# Who I Think You Are

Software engineer, Sysadmin, etc who is...

- wanting to learn about namespaces and cgroups

- intereseted in containers and how they work

- loves turtles (optional)

Core OS

# Modern Linux Server with Containers

brandon.philips@coreos.com

# Overview

# Overview

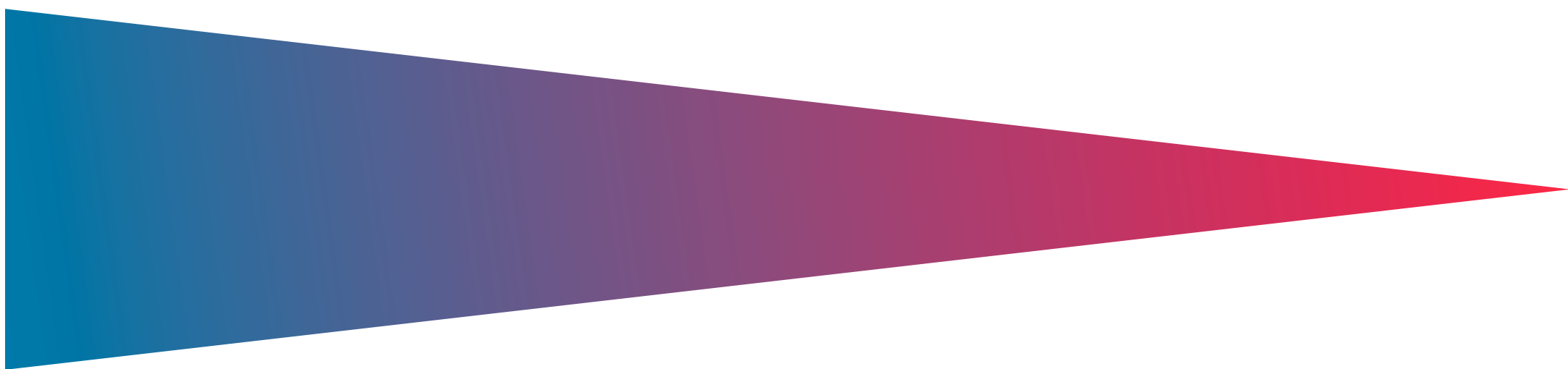- System Designs

# Overview

- System Designs

- Namespaces
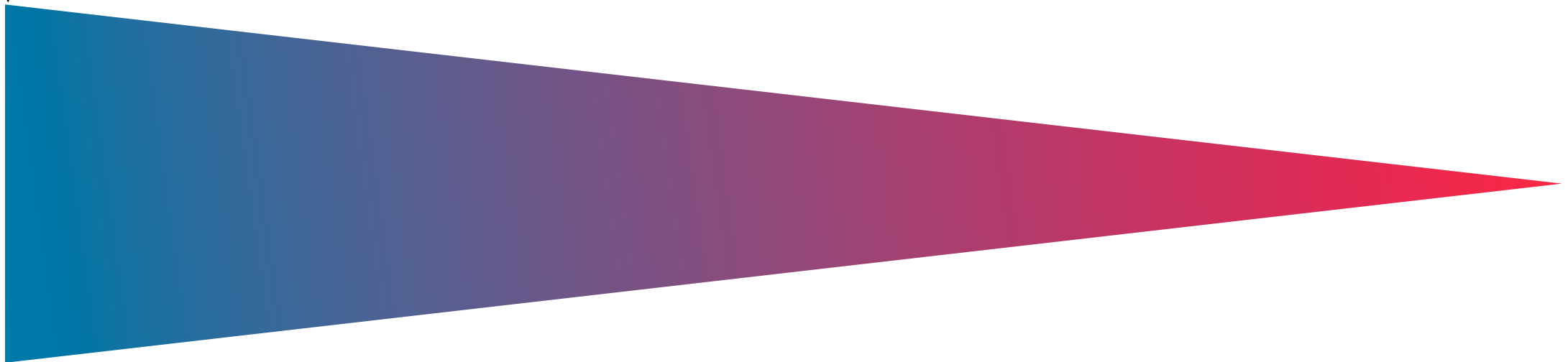
# Overview

- System Designs

- Namespaces

- Cgroups

Core OS

# Overview

- System Designs
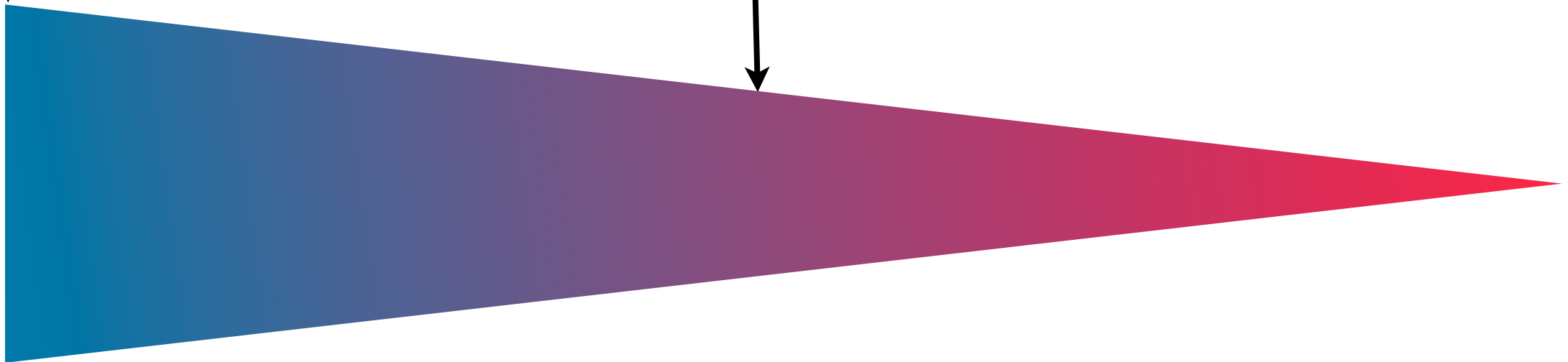
- Namespaces

- Cgroups

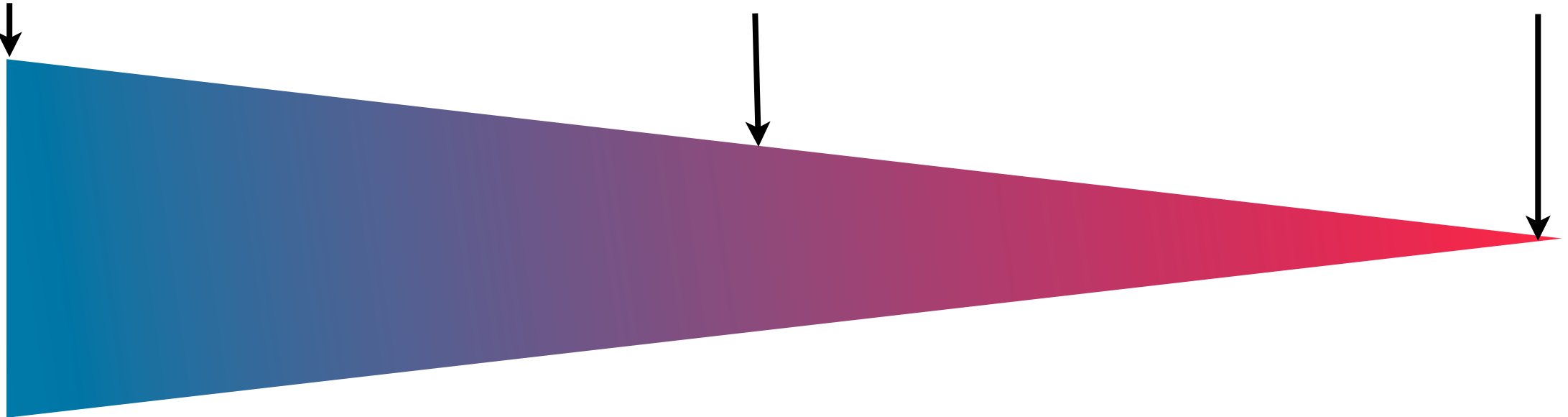- Tooling

Core OS

# The Spectrum

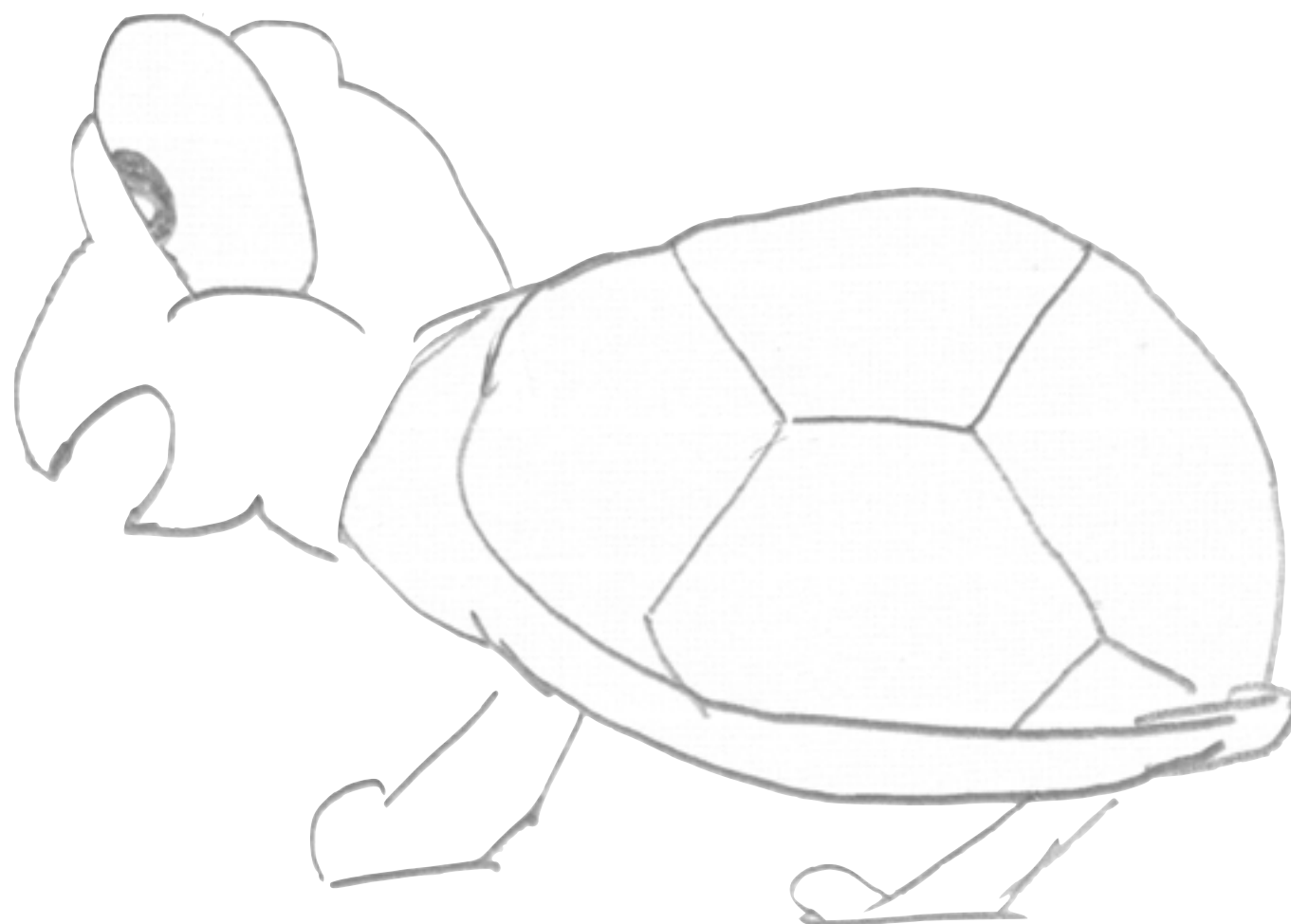Hypervisor

Core OS

Hypervisor          Container

Core OS

Hypervisor

Container

Application
Container

Core OS

# WARNING

Core OS

# System Designs

HARDWARE

LINUX

KVM
LINUX

KVM
LINUX

KVM
LINUX

KVM
LINUX

Core OS

# Hypervisor

# Hypervisor

- Host provides full hardware environment

Core OS

# Hypervisor

- Host provides full hardware environment

- Block device, ethernet device, etc

# Hypervisor

- Host provides full hardware environment

- Block device, ethernet device, etc

- Guests run a full kernel

Core OS

# Container

# Container

- Host provides Kernel

# Container

- Host provides Kernel

- Filesystem, network interface, etc are already there

Core OS

# Container

- Host provides Kernel

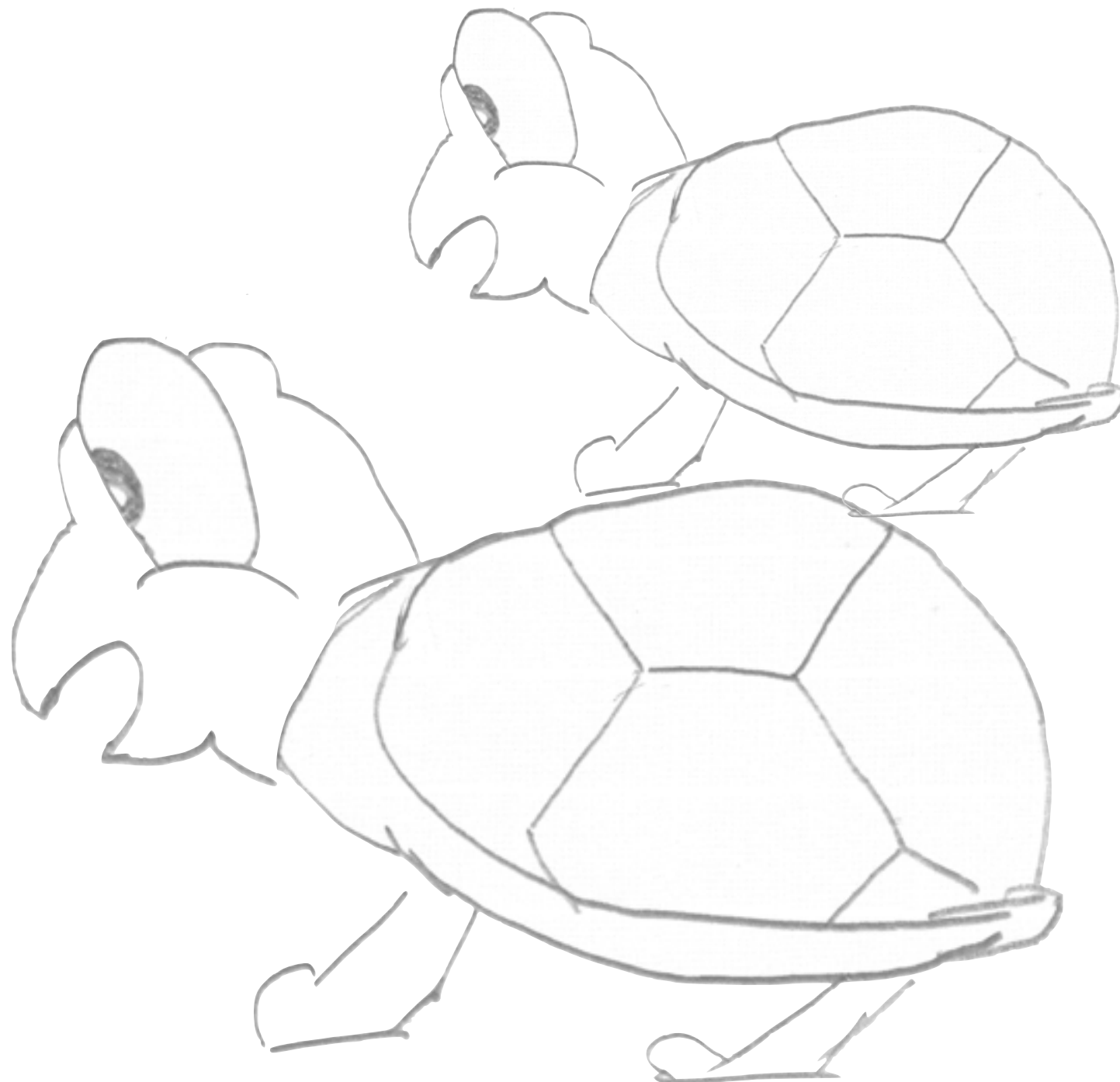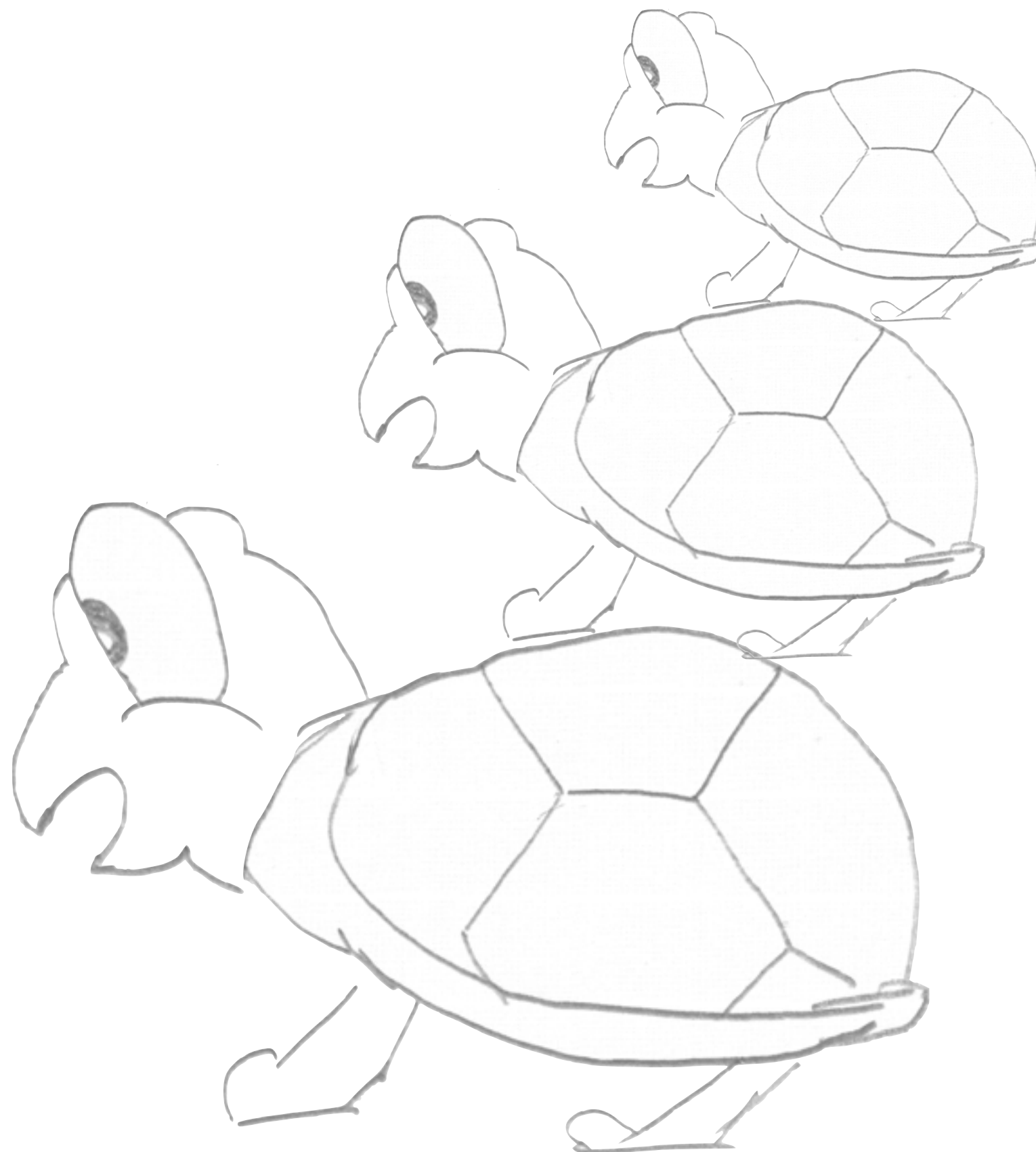- Filesystem, network interface, etc are already there

- Guest starts from /sbin/init

Core OS

# Application Container

# Application Container

- Host provides Kernel

# Application Container

- Host provides Kernel

- User data, socket fd, etc are already there

Core OS

# Application Container

- Host provides Kernel

- User data, socket fd, etc are already there

- Starts from application not init

Core OS

# Namespaces

**Imagine:** cool medieval castle photo
*perhaps fog rolling in*

# Filesystem

# Filesystem

- Read-only

# Filesystem

- Read-only

- Shared

Core OS

# Filesystem

- Read-only

- Shared

- Slave

Core OS

# Filesystem

- Read-only

- Shared

- Slave

- Private

Core OS

# Read-only

```
Private bind mount
before:
after:
source/a-file
bind/a-file

mount -t tmpfs -o size=1M tmpfs source/mnt
before:
after:
source/mnt/tmpfs-file

mount -t tmpfs -o size=1M tmpfs bind/mnt2
before:
after:
bind/mnt2/mnt2-file
```

Shared bind mount
before:
after:
source/a-file
bind/a-file

mount -t tmpfs -o size=1M tmpfs source/mnt
before:
after:
source/mnt/tmpfs-file
bind/mnt/tmpfs-file

mount -t tmpfs -o size=1M tmpfs bind/mnt2
before:
after:
source/mnt2/mnt2-file
bind/mnt2/mnt2-file

Core OS

```
Slave bind mount
before:
after:
source/a-file
bind/a-file

mount -t tmpfs -o size=1M tmpfs source/mnt
before:
after:
source/mnt/tmpfs-file
bind/mnt/tmpfs-file

mount -t tmpfs -o size=1M tmpfs bind/mnt2
before:
after:
bind/mnt2/mnt2-file
```

Core OS

# Patterns

- Mounting RO /usr inside a container

- Private /tmp per service

- Sharing data across containers via binds

# Networking

# Networking

- Root namespace

# Networking

- Root namespace

- Bridging

# Networking

- Root namespace

- Bridging

- Private namespace with socket activation

Core OS

# Root Namespace

- Full access to the machine interfaces

# Root Namespace

Core OS

# Root Namespace

- **Advantages**

Core OS

# Root Namespace

- **Advantages**

- Fast

# Root Namespace

- **Advantages**

- Fast

- Easy to get setup

# Root Namespace

- **Advantages**

- Fast

- Easy to get setup

- Network looks normal
  to the container

Core OS

# Root Namespace

- **Advantages**

- Fast

- Easy to get setup

- Network looks normal to the container

Core OS

# Root Namespace

- **Advantages**

- Fast

- Easy to get setup

- Network looks normal
  to the container

- **Disadvatages**

Core OS

# Root Namespace

- **Advantages**

- Fast

- Easy to get setup

- Network looks normal
  to the container

- **Disadvatages**

- No separation of
  concerns

Core OS

# Root Namespace

- **Advantages**

- Fast

- Easy to get setup

- Network looks normal
to the container

- **Disadvatages**

- No separation of
concerns

- Container has full
control

Core OS

# Network Bridges

# Network Bridges

- Create a bridge, like a virtual switch

CoreOS

# Network Bridges

- Create a bridge, like a virtual switch

- Create container namespace and add interface

Core OS

# Network Bridges

- Create a bridge, like a virtual switch

- Create container namespace and add interface

- Attach container interface to bridge

Core OS

# Network Bridges

# Network Bridges

- **Advantages**

Core OS

# Network Bridges

- **Advantages**

- More complex to get setup

# Network Bridges

- **Advantages**

- More complex to get setup

- Network looks normal to the container

Core OS

# Network Bridges

- **Advantages**

- More complex to get setup

- Network looks normal to the container

Core OS

# Network Bridges

- **Advantages**

- More complex to get setup

- Network looks normal to the container

Core OS

# Network Bridges

- **Advantages**

- More complex to get setup

- Network looks normal to the container

- **Disadvantages**

Core OS

# Network Bridges

- **Advantages**

- More complex to get setup

- Network looks normal to the container

- **Disadvantages**

- Less speed

Core OS

# Network Bridges

- **Advantages**

- More complex to get setup

- Network looks normal to the container

- **Disadvantages**

- Less speed

- NAT to the internet

Core OS

# Network Bridges

- **Advantages**

- More complex to get setup

- Network looks normal to the container

- **Disadvantages**

- Less speed

- NAT to the internet

- iptables to expose public socket

# Socket Activation

# Socket Activation

- No interface

# Socket Activation

- No interface

- Sockets are passed via stdin (inetd)

Core OS

# Socket Activation

- No interface

- Sockets are passed via stdin (inetd)

- systemd style listen fd API

Core OS

# inetd style

# inetd style

- **Advantages**

Core OS

# inetd style

- **Advantages**

- Fast and isolated

Core OS

# inetd style

- **Advantages**

- Fast and isolated

- Simple and well understood

# inetd style

- **Advantages**

- Fast and isolated

- Simple and well understood

- Support from existing daemons like ssh

Core OS

# inetd style

- **Advantages**

- Fast and isolated

- Simple and well understood

- Support from existing daemons like ssh

- No process running until needed

Core OS

# inetd style

- **Advantages**

  - Fast and isolated

  - Simple and well understood

  - Support from existing daemons like ssh

  - No process running until needed

- **Disadvantages**

Core OS

# inetd style

- **Advantages**

- Fast and isolated

- Simple and well understood

- Support from existing daemons like ssh

- No process running until needed

- **Disadvantages**

- One process per client (scaling problems!)

Core OS

# listen fd style

# listen fd style

- **Advantages**

# listen fd style

- **Advantages**

- Fast and isolated

# listen fd style

- **Advantages**

- Fast and isolated

- Only one process needed per service

Core OS

# listen fd style

- **Advantages**

- Fast and isolated

- Only one process needed per service

- No process running until needed

# listen fd style

- **Advantages**

- Fast and isolated

- Only one process needed per service

- No process running until needed

Core OS

# listen fd style

- **Advantages**

- Fast and isolated

- Only one process needed per service

- No process running until needed

- **Disadvantages**

Core OS

# listen fd style

- **Advantages**

- Fast and isolated

- Only one process needed per service

- No process running until needed

- **Disadvantages**

- Patches required to daemons

Core OS

# Process Namespace

- PID 1 is something else outside the namespace

# All the Rest

Core OS

# Cgroups

**Imagine**: an accountant's overflowing desk
perhaps hands on head in dispair

Core OS

# Block I/O

- **Limit**: Weight from 10 to 1000

- **Limit**: Bandwidth limits R/W

- **Metrics:** iops serviced, waiting and queued

# CPU

- **Limit**: Shares system 1024 is half of 2048

- **Metrics:** cpuacct.stats user and system

# Memory

- **Limit**: Total RSS memory limit

- **Metrics**: swap, total rss, # page ins/outs

Core OS

# Tooling

# docker

# nspawn

# nsenter

# /sys/fs/cgroup

# systemd units

# systemd-cgtop

# Recap

# Recap

- Containers are built on namespaces and cgroups

Core OS

# Recap

- Containers are built on namespaces and cgroups

- Namespaces provide isolation similar to hypervisors

# Recap

- Containers are built on namespaces and cgroups

- Namespaces provide isolation similar to hypervisors

- Cgroups provide resource limiting and accounting

Core OS

# Recap

- Containers are built on namespaces and cgroups

- Namespaces provide isolation similar to hypervisors

- Cgroups provide resource limiting and accounting

- These tools can be mixed to create hybrids

Core OS

# Future

# Thanks!

@BrandonPhilips
@CoreOSLinux