

OS^v

Dor Laor, Avi Kivity
Cloudeus Systems



Glauber Costa

KVM, Containers, Xen



Nadav Har'EL,
Nested KVM

OS^V



Avi Kivity KVM
originator



Pekka Enberg,
kvm, jvm, slab



Dor Laor, Former kvm
project mngr

Or Cohen



Dmitry Fleytman



Ronen Narkis



Guy Zana



hch



The story so far

In the beginning there was hardware

... and then they added
an application

... and then they added
an operating system

... and then they added
a hypervisor

... and then they added
managed runtime

Notice the pattern?

Typical Cloud Stack

Your App

Application Server

JVM

Operating System

Hypervisor

Hardware

**Our software stack
Congealed into existence.**

A Historical Anomaly

Your App

Application Server

JVM

provides protection and abstraction

Operating System

provides protection and abstraction

Hypervisor

provides protection and abstraction

Hardware

Too Many Layers, Too Little Value

Property/Component	VMM	OS	runtime
Hardware abstraction	✓	✓	✓
Isolation	✓	✓	✓
Resource virtualization	✓	✓	✓
Backward compatibility	✓	✓	✓
Security	✓	✓	✓
Memory management	✓	✓	✓
I/O stack	✓	✓	
Configuration		✓	

Public Edition



Virtualization

Virtualization 1.0



Transformed the
enterprise from
physical2virtual

Virtualization 2.0



Compute node
~~≠~~
virtual server

Virtualization 2.0, Massive Scale

Scalability

Content="malware,Exec Code, Overflow, ExecCode Bypass"
nt="...0day. Gh0st RAT, Mac Control...
nt="Home Hwewer, Jenwe...
Data Day Texas @DataDayTexas
Wool! @aerospike just joined as a sponsor of Data Day Texas
bit.ly/VpctJ3 #d2x03 #nosql #bigdata
Expand Reply Retweeted Escorted More
John W. Smith
Update Shopping Cart
Remove from Cart
Submit
Wish List
HashTags
AEROSPIKE © 2012 Aerospoke. All rights reserved. Confidential Pg. 1

Virtualization 2.0, agility!

Deployments at Amazon.com

11.6

Seconds mean time
between deployments
(weekday)

1,079

Max number of
deployments in a
single hour

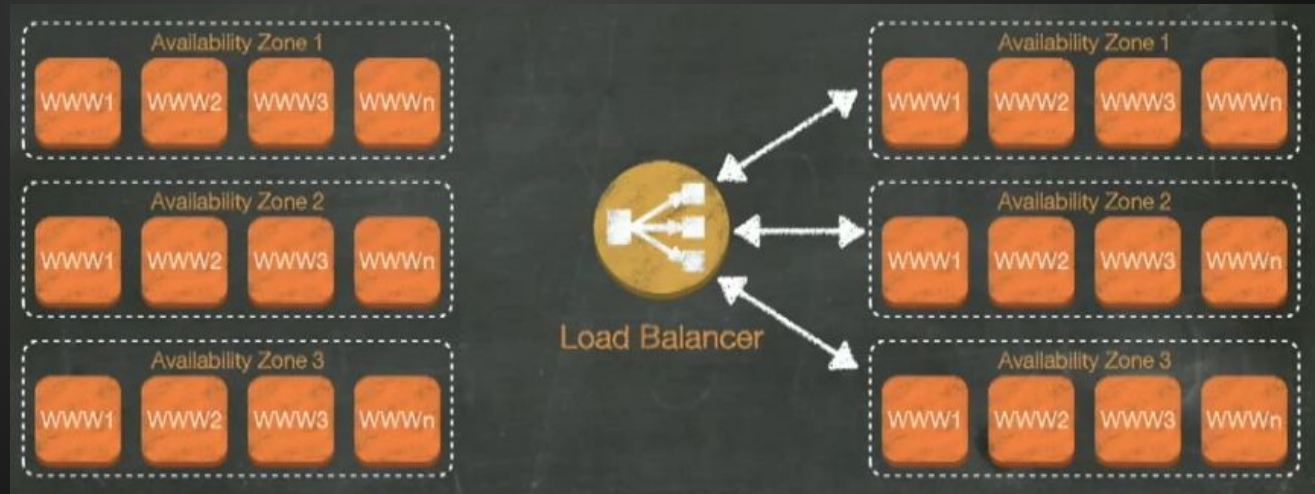
10,000

Mean number of hosts
simultaneously receiving
a deployment

30,000

Max number of hosts
simultaneously
receiving a deployment

**Rolling upgrade
within seconds and
a fall back option**



Virtualization 2.0

Architecture

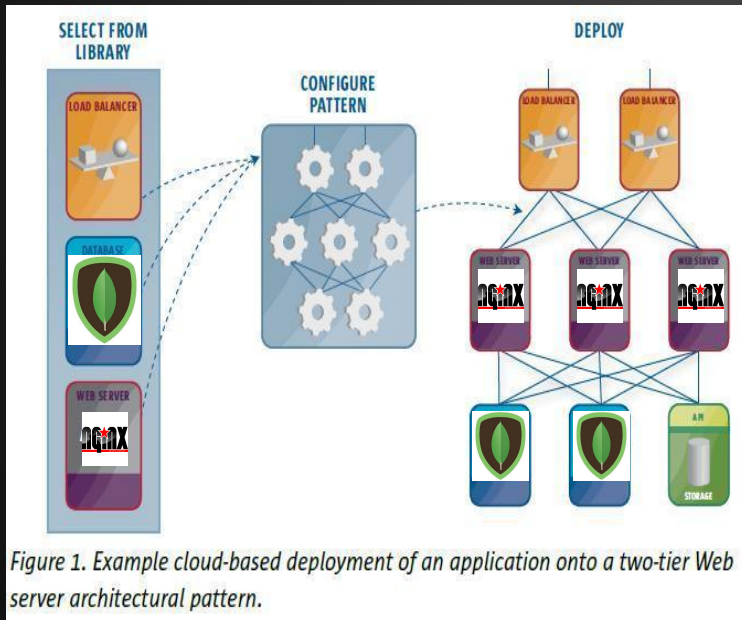


Figure 1. Example cloud-based deployment of an application onto a two-tier Web server architectural pattern.

vServer OS 1.0

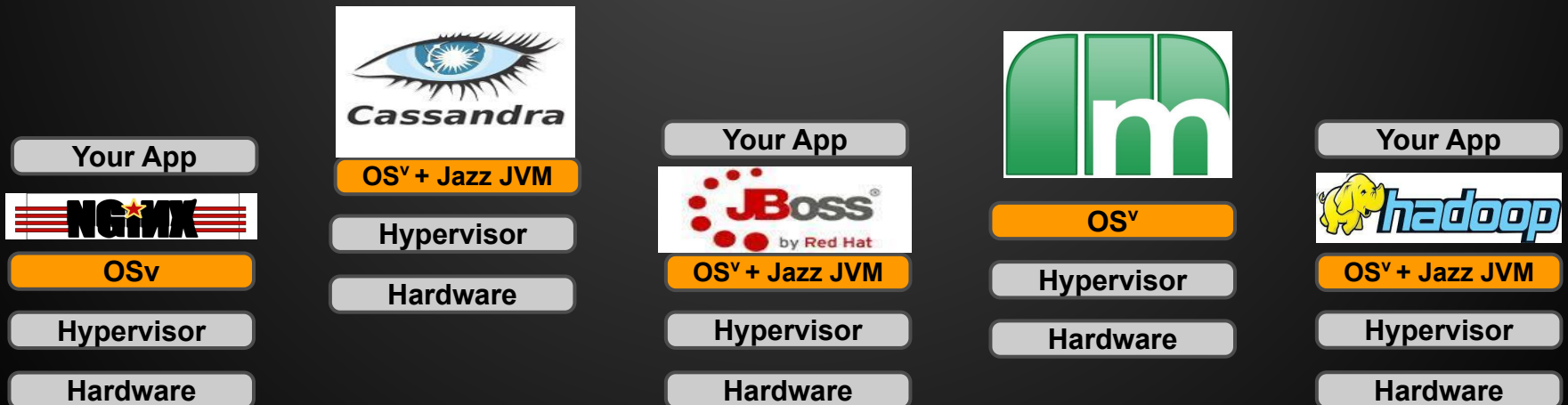
- No Hardware
 - No Users
 - No app(S)
-
- Yes Complexity



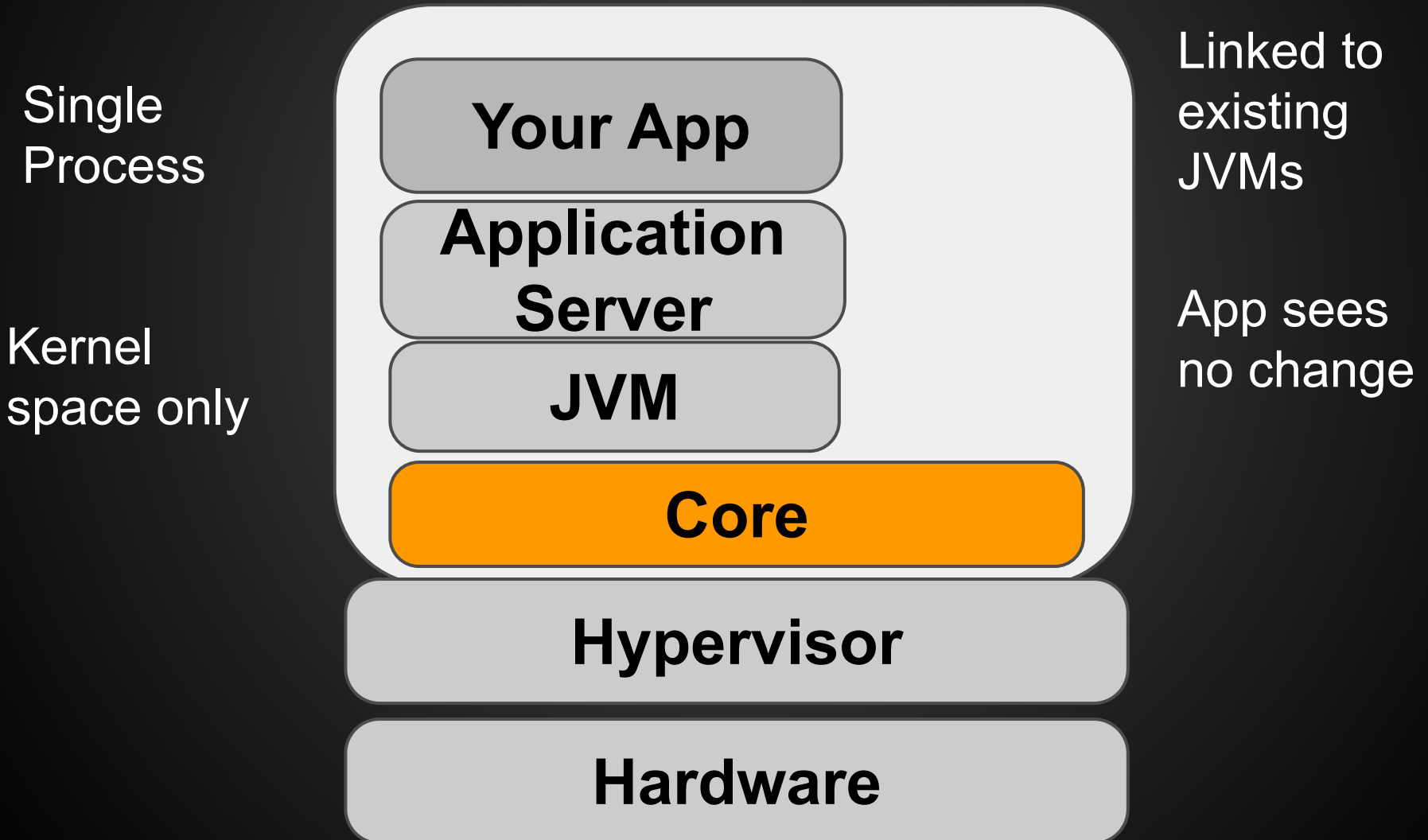
less is more.

Mission statement

Be the **best OS**
powering virtual machines
in the **cloud**



The new Cloud Stack - OS^v



The new Cloud Stack - OS^v

Memory

Huge pages, Heap vs Sys

I/O

Zero copy, full aio, batching

Scheduling

Lock free, low latency

Tuning

Out of the box, auto

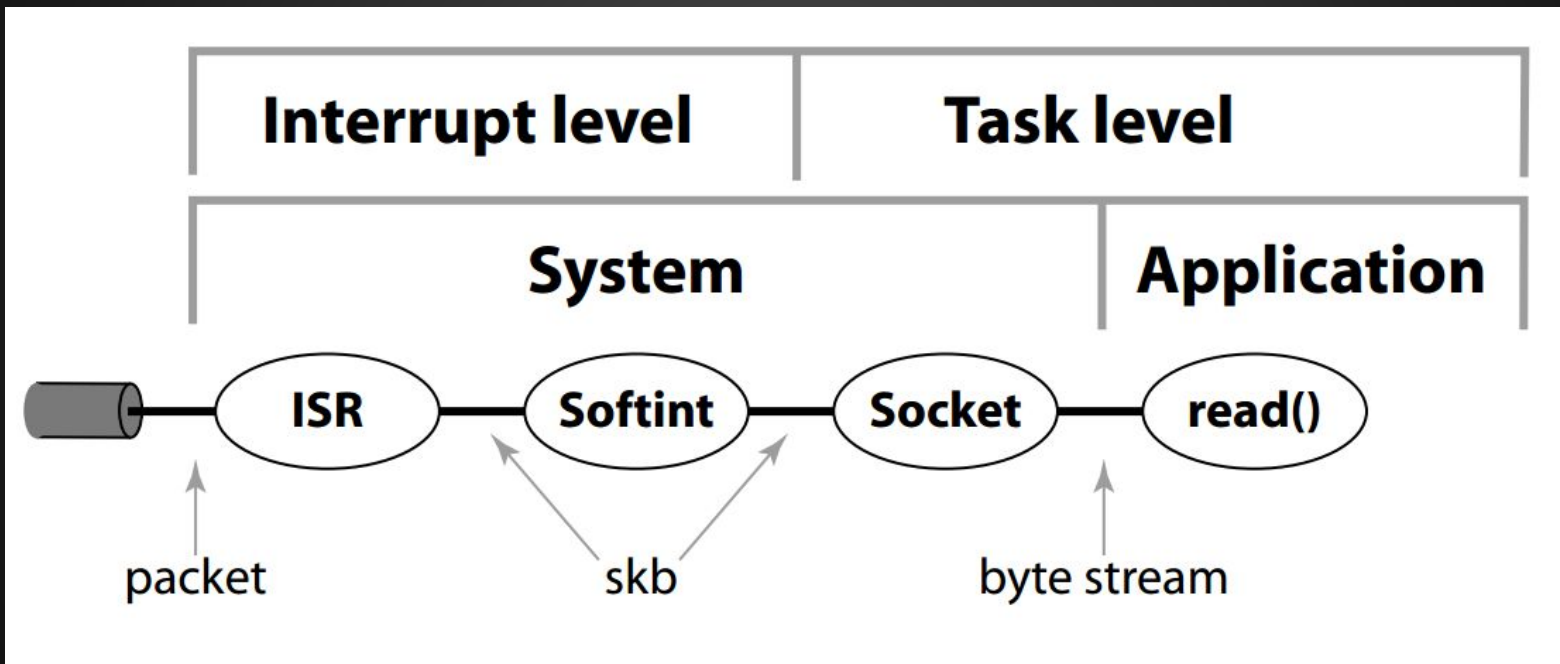
CPU

**Low cost ctx, Direct
signals,...**

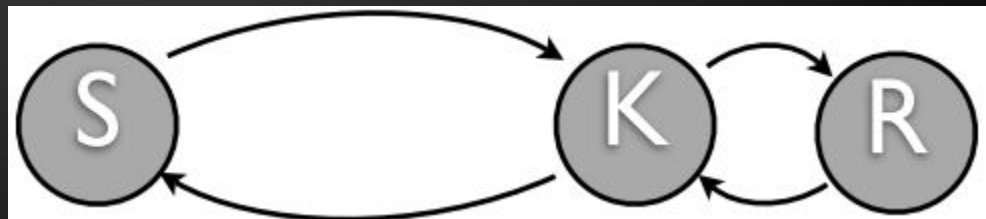
Van Jacobson == TCP/IP



Common kernel network stack



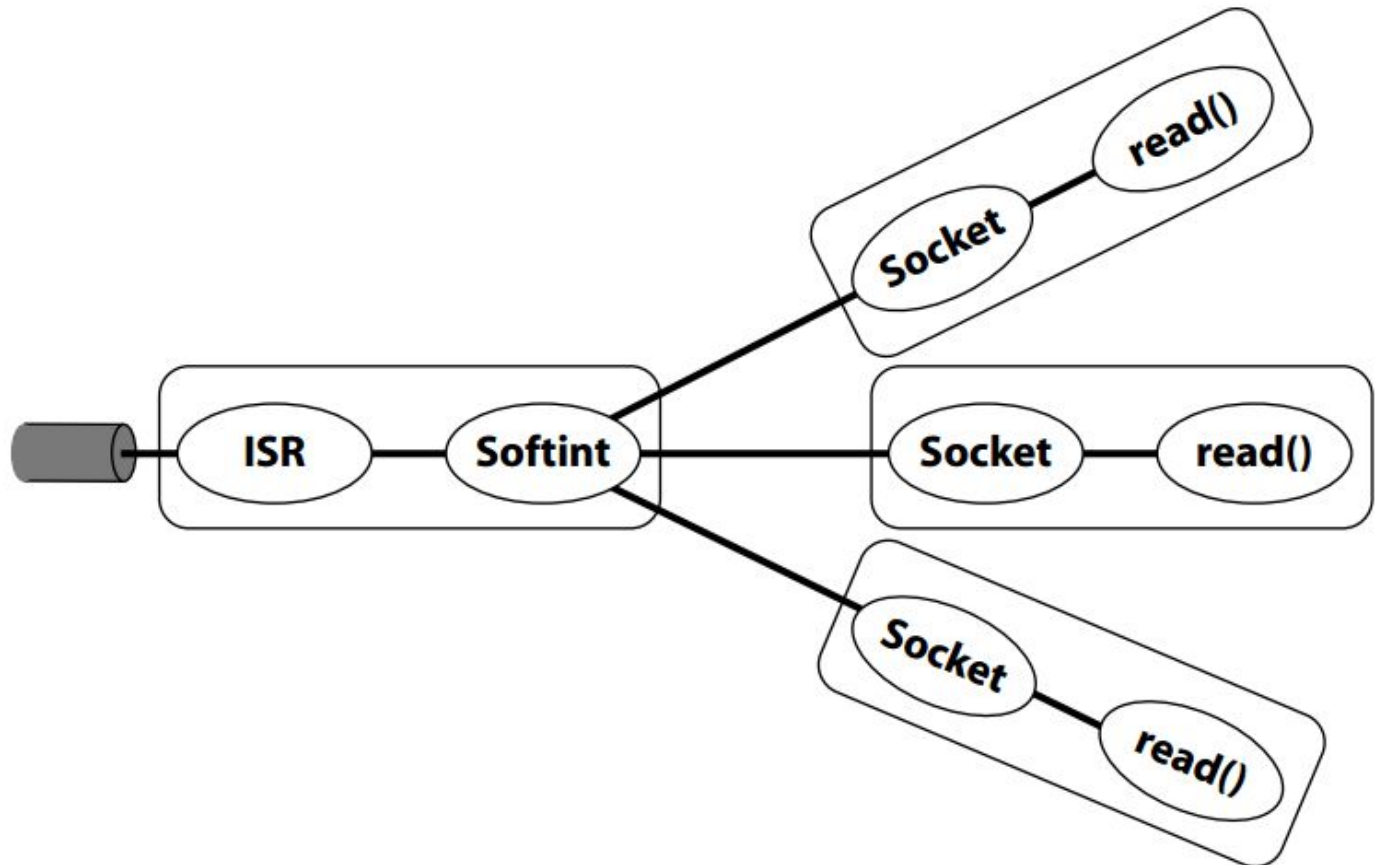
 Leads to servo-loop:



Van Jacobson == TCP/IP



Net Channel design:



Van Jacobson == TCP/IP



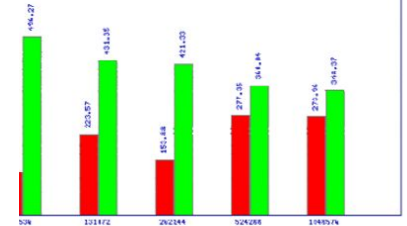
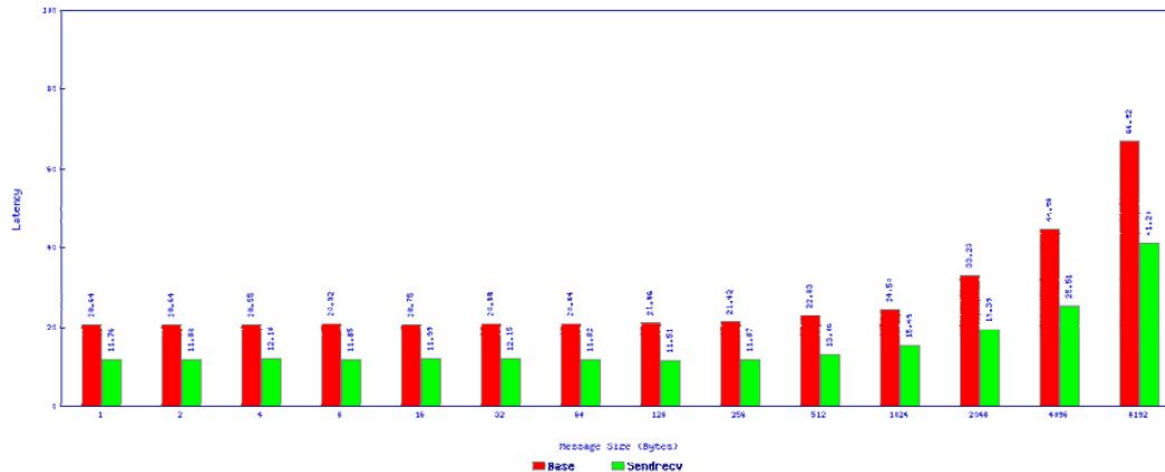
LAM MPI: Intel MPI Benchmark (IMB) using 4 boxes (8 processes) SendRecv bandwidth (bigger is better)

Intel Benchmark Absolute Bandwidth Comparison (Driver = e1000, Lib = POLL, Nodes = 4, Grid = 4(n) x 2(p))

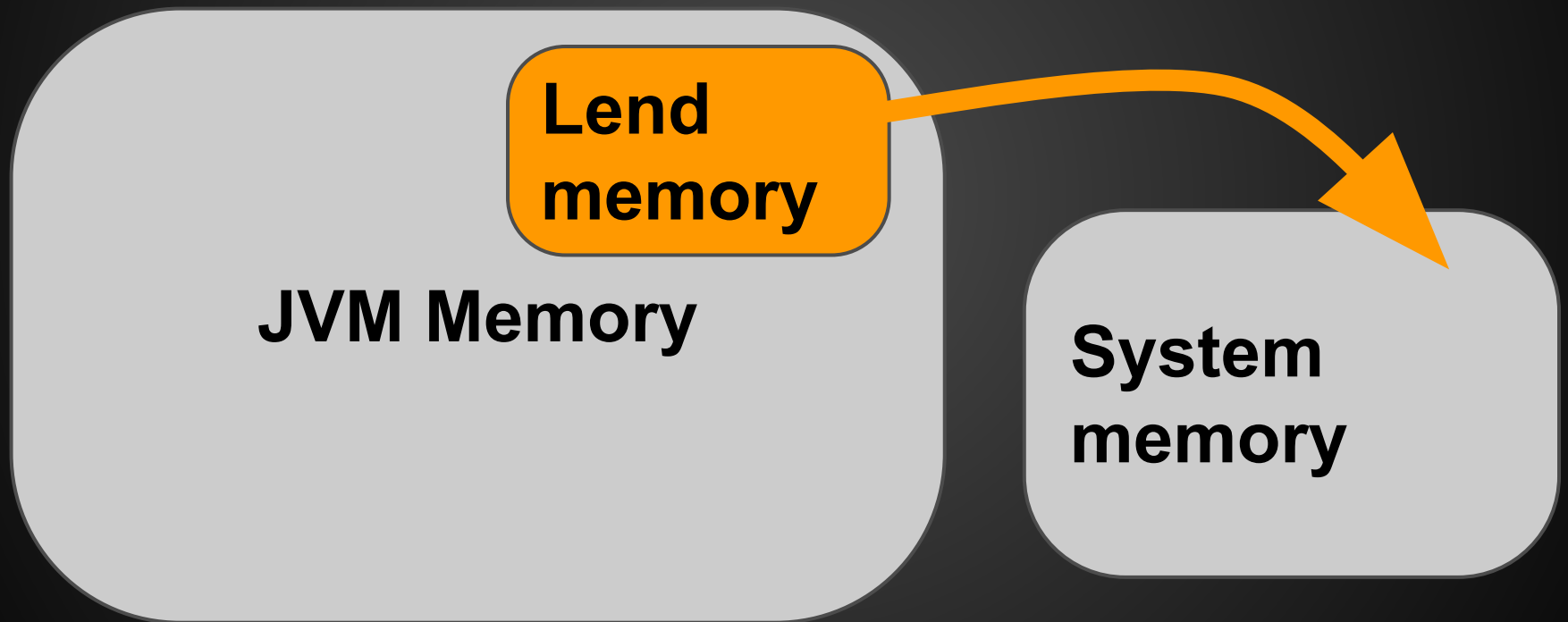


LAM MPI: Intel MPI Benchmark (IMB) using 4 boxes (8 processes) SendRecv Latency (smaller is better)

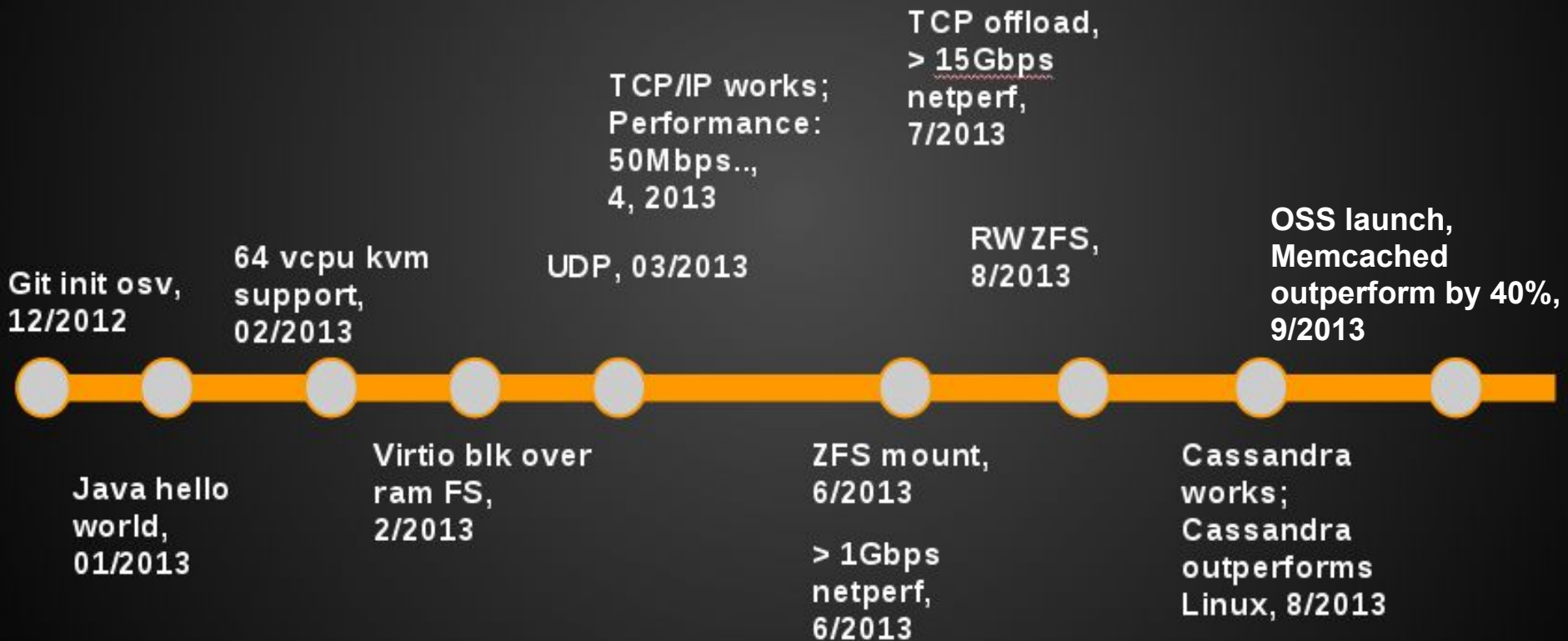
Intel Benchmark Absolute Latency Comparison (Driver = e1000, Lib = POLL, Nodes = 4, Grid = 4(n) x 2(p))



Dynamic heap, sharing is good



Milestones



Status

- Runs:
 - Java, C, JRuby, Scala, Groovy, Clojure, JavaScript
- Outperforms Linux:
 - SpecJVM, MemCacheD, Cassandra, TCP/IP
- 400% better w/ scheduler micro-benchmark
- < 1sec boot time
- ZFS filesystem
- Huge pages from the very beginning

Open Source

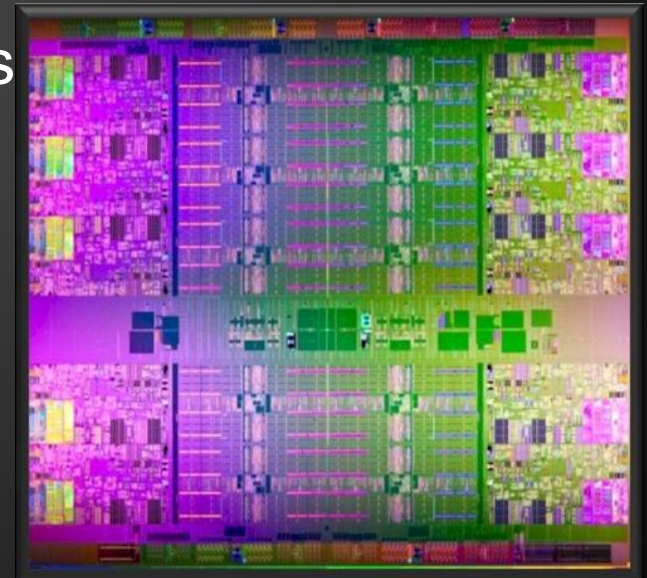
Fork me on GitHub

- These days, credibility == open source
- Looking for cooperation:
 - Kernel-level developers
 - Management stack
 - Dev/ops workflow
- BSD-style license



Architecture ports

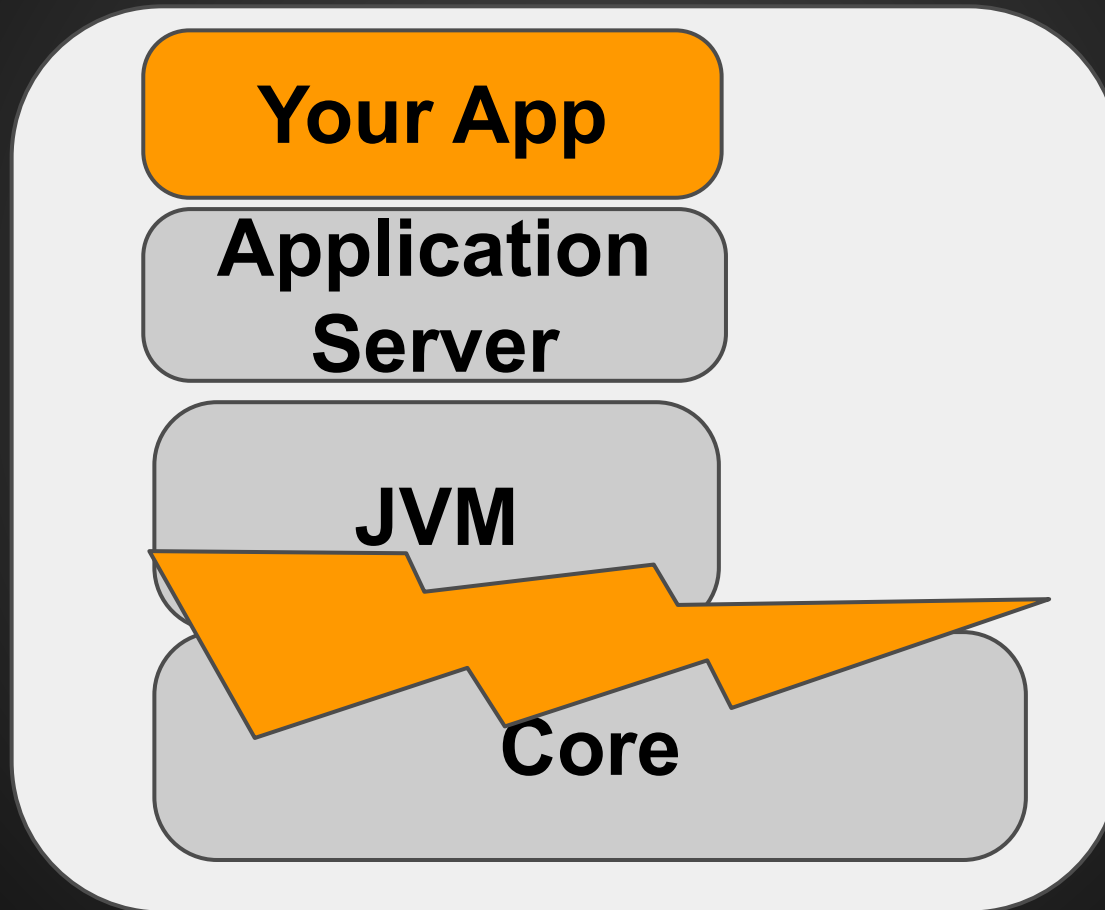
- 64-bit x86
 - KVM - running like a bat out of hell
 - Xen HVM - running (still slow :-)
 - Xen PV - in progress
 - VMware - planned in 2 months
- 64-bit ARM - planned
- Others - patches welcome



Integrating the JVM into the kernel

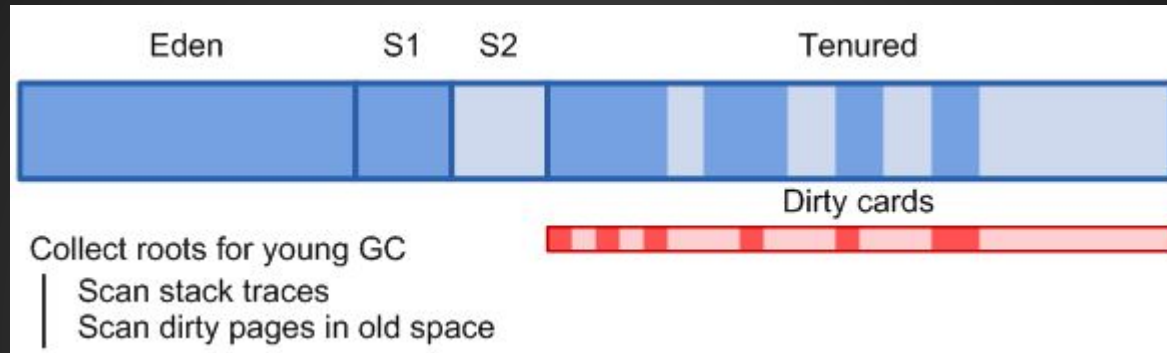
Dynamic
Heap
Memory

TCP in the
JVM + App
context



Fast inter
thread
wakeup

Integrating the JVM into the kernel



- G - Global
- D - Dirty
- A - Accessed
- C - Cache Disabled
- W - Write Through
- U - User/Supervisor
- R - Read/Write
- P - Present

Technical deep dive

- C++
- Idle time polling
- Performance and tracing
- Virtio-app

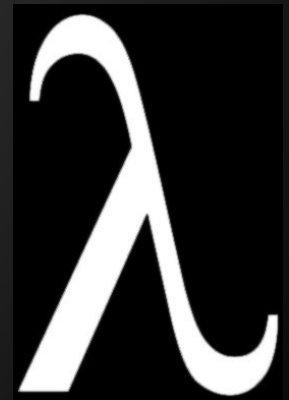


C++

```
int before(struct something *p)
{
    int r;

    r = -ENOENT;
    if (!p)
        goto out2;
    mutex_lock(&p->lock);
    r = -EINVAL;
    if (!p->y)
        goto out1;
    mutex_lock(&p->y->lock);
    r = ++p->y->n;
    mutex_unlock(&p->y->lock);
out1:
    mutex_unlock(&p->lock);
out2:
    return r;
}
```

```
int after(something* p)
{
    if (!p)
        return -ENOENT;
    WITH_LOCK(p->lock) {
        if (!p->y)
            return -EINVAL;
        WITH_LOCK(p->y->lock)
            return ++p->y->n;
    }
}
```



Idle-time polling

- Going idle is **much** more expensive on virtual machines
- So are inter-processor interrupts - IPIs
- Combine the two:
 - Before going idle, **announce** it via shared memory
 - **Delay** going idle
 - In the meanwhile, **poll** for wakeup requests from other processors
- Result: wakeups are faster, both for the processor waking, and for the wakee

Performance and tracing

```
TRACEPOINT(trace_mutex_lock, "%p", mutex *);  
TRACEPOINT(trace_mutex_lock_wait, "%p", mutex *);
```

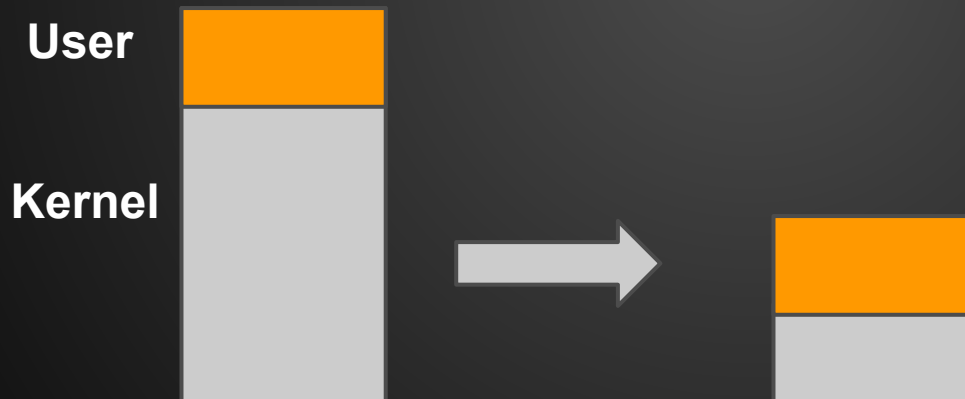
```
// ...
```

```
void mutex::lock()  
{  
    trace_mutex_lock(this);
```

```
[/]$ perf stat mutex_lock mutex_lock_wait sched_switch  
mutex_lock  mutex_lock_wait  sched_switch  
    11          0             2  
  885          0            181  
  154          0            152  
  154          0            154  
  404          0            190  
  222          0            157  
  150          0            152
```

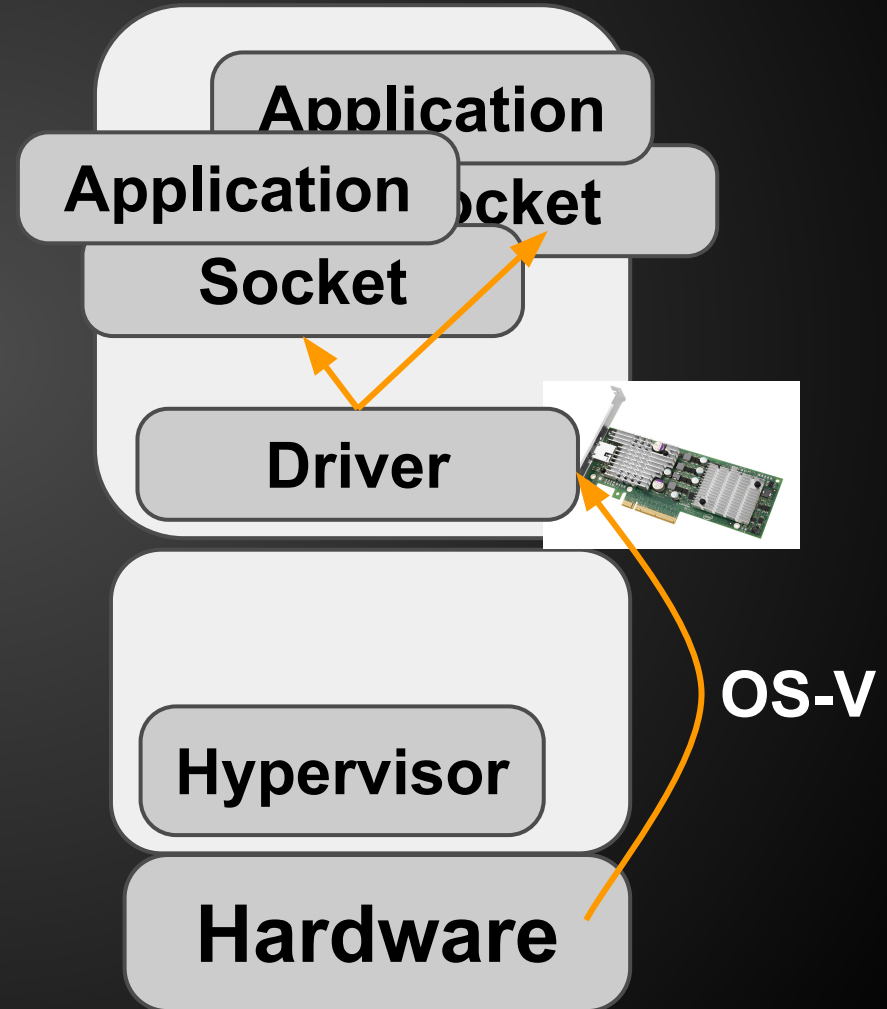
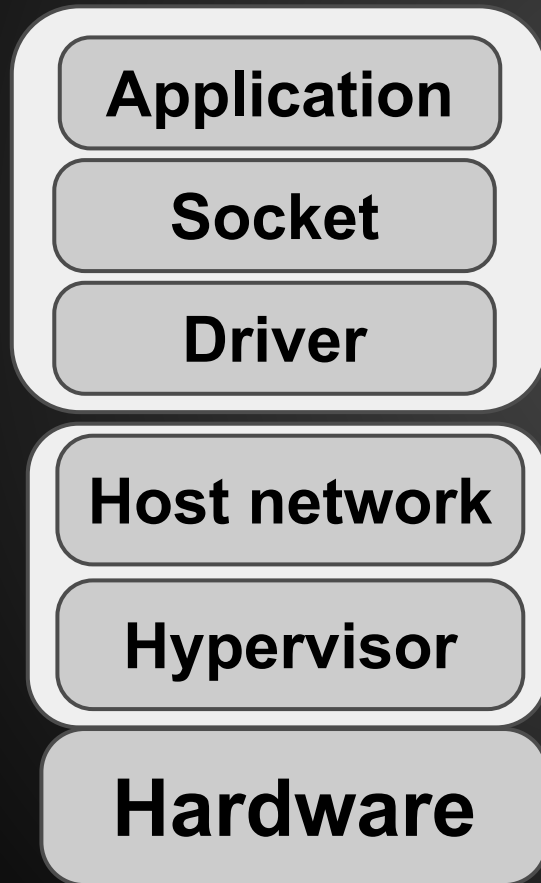
Virtio-app || Data plane

- For specialized applications, bypass the I/O stack completely
- Application consumes data from virtio rings



OS^v at the cutting edge front

Traditional



OS^v at the cutting edge front

- **Transactional Memory** (lock elision)
Better architecture match with higher transaction/sec and less contention
- Perfect match with **NVRam** abundance
In the near future we'll see NVRam reaches mainstream adoption. The importance of traditional filesystems will decrease, applications will manage their IO directly using NVRam

OS that doesn't get in the way

NO Tuning

NO State

NO Patching

X4 VMs per sys

admin ratio

Management

192.168.122.89:8080/upload

OSv Home **Deploy** Manage Monitor About Contact

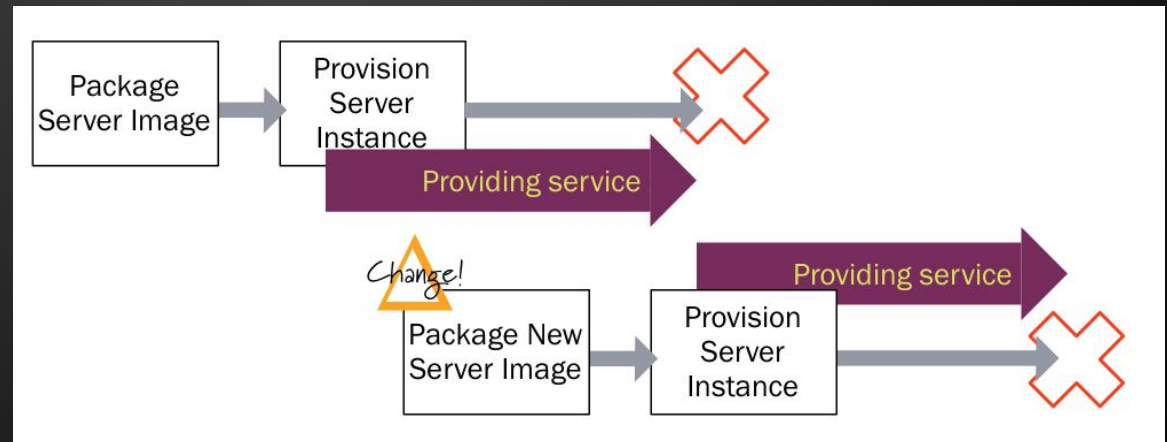
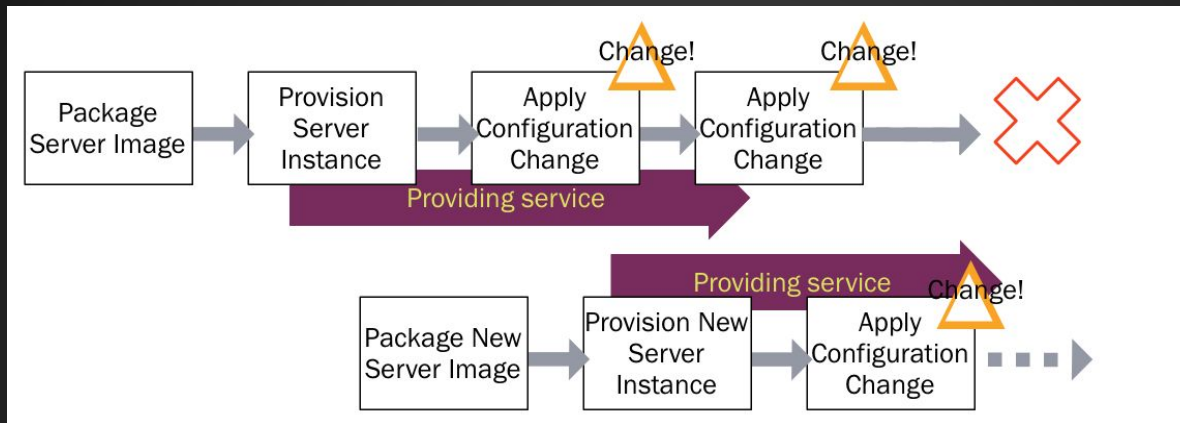
OSv application deployment

Deploy your Java applications into OSv by following these steps:

- Upload your application zip file (see [example](#) project).
- Activate the uploaded application by [starting](#) it.

+ Add files...
Choose Files No file chosen

Virtualization 2.0: Stateless servers



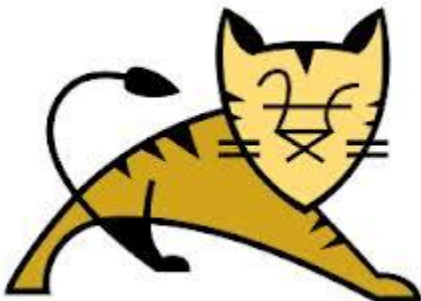
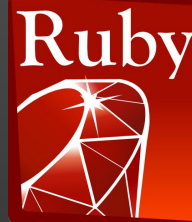
Let's Build A COMMUNITY



nodeJS



jetty://



zendServer
Community Edition

Porting a JVM application to OS^v

1. Done^{*}

* well, unless the application fork()^s

Porting a C application to OS^V

1. Must be a single-process application
2. May not fork() or exec()
3. Need to rebuild as a shared object (.so)
4. Other API limitations apply

Resources



<http://osv.io>



<https://github.com/cloudeius-systems/osv>



@CloudeiusSystems



osv-dev@googlegroups.com

OS^v@Clouddius

Clou dius Systems, OS Comparison

Feature/Property	OS ^v	Traditional OS
Good for:	Machete: Cloud/Virtualization	Swiss knife: anything goes
Typical workload	Single app * VMs	Multiple apps/users, utilities, anything
kernel vs app	Cooperation	distrust
API, compatibility	JVM, POSIX	Any, but versions/releases..
# Config files	0	1000
Tuning	Auto	Manual, requires certifications
Upgrade/state	Stateless, just boots	Complex, needs snapshots, hope..
JVM support	Tailored GC/STW solution	Yet another app
Lines of code	Few	Gazillion
License	BSD	GPL / proprietary