# A Network in a Laptop: Rapid Prototyping for Software-Defined Networks

Bob Lantz, Brandon Heller, Nick McKeown

Stanford University

HotNets 2010, 10/20/10

# Wouldn't it be amazing...
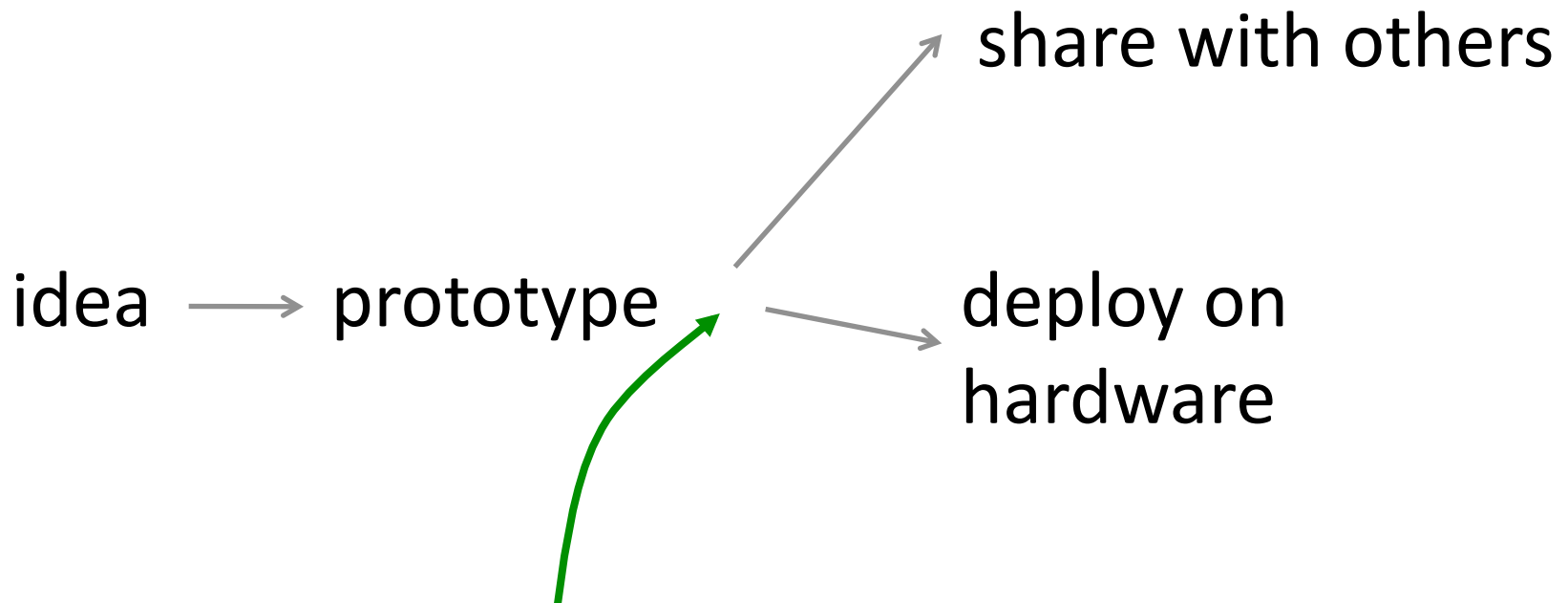
if systems papers were *runnable*.

# Wouldn't it be amazing...

If systems papers made replicating their results, modifying the described system, and sharing it with others...

*... as easy as downloading a file.*

# Wouldn't it be amazing...

if network systems papers were *more* than runnable.

share with others

idea → prototype → deploy on hardware

with no code changes!?!

Mininet: a platform for rapid *network* prototyping.

scales to **usefully large nets**
runs **unmodified applications**
provides **path to hardware**
facilitates **sharing**
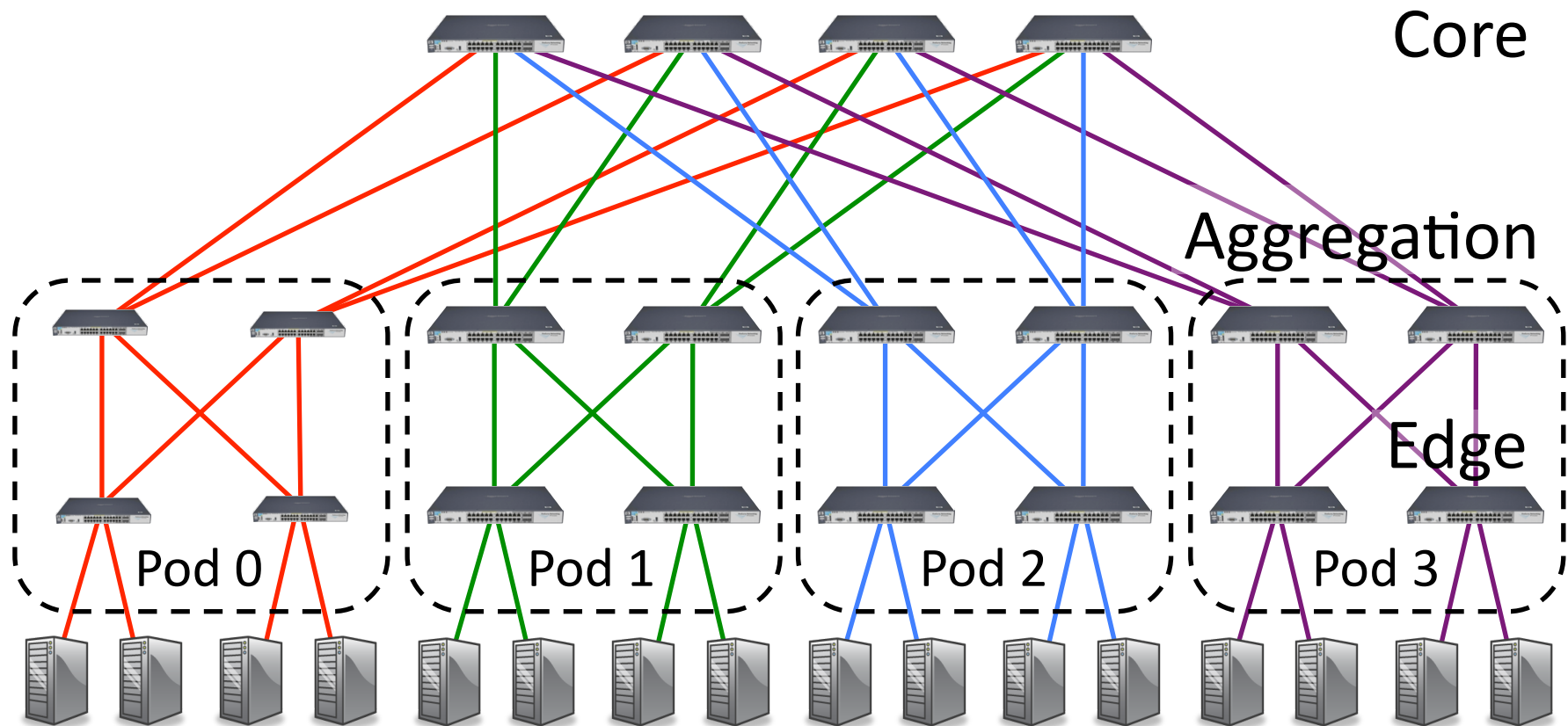
# openflow.org/mininet

140+ users

45+ on mailing list

20+ institutions

open source (BSD license)

[don't download now! save the WiFi!]

# Demo

# Demo Topology: Fat Tree



Core

Aggregation

Edge

Pod 0    Pod 1    Pod 2    Pod 3

described in Scalable Commodity Data Center, SIGCOMM 2008, Al Fares et al.
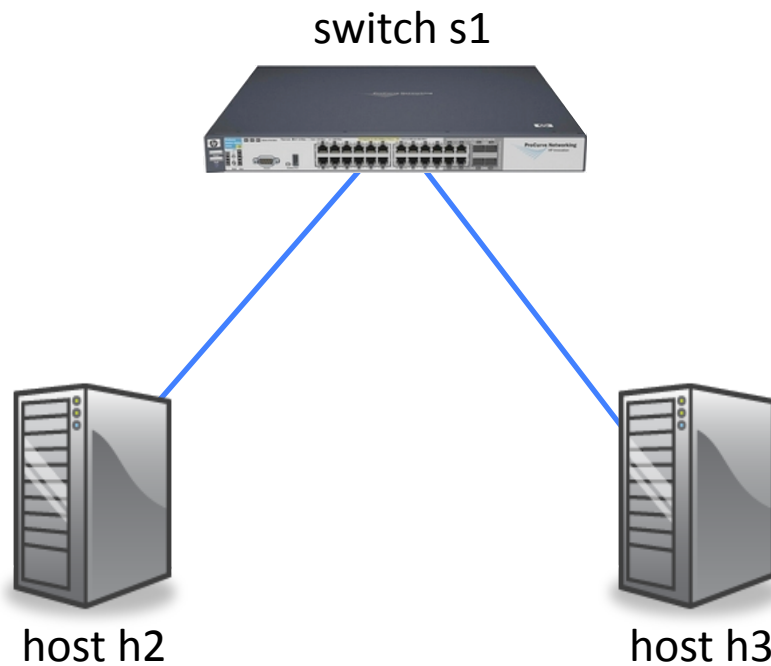
(1) share-able

(2) runs on hardware

Date
Nov 2009: deadline in 3 months

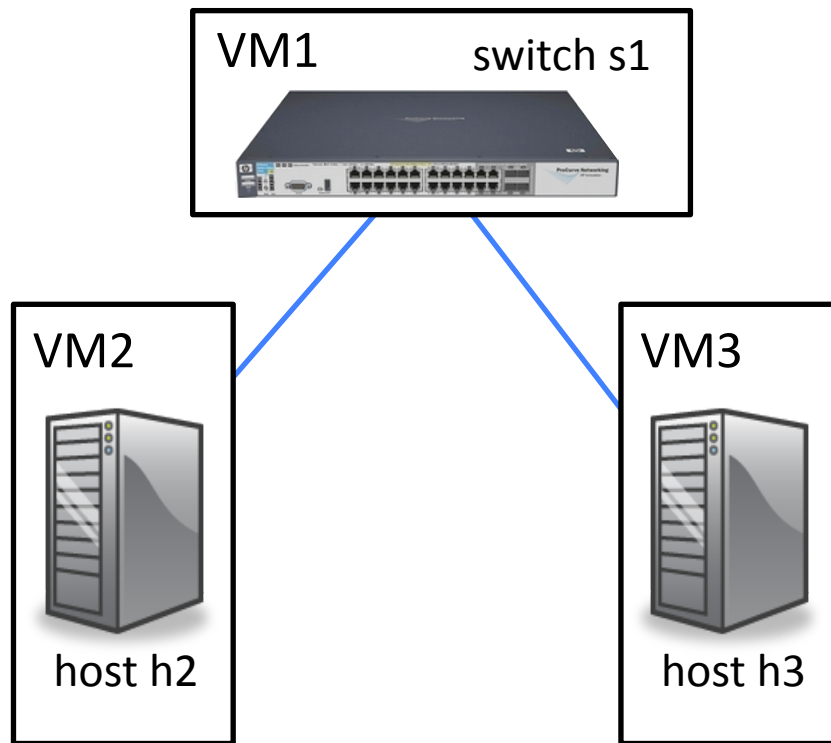[based on a true story]

# Resources: a laptop

Goal:
build/eval/demo a
realistic new networked
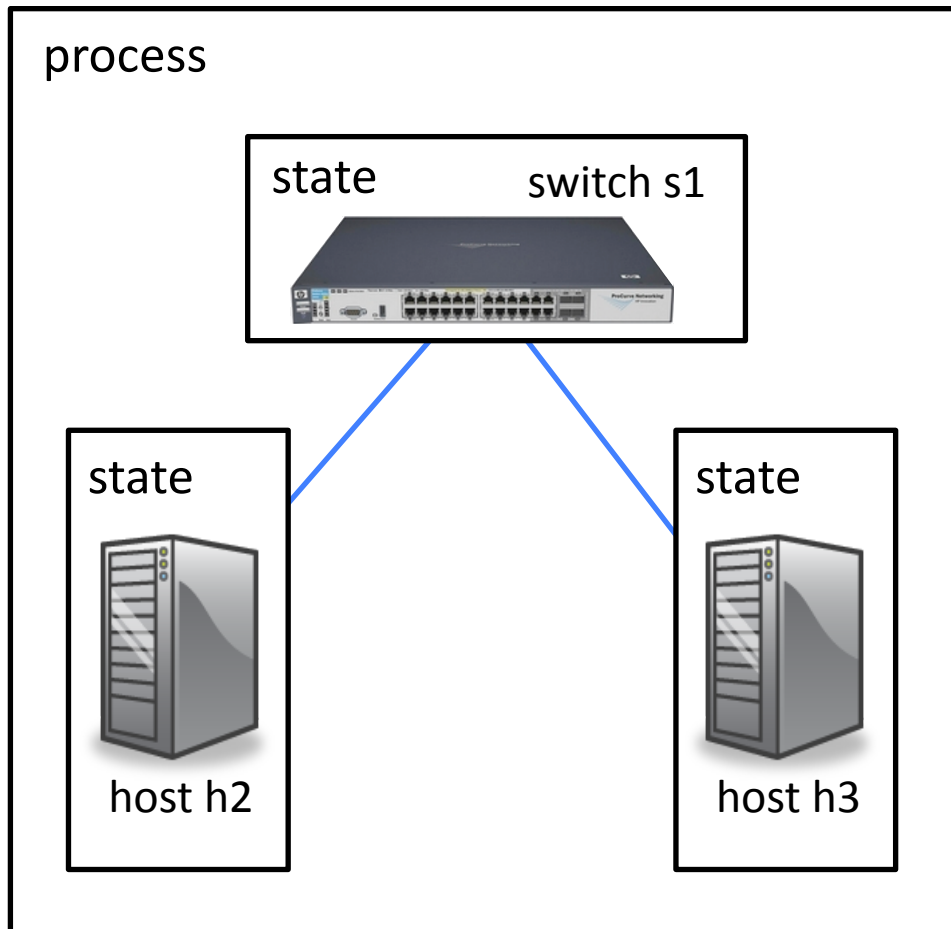system

# Why not a real system?

switch s1

host h2          host h3

+ as real as it gets
- a pain to reconfigure

# Why not networked virtual machines?

VM1  switch s1

VM2  host h2

VM3  host h3

+ easier topology changes
- scalability

# Why not a simulator?



+ good visibility
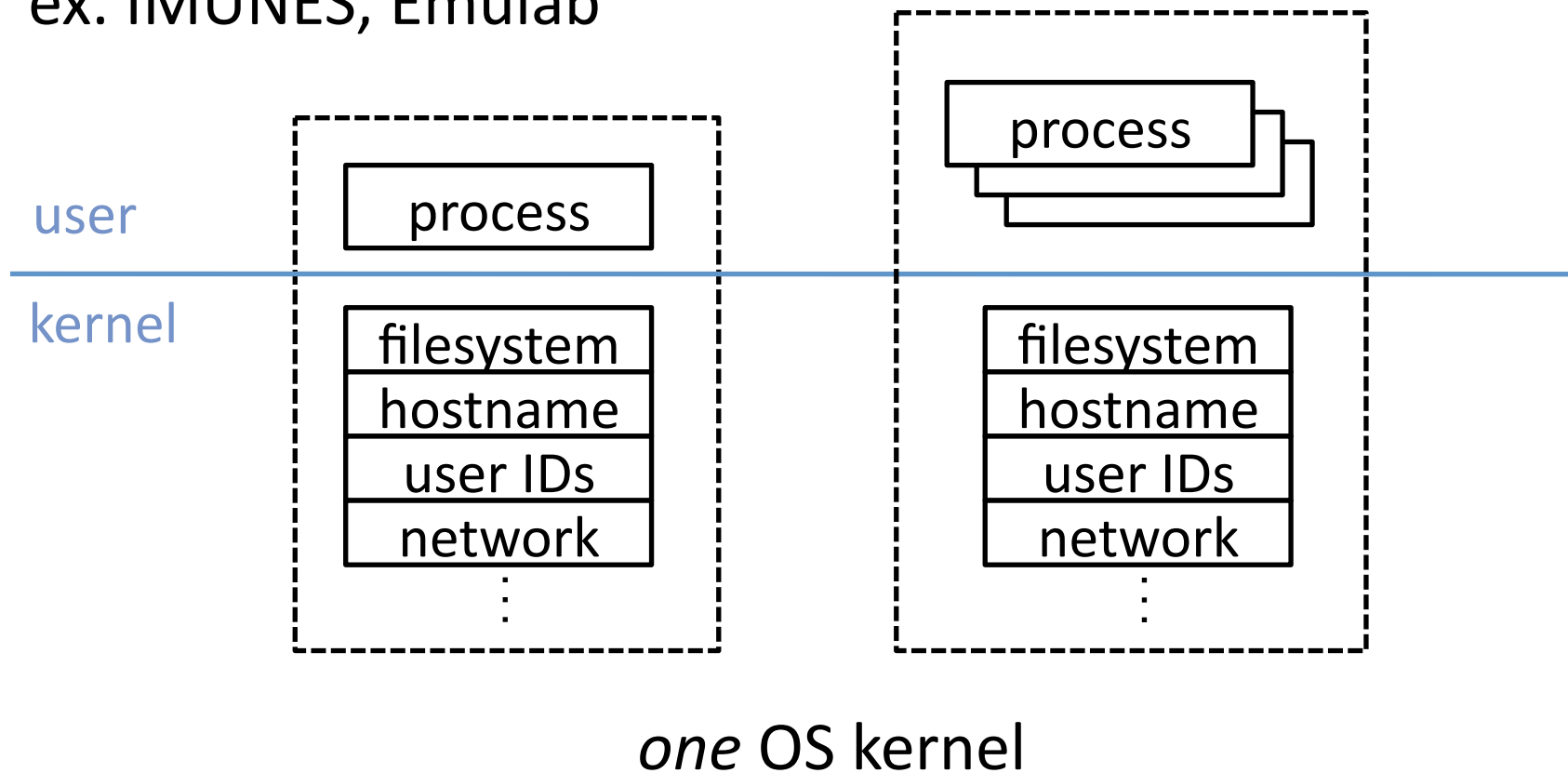- no path to hardware

Problem 1:
Want scale with unmodified applications.

→ Use lightweight, OS-level virtualization.

# OS-level Virtualization

Same system, different view.  Almost zero overhead.

ex. IMUNES, Emulab

user

kernel

process

process

| filesystem |
|---|
| hostname |
| user IDs |
| network |

⋮

| filesystem |
|---|
| hostname |
| user IDs |
| network |

⋮

*one* OS kernel
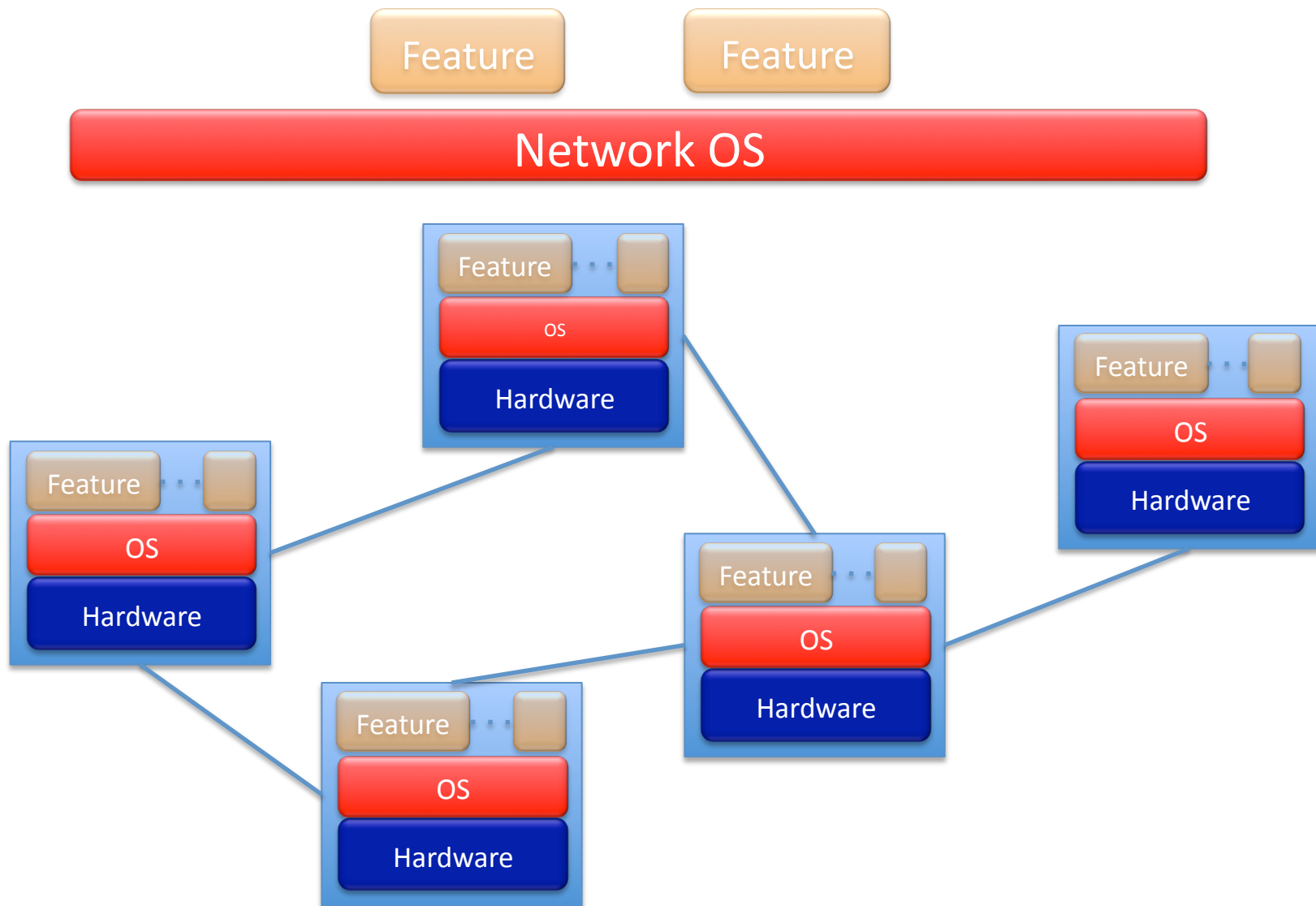
Problem 2:
Want a smooth path to hardware deployment.

→ Use Software-Defined Networking.

Feature    Feature

Network OS

Feature · · ·    OS    Hardware

Feature · · ·    OS    Hardware

Feature · · ·    OS    Hardware

Feature · · ·    OS    Hardware

Feature · · ·    OS    Hardware

# Software-Defined Network



Feature   Feature

Network OS

OpenFlow

Packet Forwarding

Packet Forwarding

Packet Forwarding

Packet Forwarding

Packet Forwarding

22

# Mininet Walkthrough

controller

switch s1

host h2

host h3

```
$> mn --topo minimal \
      --switch ovsk \
      --controller ref
```

# run Mininet launcher

```
┌─────────────────────────────────────────┐
│         root network namespace          │
│                                         │
│   ┌─────────────┐                       │
│   │     mn      │                       │
│   └─────────────┘                       │
│                                         │
│                                         │
│                                         │
└─────────────────────────────────────────┘
```

```
$> mn --topo minimal \
      --switch ovsk \
      --controller ref
```

# Hosts

## create bash processes

root network namespace

mn

pipes

/bin/bash

host h2

/bin/bash

host h3

```
$> mn --topo minimal \
       --switch ovsk \
       --controller ref
```

Hosts

unshare(CLONE_NEWNET)

root network namespace

mn

pipes

/bin/bash

h2 namespace

/bin/bash

h3 namespace

```
$> mn --topo minimal \
      --switch ovsk \
      --controller ref
```

# Links

## ip link add

root network namespace

| veth0 | | veth2 |

veth pairs

| veth1 | | veth3 |

mn

pipes

/bin/bash

h2 namespace

/bin/bash

h3 namespace

```
$> mn --topo minimal \
      --switch ovsk \
      --controller ref
```

Links

ip link set name

root network namespace

s1-eth1        s1-eth2

veth pairs

mn

h2-eth0        h3-eth0

pipes

/bin/bash                      /bin/bash

h2 namespace                   h3 namespace

```
$> mn --topo minimal \
       --switch ovsk \
       --controller ref
```

Links

`ip link set netns`



root network namespace

mn

pipes

s1-eth1   s1-eth2

veth pairs

h2-eth0

h3-eth0

/bin/bash

/bin/bash

h2 namespace

h3 namespace

$> mn --topo minimal \
    --switch ovsk \
    --controller ref

# Switch

## create OpenFlow Switch

root network namespace

ofdatapath —— ofprotocol

unix socket

mn

raw sockets

pipes

s1-eth1    s1-eth2

veth pairs

h2-eth0

/bin/bash

h2 namespace

h3-eth0

/bin/bash

h3 namespace

```
$> mn --topo minimal \
      --switch ovsk \
      --controller ref
```

# Controller

## create controller

root network namespace

| ofdatapath |   | ofprotocol |

unix socket

controller

mn

raw sockets

pipes

| s1-eth1 | | s1-eth2 |

veth pairs

| h2-eth0 |

| h3-eth0 |

| /bin/bash |

| /bin/bash |

h2 namespace

h3 namespace

Virtual Machine

root network namespace

ofdatapath — unix socket — ofprotocol

controller

mn

raw sockets

pipes

s1-eth1    s1-eth2

veth pairs

h2-eth0    h3-eth0

/bin/bash    /bin/bash

h2 namespace    h3 namespace

# Mininet example commands

Create a network using mn launcher:
```
mn --switch ovsk --controller nox --topo \ tree,depth=2,fanout=8
    --test pingAll
```
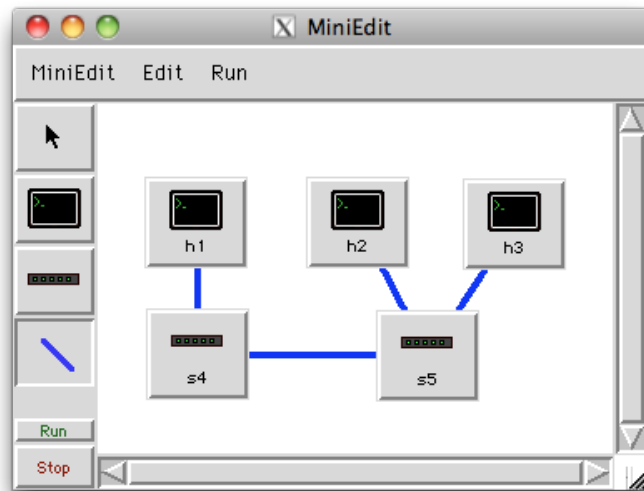
Interact with a network using CLI:
```
mininet> h2 ping h3
mininet> h2 py dir(locals())
```

Customize a network w/API:
```
from mininet.net import Mininet
from mininet.topolib import TreeTopo
tree4 = TreeTopo(depth=2,fanout=2)
net = Mininet(topo=tree4)
net.start()
h1, h4 = net.hosts[0], net.hosts[3]
print h1.cmd('ping -c1 %s' % h4.IP())
net.stop()
```

# Apps made with the Mininet API

# Evaluation

# Startup/Shutdown/Memory

| Topology | $H$ | $S$ | Setup(s) | Stop(s) | Mem(MB) |
|---|---|---|---|---|---|
| Minimal | 2 | 1 | 1.0 | 0.5 | 6 |
| Linear(100) | 100 | 100 | 70.7 | 70.0 | 112 |
| VL2(4, 4) | 80 | 10 | 31.7 | 14.9 | 73 |
| FatTree(4) | 16 | 20 | 17.2 | 22.3 | 66 |
| FatTree(6) | 54 | 45 | 54.3 | 56.3 | 102 |
| Mesh(10, 10) | 40 | 100 | 82.3 | 92.9 | 152 |
| Tree(4^4) | 256 | 85 | 168.4 | 83.9 | 233 |
| Tree(16^2) | 256 | 17 | 139.8 | 39.3 | 212 |
| Tree(32^2) | 1024 | 33 | 817.8 | 163.6 | 492 |

lots of switches & hosts
w/reasonable amounts of memory

# Microbenchmarks

| Operation | Time (ms) |
|---|---|
| Create a node (host/switch/controller) | 10 |
| Run command on a host ('echo hello') | 0.3 |
| Add link between two nodes | 260 |
| Delete link between two nodes | 416 |
| Start user space switch (OpenFlow reference) | 29 |
| Stop user space switch (OpenFlow reference) | 290 |
| Start kernel switch (Open vSwitch) | 332 |
| Stop kernel switch (Open vSwitch) | 540 |

link management is slow

# Bandwidth

| $S$ (Switches) | User(Mbps) | Kernel(Mbps) | Ratio |
|---|---|---|---|
| 1 | 445 | 2120 | ~5x |
| 10 | 49.9 | 940 | |
| 20 | 25.7 | 573 | |
| 40 | 12.6 | 315 | |
| 60 | 6.2 | 267 | |
| 80 | 4.15 | 217 | |
| 100 | 2.96 | 167 | ~50x |

## usable amount of bandwidth

# Case Studies

# Research Examples

- Ripcord: modular data center

- Asterix: wide-area load balancing

- SCAFFOLD: new internet architecture
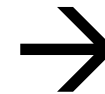
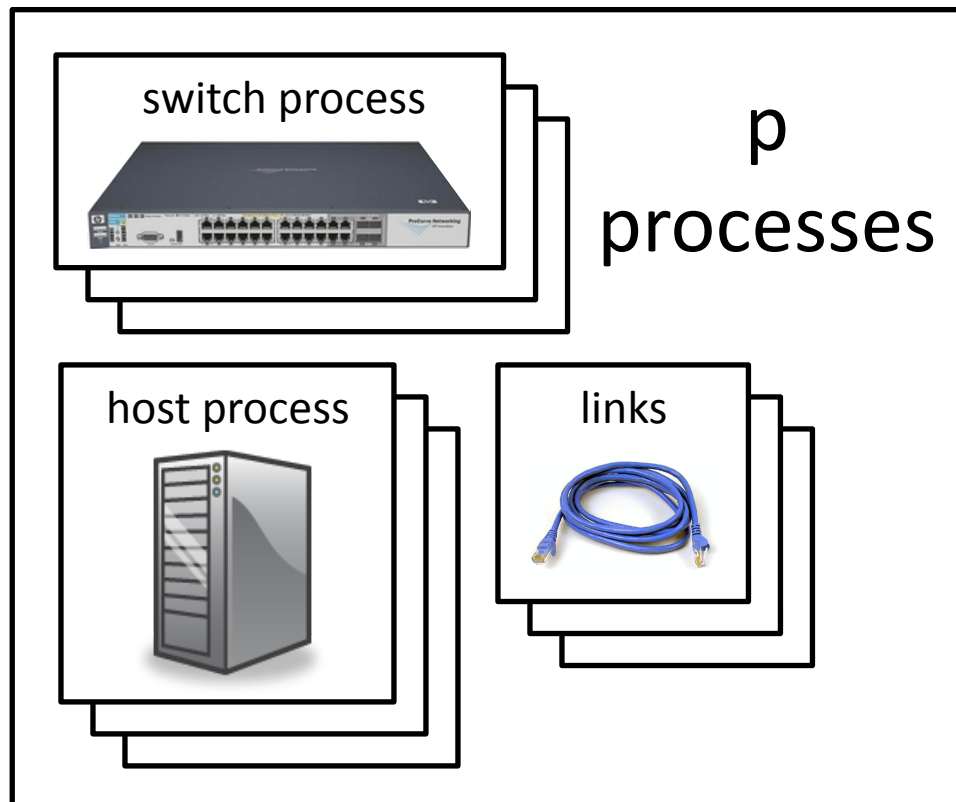- Distributed snapshot demo

# Unexpected Uses

- Tutorials
- Whole-network regression suites
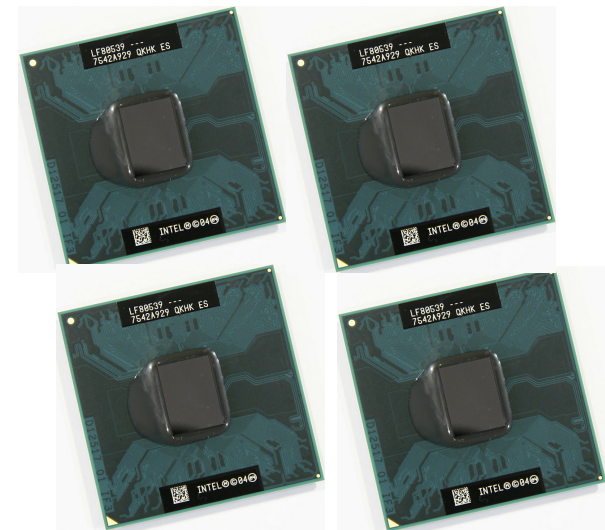- Bug replication

# Limitations

# Inherent Limitations

- OS-level virtualization → one kernel only
- Linux containers → Linux programs only
- Cannot match the introspection of an event-driven simulation

# Performance Fidelity



switch process

p processes

c cores

host process

links

→

p >> c ? time multiplexing
Issues: performance predictability, isolation

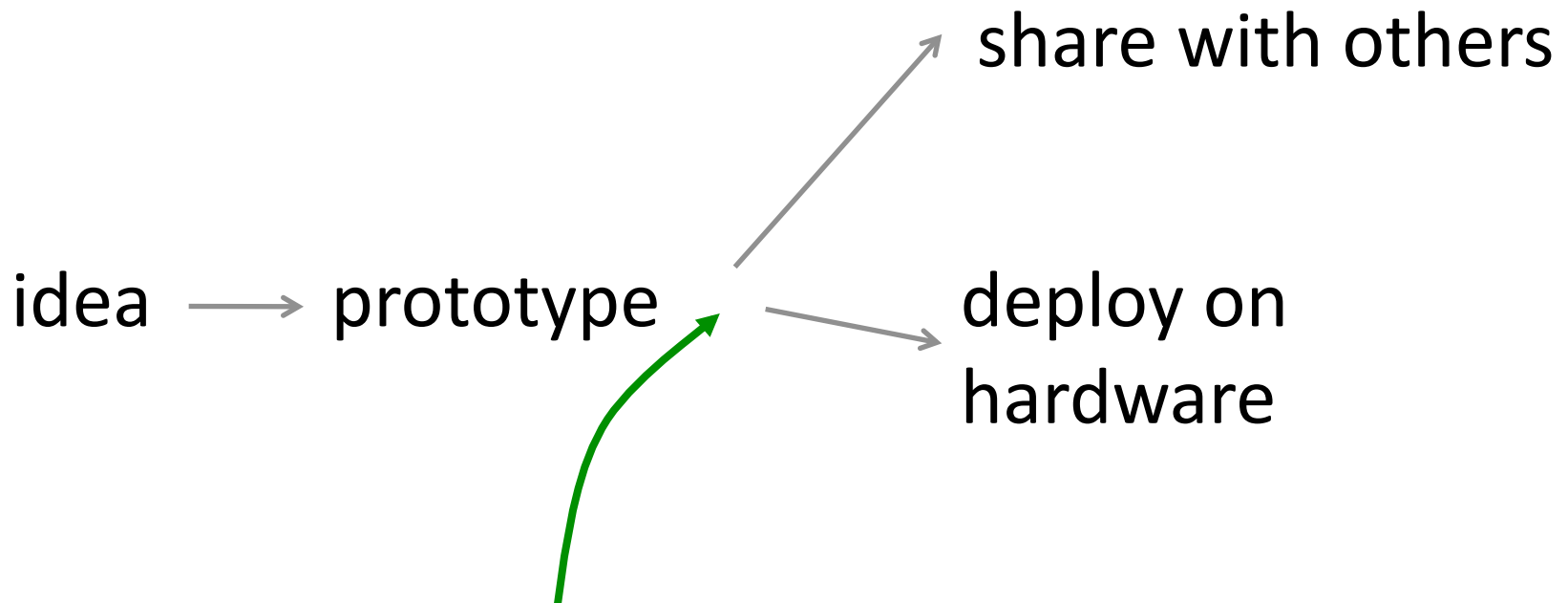# Wouldn't it be amazing…

if systems papers were *runnable*.

# Wouldn't it be amazing…

If systems papers made replicating their results, modifying the described system, and sharing it with others…

*… as easy as downloading a file.*

# Wouldn't it be amazing…

if network systems papers were *more* than runnable.

idea → prototype → share with others

prototype → deploy on hardware

with no code changes!?!

"A Network in a Laptop…"
is a runnable paper

…which itself describes how
to make other runnable
papers.

# Mininet

- Rapid prototyping

- Scalable

- Shareable

- Functionally correct

- Path to hardware

enables
"runnable papers"
for a subset of
networking

# openflow.org/mininet

# The SDN Approach

Separate control from the datapath

   – i.e. separate policy from mechanism


**Datapath**: Define minimal network instruction set

   – A set of "plumbling primitives"

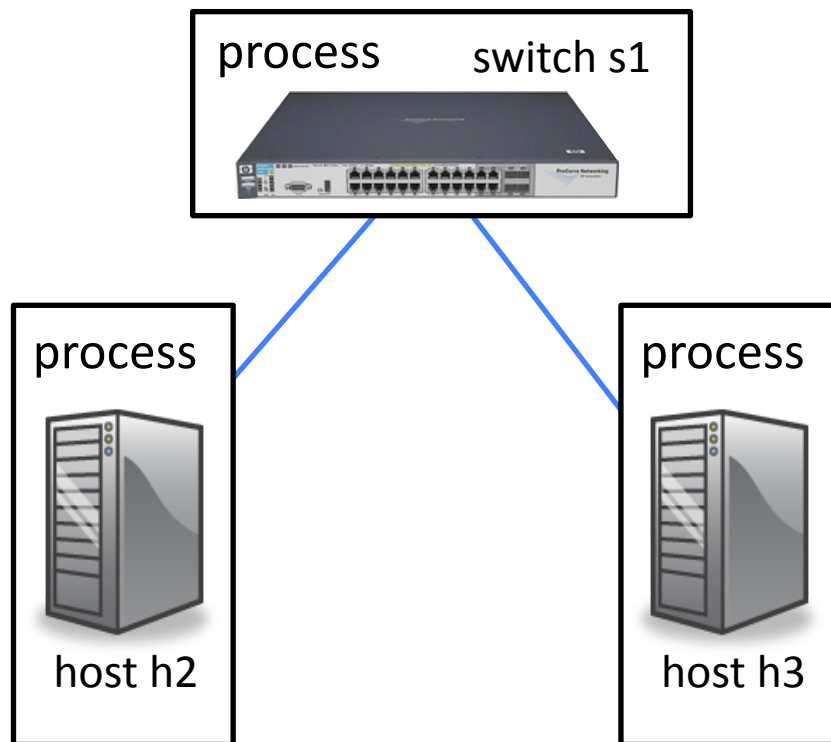   – A narrow interface: e.g. OpenFlow

**Control**: Define a network-wide OS

   – An API that others can develop on

# How to get performance fidelity?

- Careful process-to-core allocation
- Bandwidth limits
- Scheduling priorities
- Real-time scheduling
- Scheduling groups w/resource isolation

# Why not processes?

process     switch s1



process

host h2

process

host h3

+ scales better
-  breaks applications