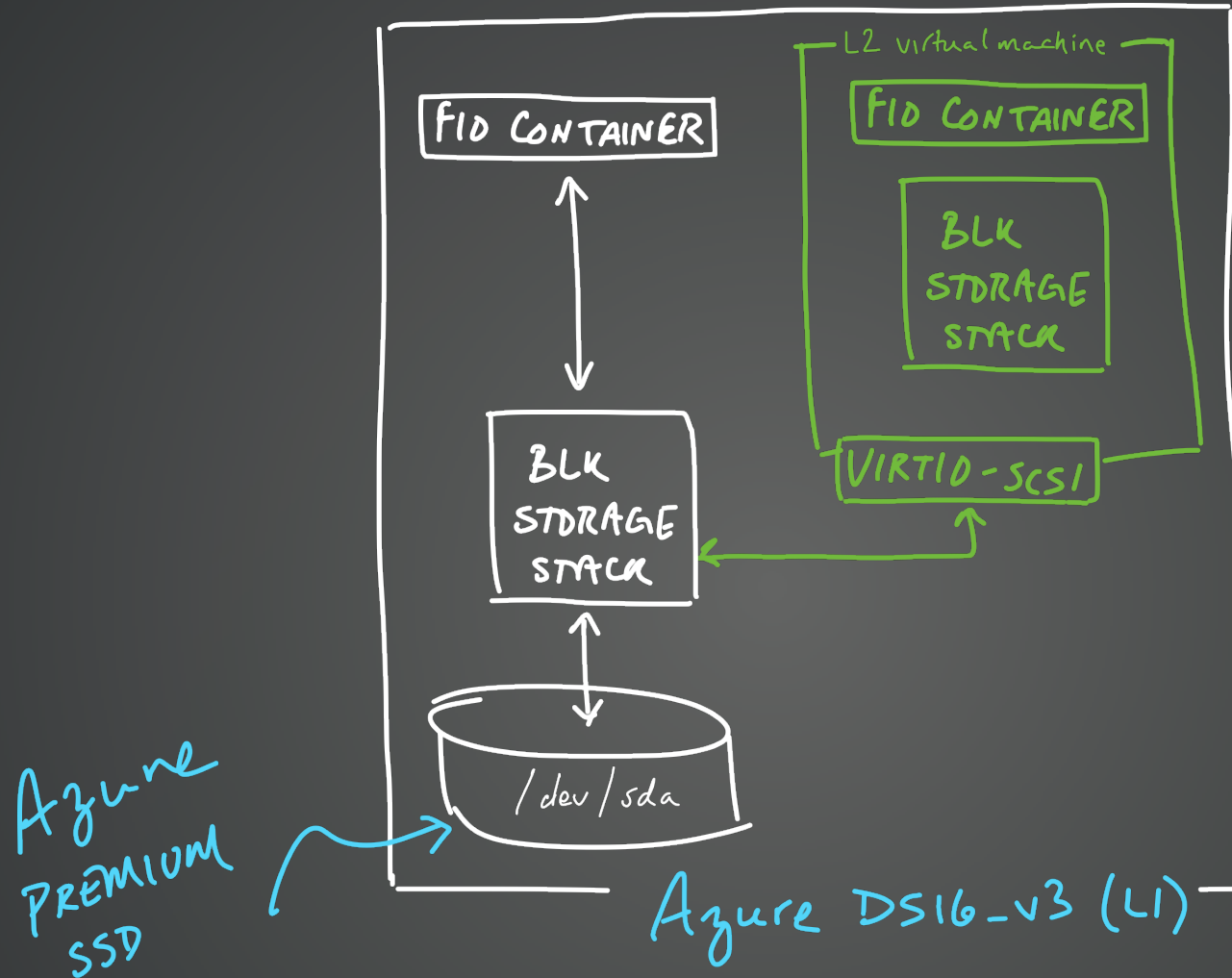


CORES USED

	<u>10Gbps</u>	<u>25Gbps</u>
L1:	3.2	9.1
L2:	4.7	12.2

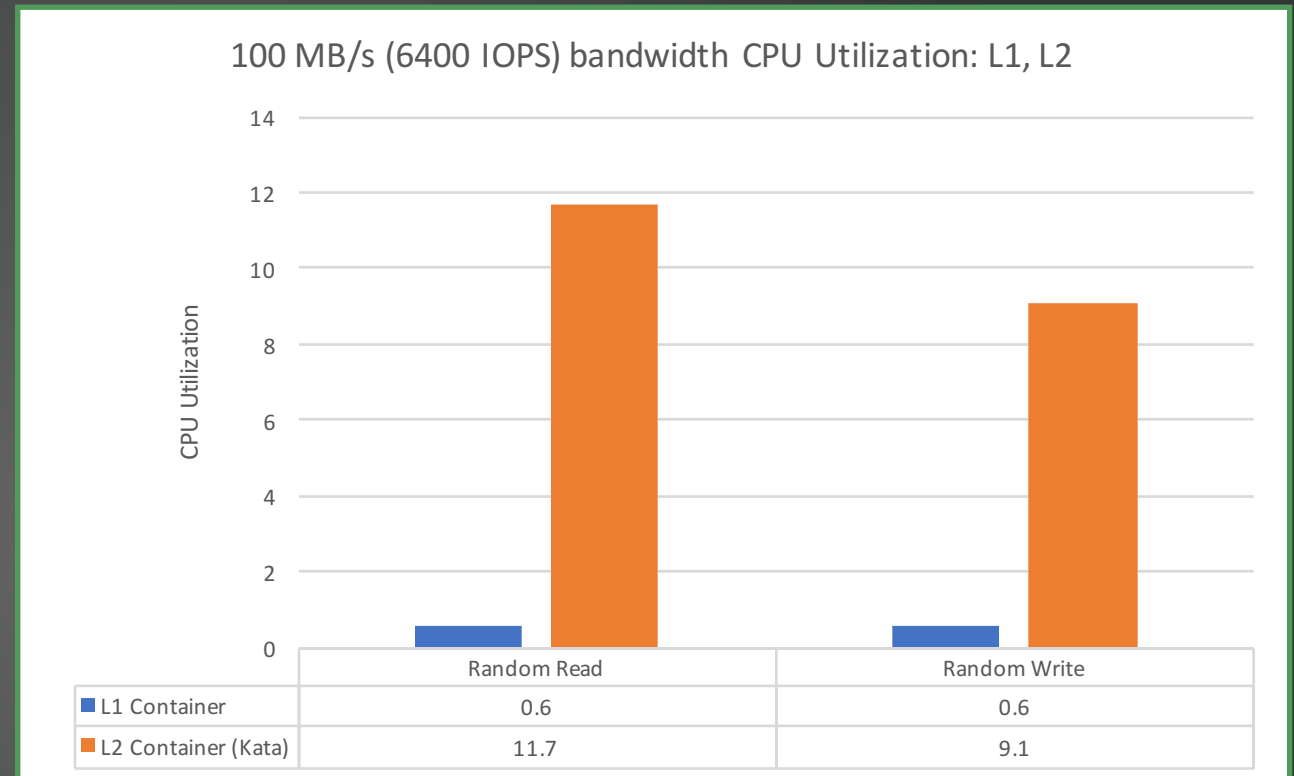
\* measured utilization in L0

# Nested Kata: Storage I/O - setup



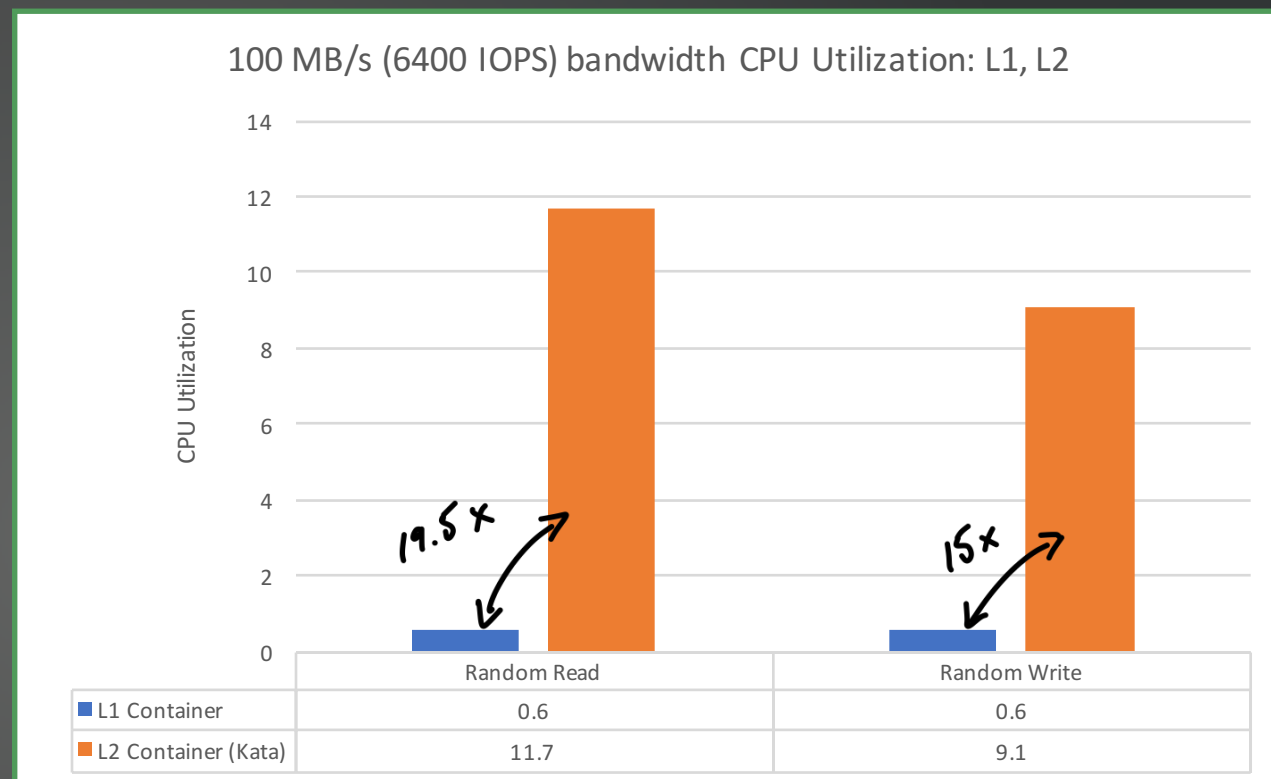
# Nested Kata Storage I/O

- Nested is relatively expensive
- High amount of iowait observed in L1 during L2 random read testing



# Nested Kata Storage I/O

- Nested is relatively expensive
- High amount of iowait observed in L1 during L2 random read testing



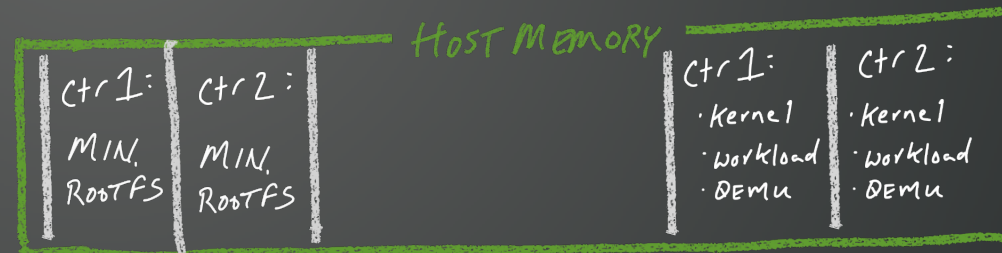
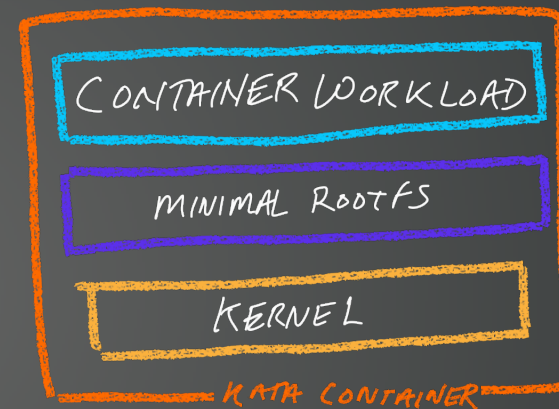
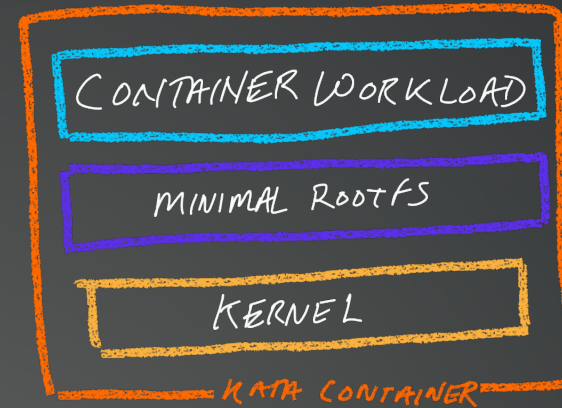
\* measured utilization  
in L1

RR  
RW  
L2: 1.8 cores, 1.5 cores



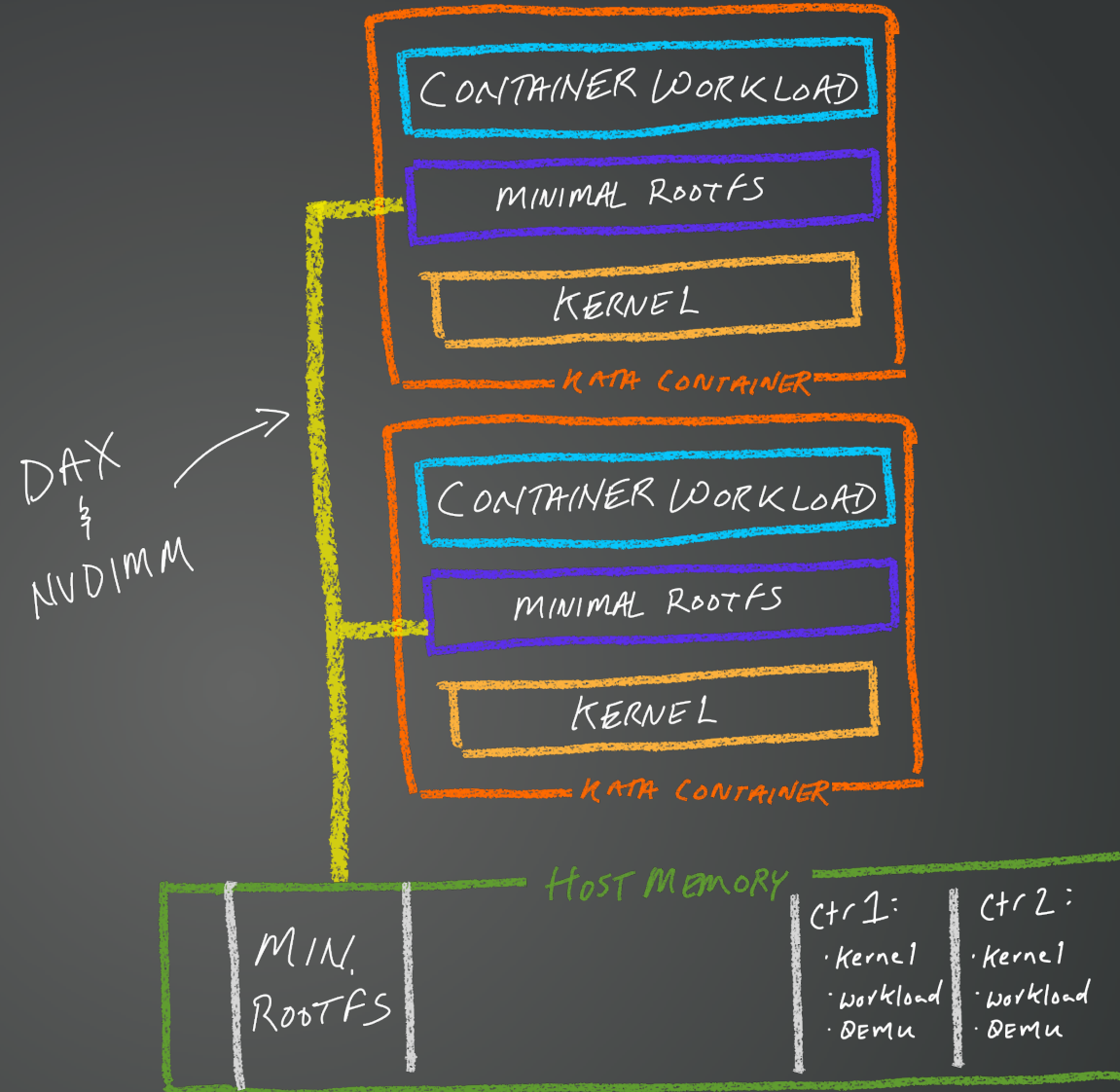
# Reducing Kata's footprint

- Minimal kernel
- Minimal rootfs
- Minimally configured QEMU



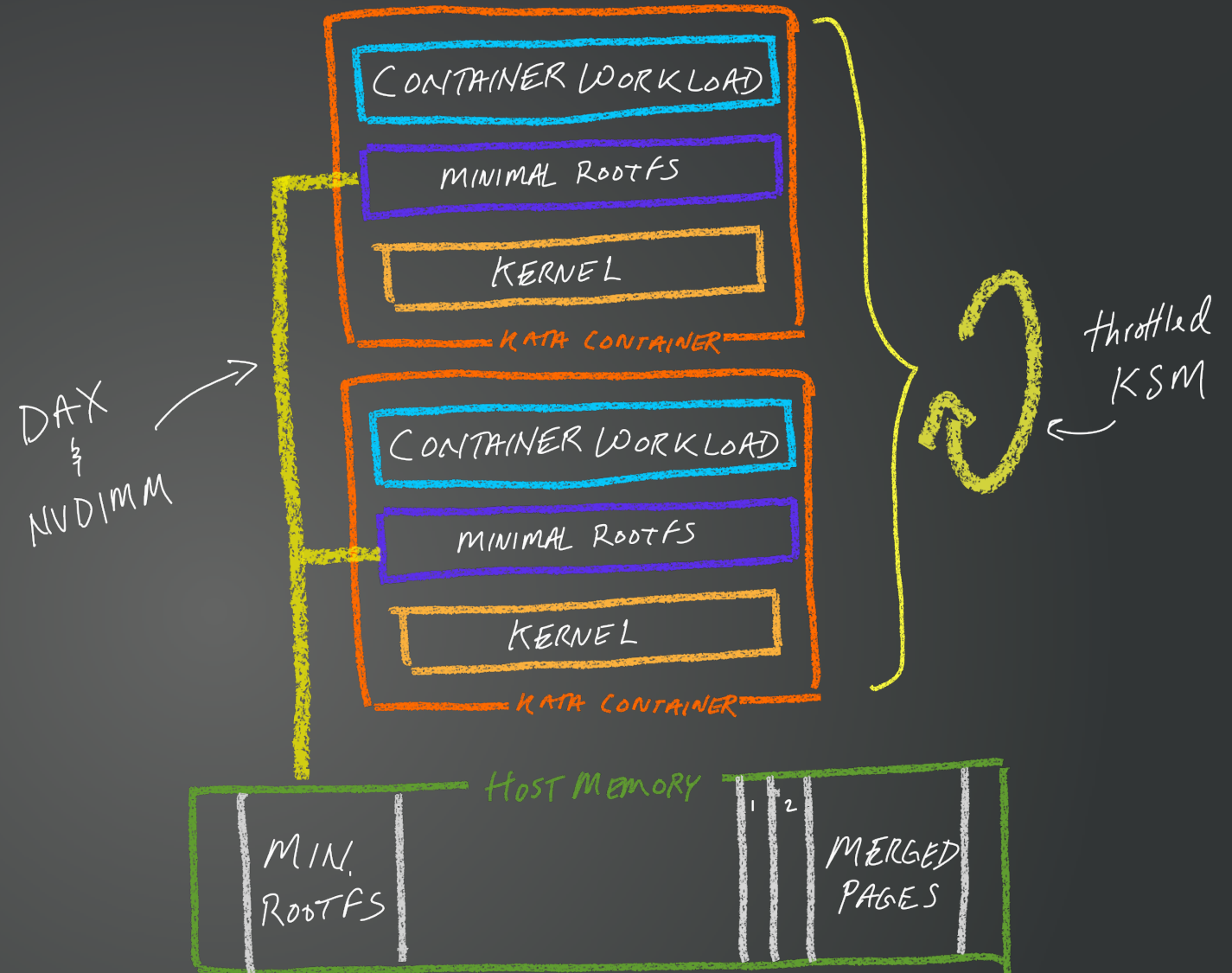
# Reducing Kata's footprint

- Minimal kernel
- Minimal rootfs
- Minimally configured QEMU
- DAX/NVDIMM



# Reducing Kata's footprint

- Minimal kernel
- Minimal rootfs
- Minimally configured QEMU
- DAX/NVDIMM
- De-duplicating memory



# Nested Kata container density

Dockerhub workload	Memory footprint		Containers/GB	
	Kata	runc	Kata	runc
busybox (small)	93.2 MB	682 KB	11	1535.7
mysql (medium)	135.5 MB	160.8 MB	7.6	6.5
elasticsearch (large)	2.5 GB	2.2 GB	0.4	0.5

“it depends”

A decorative graphic on the left side of the slide, consisting of several overlapping circles in various shades of green. The circles are partially cut off by the left edge of the frame. The background of the slide is a dark grey gradient.

# Summary, next steps

# Next Steps

- Nesting:
  - Continued improvements for KVM on Hyper-V
  - Optimizations for L2:
    - Investigate more efficient L2 storage options
    - General efficiency improvements to minimize nested cost
- Kata:
  - Improvements on density as well as security
  - Released support for NEMU

# Where can you get Kata?

- Dockerhub:
  - [katadocker/kata-deploy](#)
- Packages:
  - Clear linux, Snap
  - Built for CentOS, Fedora, SLES, RHEL, Ubuntu
- Running on public Cloud:
  - ACS-engine support in Azure
  - Anywhere bare-metal or nested virtualization is supported, including AWS, Azure, GCP, packet.net, Vexxhost



# OPEN SOURCE SUMMIT

EUROPE

THE LINUX FOUNDATION

