



KubeCon



CloudNativeCon

Europe 2020

*Virtual*

# Container Isolation via Virtualization: Don't Forget to Shrink the Guest

*Dan Williams, IBM & Hsuan-Chi (Austin) Kuo, UIUC*

# Virtualization + containers = security?



[LEARN](#)

[SOFTWARE](#)

[DOCS](#)

[COMMUNITY](#)

[SUPPORTERS](#)

[BLOG](#)

## The speed of containers, the security of VMs

Kata Containers is an open source container runtime, building lightweight virtual machines that seamlessly plug into the containers ecosystem.

[GET THE LATEST RELEASE](#) >



# Virtualization + containers = security?

The screenshot shows the Firecracker website. At the top left is the 'katacontainers' logo. To its right are navigation links for 'LEARN' and 'SOFTWARE'. Further right is the 'Firecracker' logo, followed by 'BENEFITS', 'HOW IT WORKS', and 'FAQS'. The main content area features a dark blue background with a grid pattern. The headline reads 'The speed and security of containers'. Below this, it states 'Kata Containers is an open source... seamlessly plug into the containers ecosystem.' An orange button at the bottom says 'GET THE LATEST RELEASE >'. An orange overlay on the right side of the image contains the text 'Secure and fast microVMs for serverless computing'.

# Virtualization + containers = security?



LEARN

SOFTWARE



Firecracker

BENEFITS

HOW IT WORKS

FAQS

## Weave Ignite

Weave Ignite is an open source Virtual Machine (VM) manager with a container UX and built-in GitOps management.

- Combines [Firecracker MicroVMs](#) with Docker / [OCI images](#) to unify containers and VMs.
- Works in a [GitOps](#) fashion and can manage VMs declaratively and automatically like Kubernetes and Terraform.

Ignite is fast and secure because of Firecracker. Firecracker is an [open source KVM implementation](#) from AWS that is optimised for [high security](#)



and fast microVMs  
erless computing

system.



# Virtualization + containers = security?

The screenshot shows a web page with a dark blue header on the left containing the 'katacontainers' logo and navigation links for 'LEARN' and 'SOFTWARE'. A large orange banner on the right features the 'Firecracker' logo and the text 'and fast microVMs'. Below the banner, a green button says 'SIGN UP / LOGIN'. The main content area has a white background with a blue 'InfoQ' logo and a navigation menu with categories: Streaming, Machine Learning, Reactive, Microservices, Containers, Observability, and Security. The article title is 'Containers In 2019: They're Calling It A [Hypervisor] Comeback' with a 'DEVOPS' tag. The article text includes: 'Weave Ignite is an open source container manager with a container U management.', 'Combines Firecracker and OCI images to unify container management.', 'Works in a GitOps fashion declaratively and automatically and Terraform.', and 'Ignite is fast and secure because Firecracker is an open source project from AWS that is optimised for security'.

**Weave Ignite**

Weave Ignite is an open source container manager with a container U management.

- Combines [Firecracker](#) and [OCI images](#) to unify container management.
- Works in a [GitOps](#) fashion declaratively and automatically and Terraform.

Ignite is fast and secure because Firecracker is an [open source](#) project from AWS that is optimised for security

**InfoQ**

Streaming Machine Learning Reactive Microservices Containers Observability Security

InfoQ Homepage > Articles > Containers In 2019: They're Calling It A [Hypervisor] Comeback

**DEVOPS**

## Containers in 2019: They're Calling it a [Hypervisor] Comeback

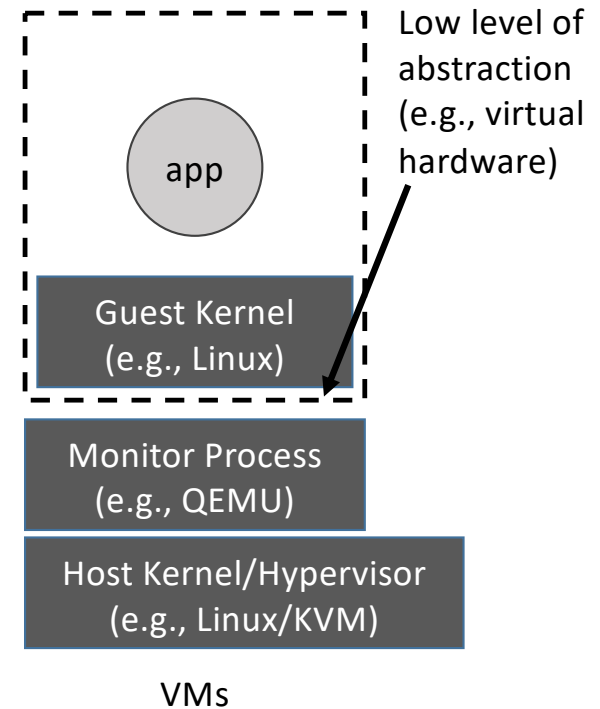
# But wait? Aren't VMs slow and heavyweight?



- Boot time?
- Memory footprint?
- Especially for environments like serverless??!!

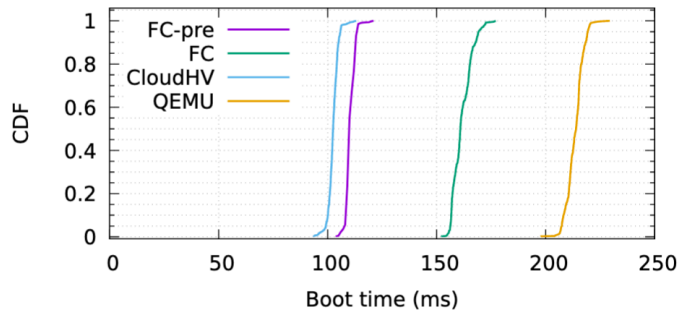
# VMs are becoming lightweight

- Thin monitors
  - e.g., AWS Firecracker
  - Reduce complexity for performance (e.g., no PCI)

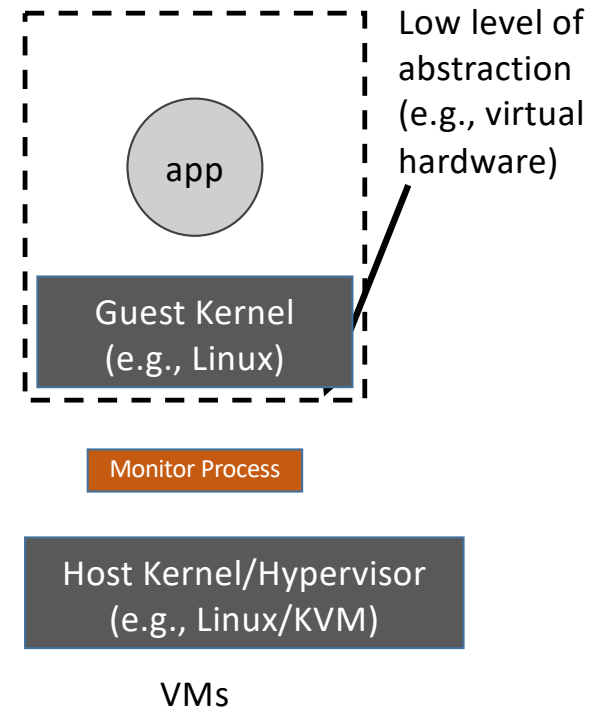


# VMs are becoming lightweight

- Thin monitors
  - e.g., AWS Firecracker
  - Reduce complexity for performance (e.g., no PCI)



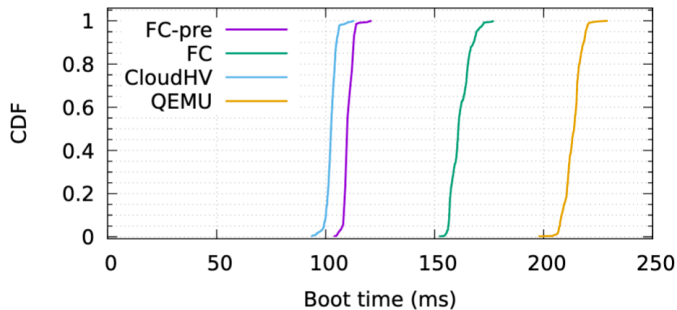
Firecracker boot times as reported in Agache et al., NSDI 2020





# VMs are becoming lightweight

- Thin monitors
  - e.g., AWS Firecracker
  - Reduce complexity for performance (e.g., no PCI)



Firecracker boot times as reported in Agache et al., NSDI 2020

## My VM is Lighter (and Safer) than your Container

**Filipe Manco**  
 NEC Laboratories Europe  
 filipe.manco@gmail.com  
**Costin Lupu**  
 Univ. Politehnica of Bucharest  
 costin.lupu@cs.pub.ro  
**Florian Schmidt**  
 NEC Laboratories Europe  
 florian.schmidt@nec-lab.eu  
**Jose Mendes**  
 NEC Laboratories Europe  
 jose.mendes@nec-lab.eu  
**Simon Kuenzer**  
 NEC Laboratories Europe  
 simon.kuenzer@nec-lab.eu  
**Sumit Sati**  
 NEC Laboratories Europe  
 sativishy@gmail.com  
**Kenichi Yasukata**  
 NEC Laboratories Europe  
 kenichi.yasukata@nec-lab.eu  
**Costin Raicu**  
 Univ. Politehnica of Bucharest  
 costin.raicu@cs.pub.ro  
**Felipe Huici**  
 NEC Laboratories Europe  
 felipe.huici@nec-lab.eu

**ABSTRACT**  
 Containers are in great demand because they are lightweight when compared to virtual machines. On the downside, containers offer weaker isolation than VMs, to the point where people run containers in virtual machines to achieve proper isolation. In this paper, we examine whether there is indeed a strict tradeoff between isolation (VMs) and efficiency (containers). We find that VMs can be as amiable as containers, as long as they are small and the toolstack is fast enough. We achieve lightweight VMs by using unikernels for specialized applications and with TinyT, a tool that enables creating tailor-made, trimmed-down Linux virtual machines. By themselves, lightweight virtual machines are not enough to ensure good performance since the virtualization control plane (the toolstack) becomes the performance bottleneck. We present LightVM, a new virtualization solution based on Xen that is optimized to offer fast boot-times regardless of the number of active VMs. LightVM features a complete redesign of Xen's control-plane, transforming its centralized operation to a distributed one where interactions with the hypervisor are reduced to a minimum. LightVM can boot a VM in 2.3ms, comparable to `forkexec` on Linux (1ms), and two orders of magnitude faster than Docker. LightVM can pack thousands of LightVM guests on modest hardware with memory and CPU usage comparable to that of processes.

**CCS CONCEPTS**  
 • Software and its engineering → Virtual machines; Operating systems.

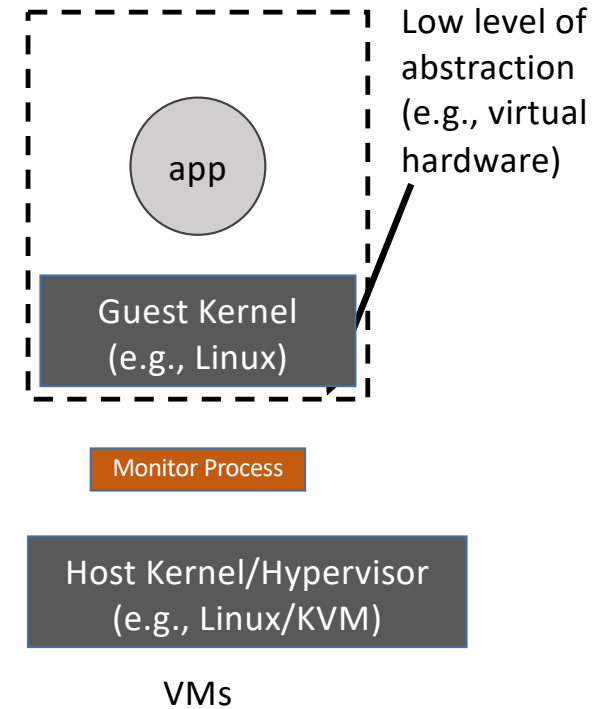
**KEYWORDS**  
 Virtualization, unikernels, specialization, operating systems, Xen, containers, hypervisor, virtual machine.

**ACM Reference Format:**  
 Filipe Manco, Costin Lupu, Florian Schmidt, Jose Mendes, Simon Kuenzer, Sumit Sati, Kenichi Yasukata, Costin Raicu, and Felipe Huici. 2017. My VM is Lighter (and Safer) than your Container. In *Proceedings of SOSP '17: ACM SIGOPS 20th Symposium on Operating Systems Principles*, Shanghai, China, October 28, 2017. SOSP '17, 14 pages.  
<https://doi.org/10.1145/312747.312783>

**1 INTRODUCTION**  
 Lightweight virtualization technologies such as Docker [6] and LXC [15] are gaining enormous traction. Google, for instance, is reported to run all of its services in containers [4], and Container as a Service (CaaS) products are available from a number of major players including Azure's Container Service [32], Amazon's EC2 Container Service and Lambda offerings [1, 2], and Google's Container Engine service [10]. Beyond these services, lightweight virtualization is crucial in a wide range of use cases, including just-in-time instantiation of services [23, 26] (e.g., fibers against DDoS attacks, TCP acceleration proxies, content caches, etc.) and NFV [6, 11] in which *resilience, confidentiality, cost, and performance*.

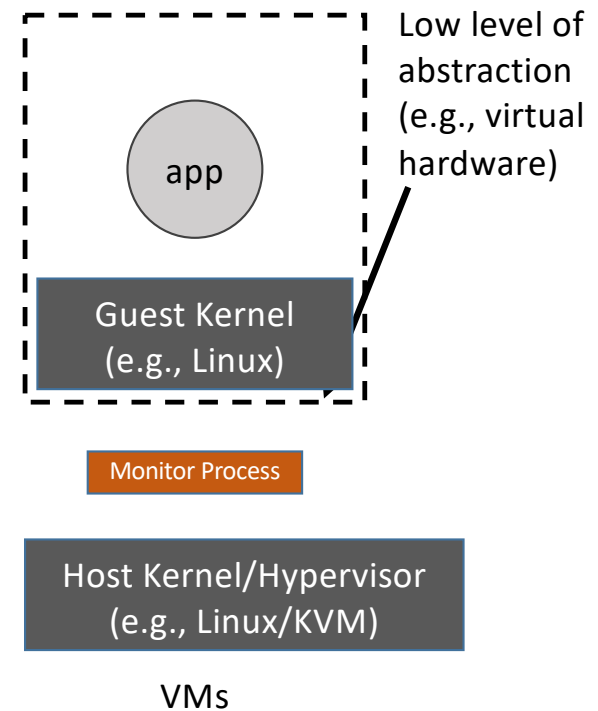
Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted by ACM, provided that the fee code of each article is paid to ACM. For more information, contact [permissions@acm.org](mailto:permissions@acm.org).

Manco et al., SOSP 2017



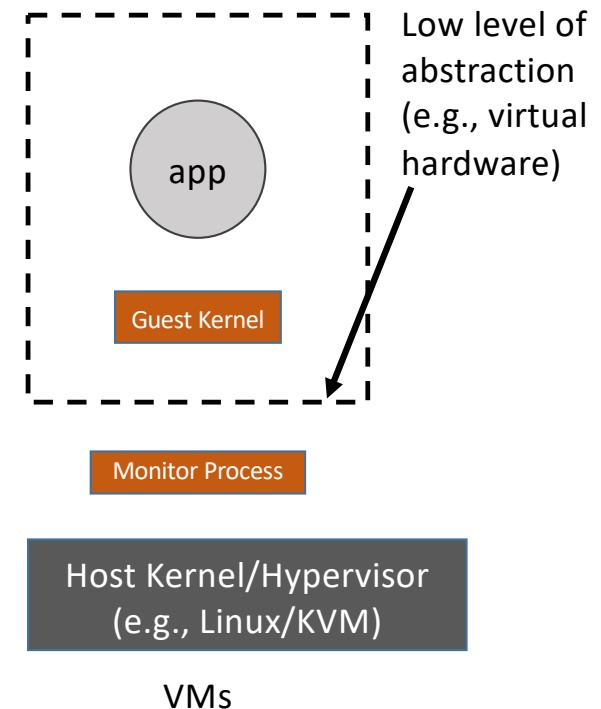
# VMs are becoming lightweight

- Thin monitors
  - e.g., AWS Firecracker
  - Reduce complexity for performance (e.g., no PCI)
- What about thin guests?



# VMs are becoming lightweight

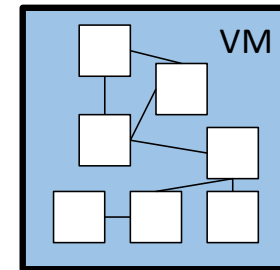
- Thin monitors
  - e.g., AWS Firecracker
  - Reduce complexity for performance (e.g., no PCI)
- What about thin guests?
  - Userspace: (e.g., Ubuntu --> Alpine Linux)
  - Kernel configuration (e.g., TinyX)
  - **How thin can you go?**



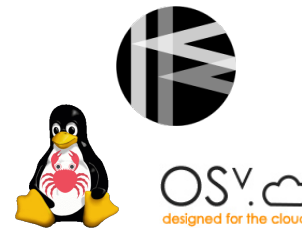
# Unikernels are thin guests to the extreme

- An application linked with **library OS** components
- Run on **virtual hardware** (like) abstraction
- Single CPU

- Language-specific
  - MirageOS (OCaml)
  - IncludeOS (C++)



- Legacy-oriented
    - Rumprun (NetBSD-based)
    - Hermitux
    - OSv
- } Claim binary compatibility with Linux



# Unikernels are great

- Small kernel size
- Fast boot time
- Performance
- Security

# Unikernels are great... but

- Small kernel size
- Fast boot time
- Performance
- Security
- Lack full Linux support
- Hermitux: supports only 97 system calls
- OSv:
  - application needs to be compiled with `-PIE`, can't use TLS
  - Static-linked applications are not supported
  - `Fork()` , `execve()` are not supported
  - Special files are not supported such as `/proc`
  - Signal mechanism is not complete
- Rumprun: only 37 curated applications
- Community is too small to keep it rolling



Lupine Linux  
"Unikernel"

Can Linux

> **be as small as**

> **boot as fast as**

> **outperform**

unikernels?



Lupine Linux  
"Unikernel"

Can Linux

> **be as small as**

> **boot as fast as**

> **outperform**

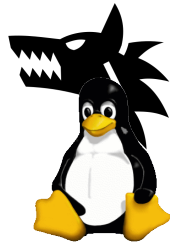
**unikernels?**

- Spoiler alert: Yes!
  - 4MB image size
  - 23 ms boot time
  - Up to 33% higher throughput



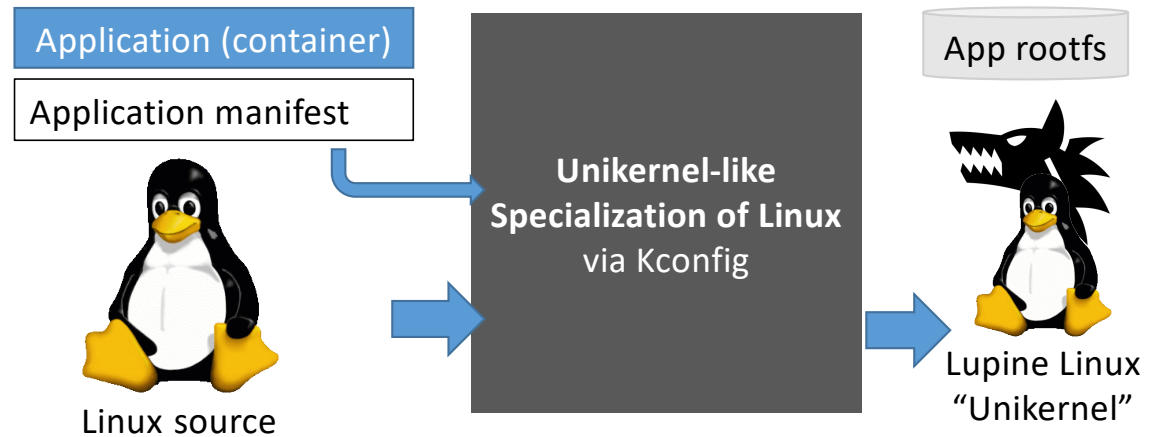
# Segue to Austing talking about...

- Lupine Linux



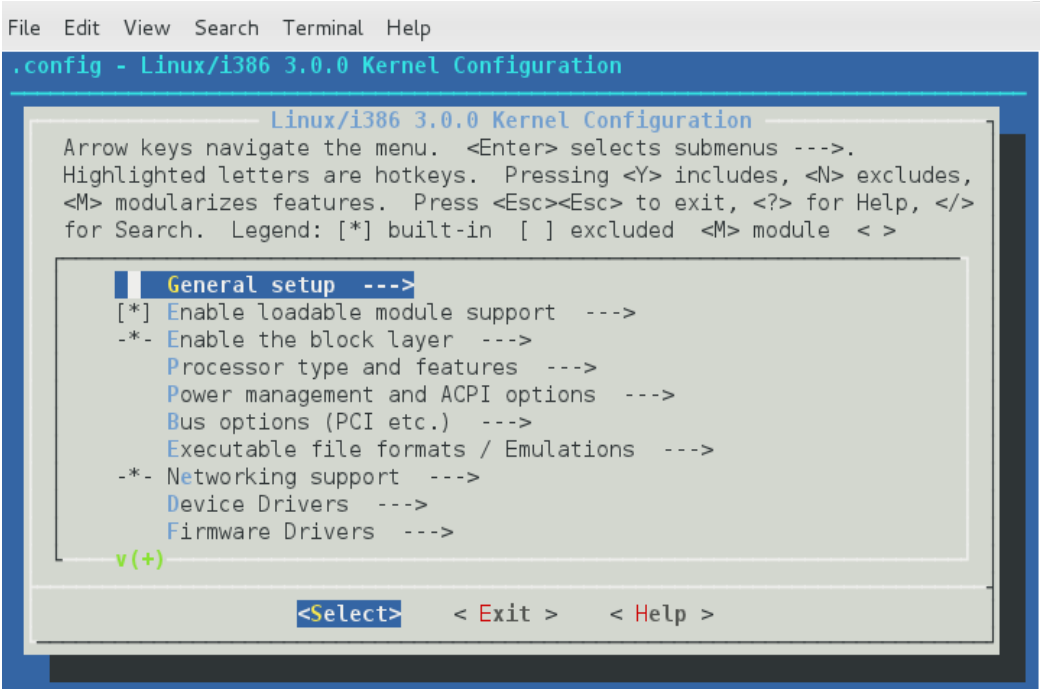
# Lupine Linux Overview and Roadmap

- Introduction
- Lupine Linux
- Evaluation
- Related Work



# Unikernels are all about specialization

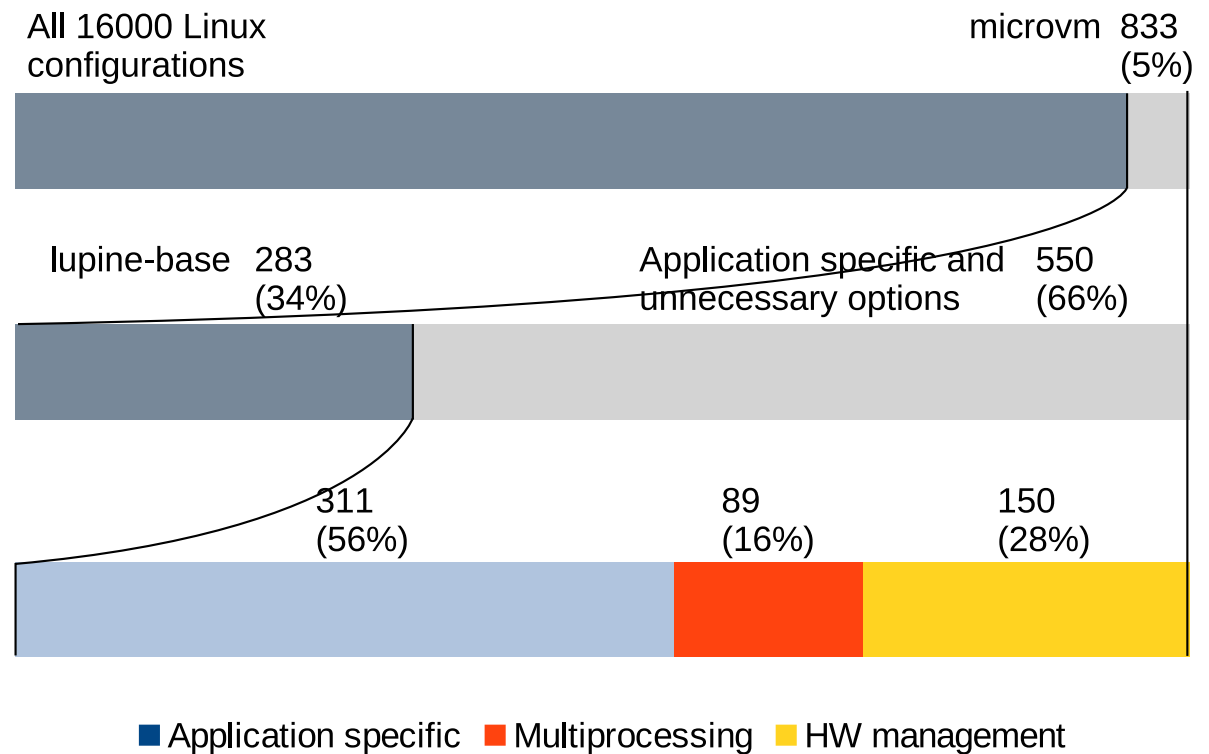
- Unikernels include only what is needed
- Linux is very configurable
  - Kconfig
  - 16,000 options
    - Drivers
    - Filesystems
    - Processor features
    - ...



The image shows a terminal window titled ".config - Linux/i386 3.0.0 Kernel Configuration". The window displays the "Linux/i386 3.0.0 Kernel Configuration" menu. At the top, it provides instructions: "Arrow keys navigate the menu. <Enter> selects submenus --->. Highlighted letters are hotkeys. Pressing <Y> includes, <N> excludes, <M> modularizes features. Press <Esc><Esc> to exit, <?> for Help, </> for Search. Legend: [\*] built-in [ ] excluded <M> module < >". The main menu is titled "General setup --->" and lists several options: "[\*] Enable loadable module support --->", "-\*- Enable the block layer --->", "Processor type and features --->", "Power management and ACPI options --->", "Bus options (PCI etc.) --->", "Executable file formats / Emulations --->", "-\*- Networking support --->", "Device Drivers --->", and "Firmware Drivers --->". A green cursor "v(+)" is positioned at the bottom left of the menu. At the bottom of the window, there are three buttons: "<Select>", "< Exit >", and "< Help >".

# Specializing Linux through configuration

- Start with Firecracker MicroVM configuration
- Can we remove even more?
  - Application-specific options
  - Multiprocessing
  - HW management



# Specializing for lightweight VMs

- Do we need support for multiple trust domains?
  - Related to isolating, accounting for processes
    - Cgroups, namespaces, SELinux, seccomp, KPTI
  - SMP, NUMA
  - Module support
- Do we need support for general hardware?
  - Intended to run as VMs in the cloud
  - MicroVM removes many drivers and arch-specific configs
  - Lupine removes more, including power mgmt

# Application-specific options

- Example: system calls

- Kernel services

- e.g., /proc, sysctl

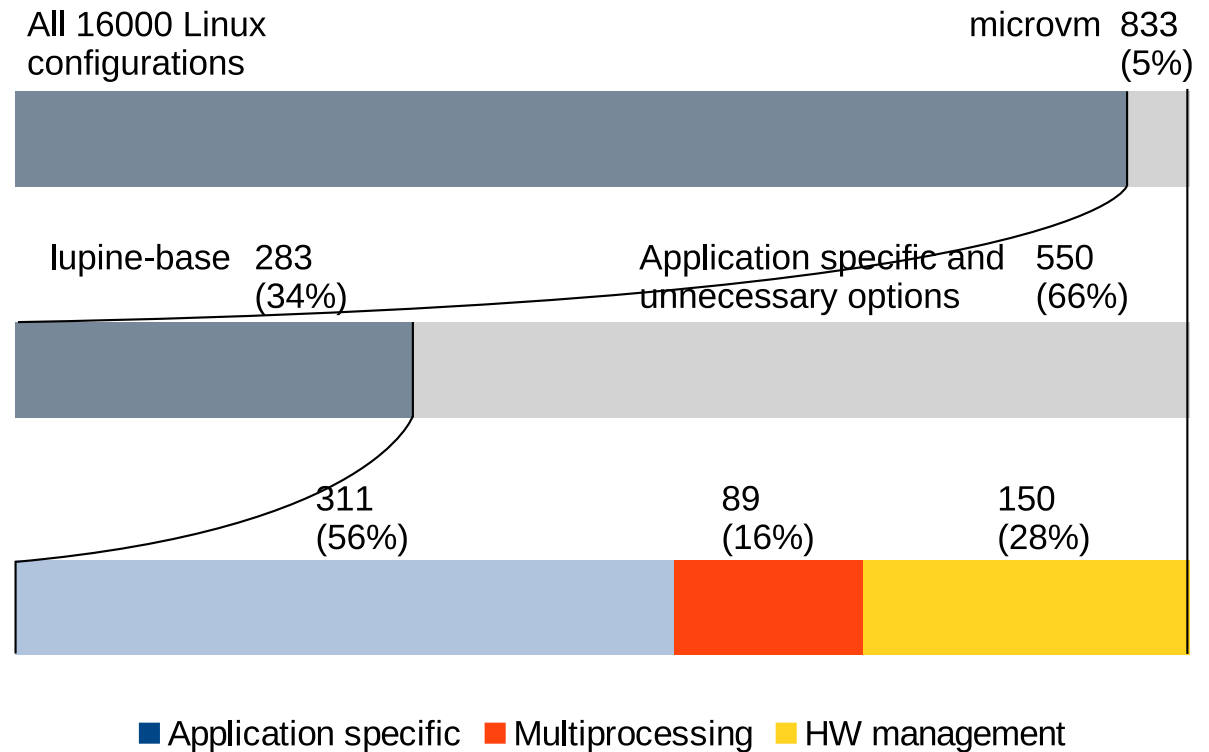
- Kernel library

- Crypto routines
- Compression routines

Option	Enabled System Call(s)
ADVISE_SYSCALLS	madvise, fadvise64
AIO	io_setup, io_destroy, io_submit, io_cancel, io_getevents
BPF_SYSCALL	bpf
EPOLL	epoll_ctl, epoll_create, epoll_wait, epoll_pwait
EVENTFD	eventfd, eventfd2
FANOTIFY	fanotify_init, fanotify_mark
FHANDLE	open_by_handle_at, name_to_handle_at
FILE_LOCKING	flock
FUTEX	futex, set_robust_list, get_robust_list
INOTIFY_USER	inotify_init, inotify_add_watch, inotify_rm_watch
SIGNALFD	signalfd, signalfd4
TIMERFD	timerfd_create, timerfd_gettime, timerfd_settime

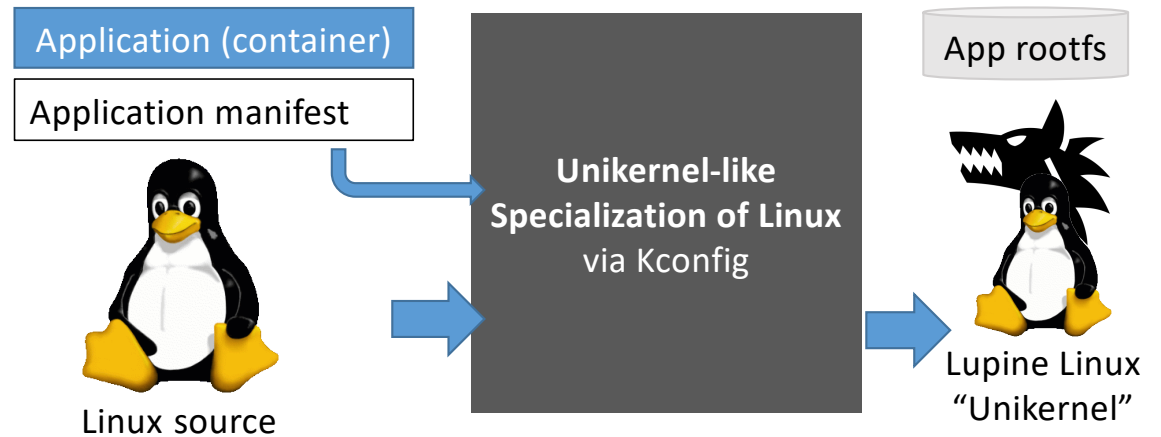
# How to get an app-specific kernel config

- Start with lupine-base
- Manual trial and error
  - Guided by application output
  - E.g., *the futex facility returned an unexpected error code*  
=> CONFIG\_FUTEX
- In general, this is a hard problem



# Lupine Linux Overview and Roadmap

- Introduction
- Lupine Linux
- Evaluation
- Related Work



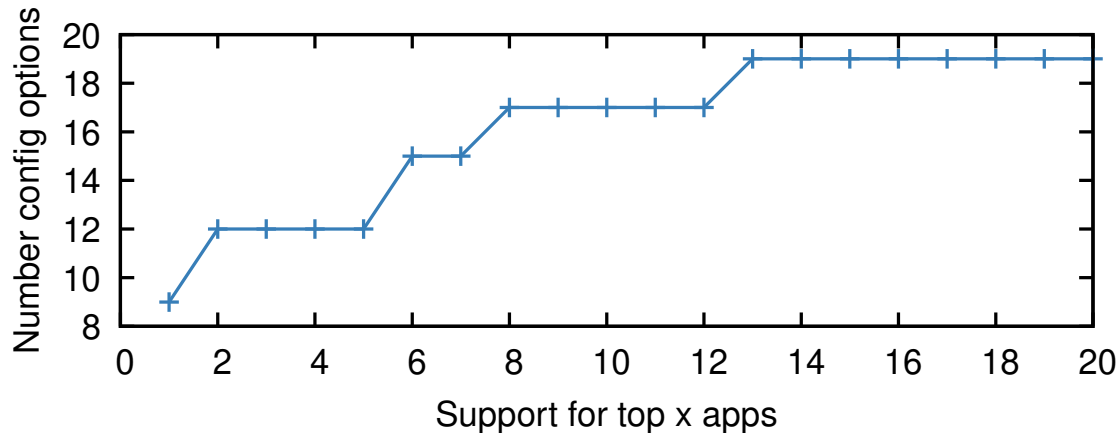


# Evaluation setup

- Machine setup
  - CPU: Intel(R) Xeon(R) CPU E3-1270 v6 @ 3.80GHz
  - Mem: 16 GB
- VM setup
  - Hypervisor : firecracker
  - 1 VCPU, 512 MB Mem
  - Guest: Linux 4.0

# Configuration Diversity

- Manually determined app-specific configurations
- 20 top apps on Docker hub (83% of all downloads)
- Only 19 configuration options required to run all 20 applications: *lupine-general*



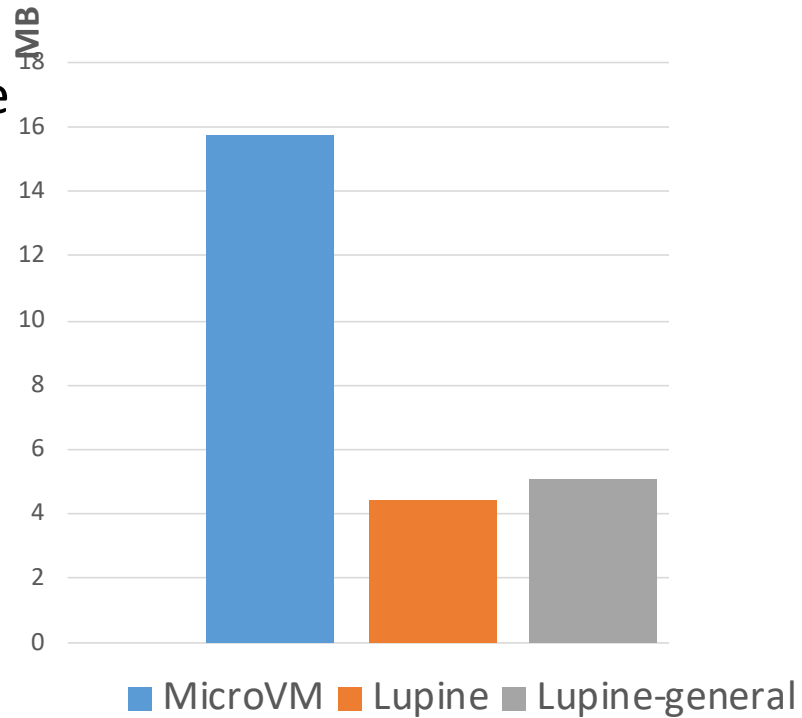
Name	Downloads	Description	# Options atop <i>lupine-base</i>
nginx	1.7	Web server	13
postgres	1.6	Database	10
httpd	1.4	Web server	13
node	1.2	Language runtime	5
redis	1.2	Key-value store	10
mongo	1.2	NOSQL database	11
mysql	1.2	Database	9
traefik	1.1	Edge router	8
memcached	0.9	Key-value store	10
hello-world	0.9	C program "hello"	0
mariadb	0.8	Database	13
golang	0.6	Language runtime	0
python	0.5	Language runtime	0
openjdk	0.5	Language runtime	0
rabbitmq	0.5	Message broker	12
php	0.4	Language runtime	0
wordpress	0.4	PHP/mysql blog tool	9
haproxy	0.4	Load balancer	8
influxdb	0.3	Time series database	11
elasticsearch	0.3	Search engine	12

**Table 3.** Top twenty most popular applications on Docker Hub (by billions of downloads) and the number of additional configuration options each requires beyond the *lupine-base* kernel configuration.<sup>9</sup>



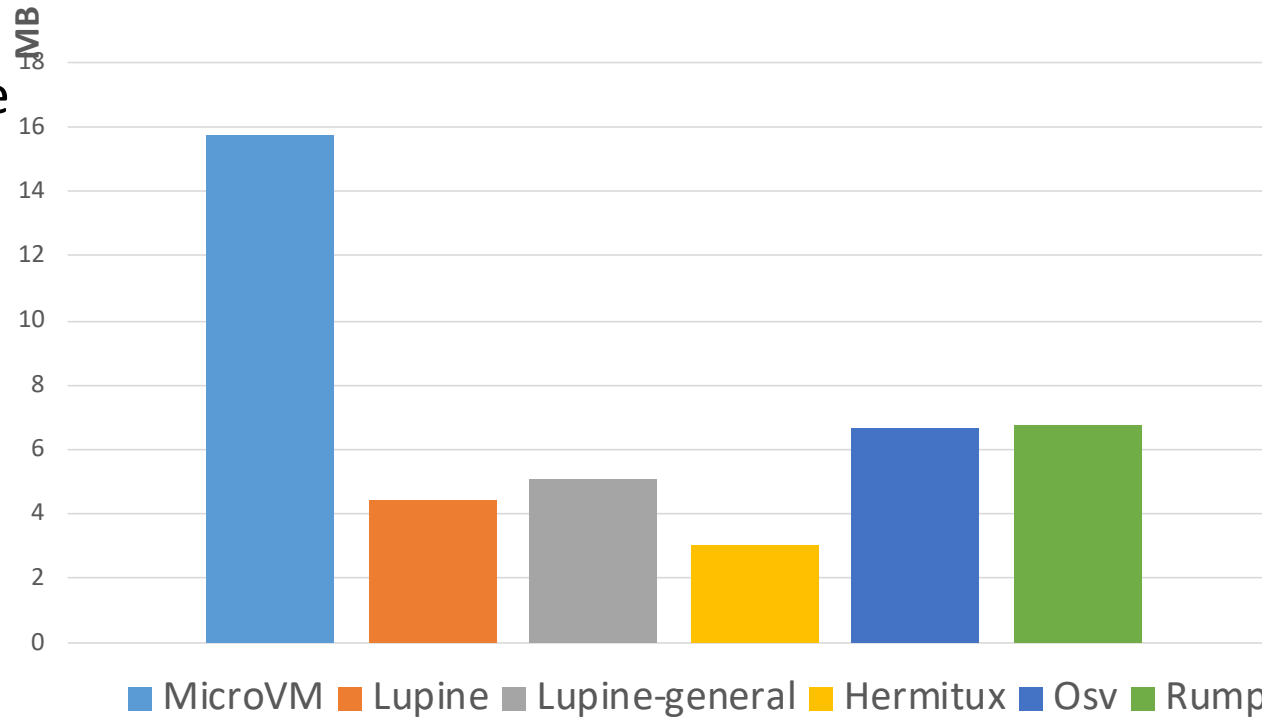
# Kernel image size

- Configuration is effective
- 4 MB
- 27% - 33% of MicroVM



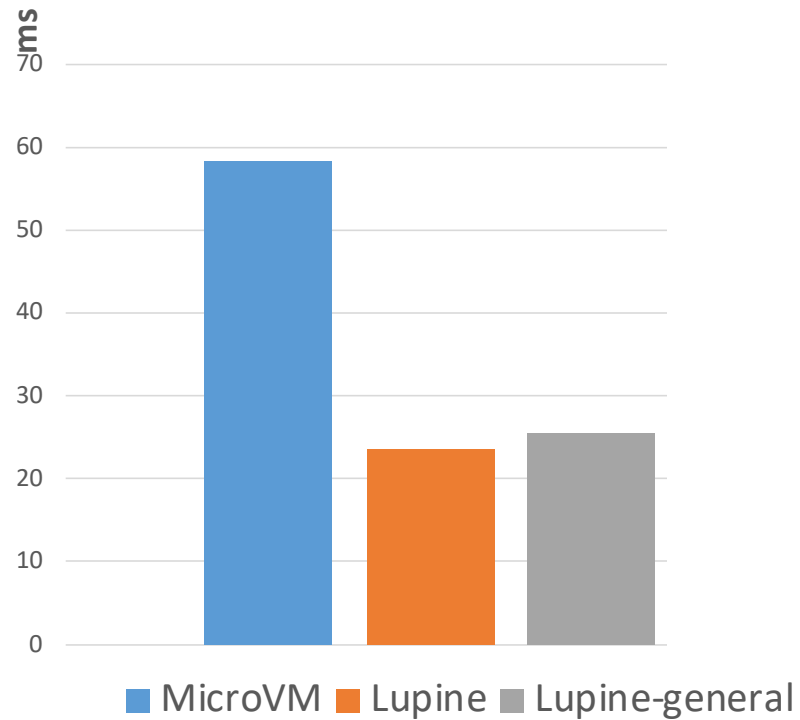
# Kernel image size

- Configuration is effective
- 4 MB
- 27% - 33% of MicroVM
- *lupine-general* is comparable with unikernels! (Rump, OSv)



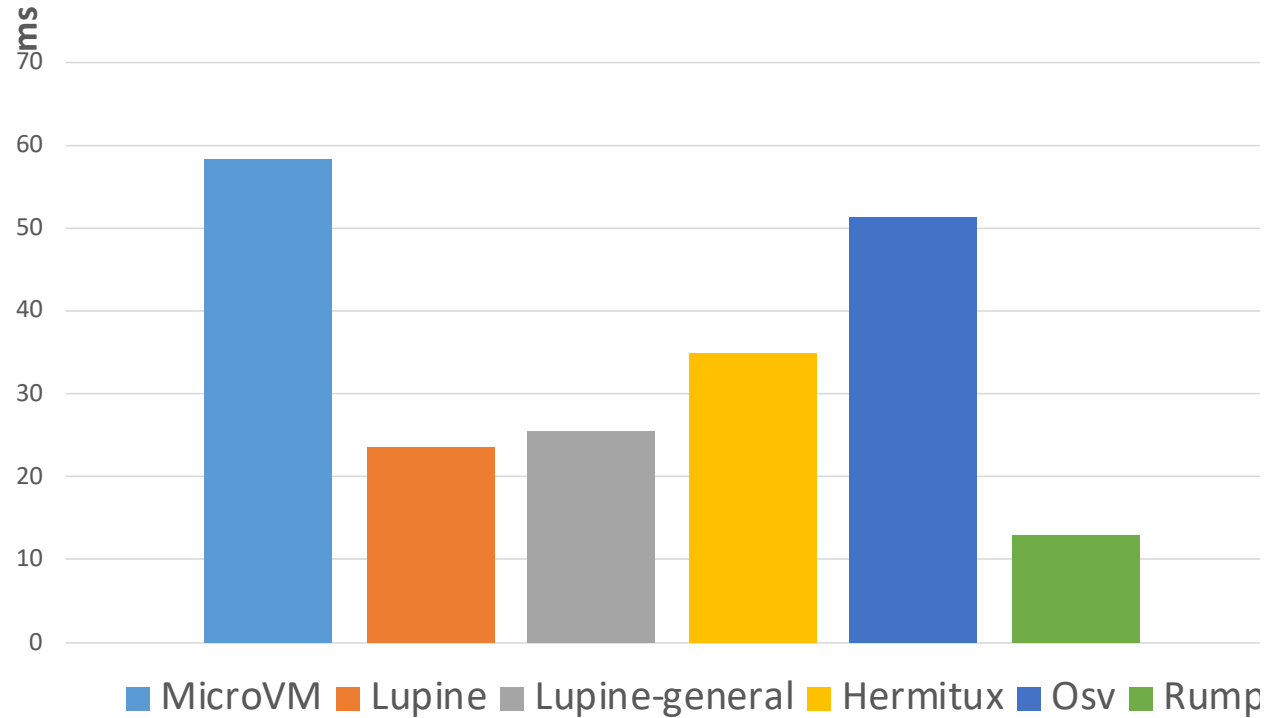
# Boot time

- Measured via I/O port write from guest
- Way better than MicroVM! (59%)



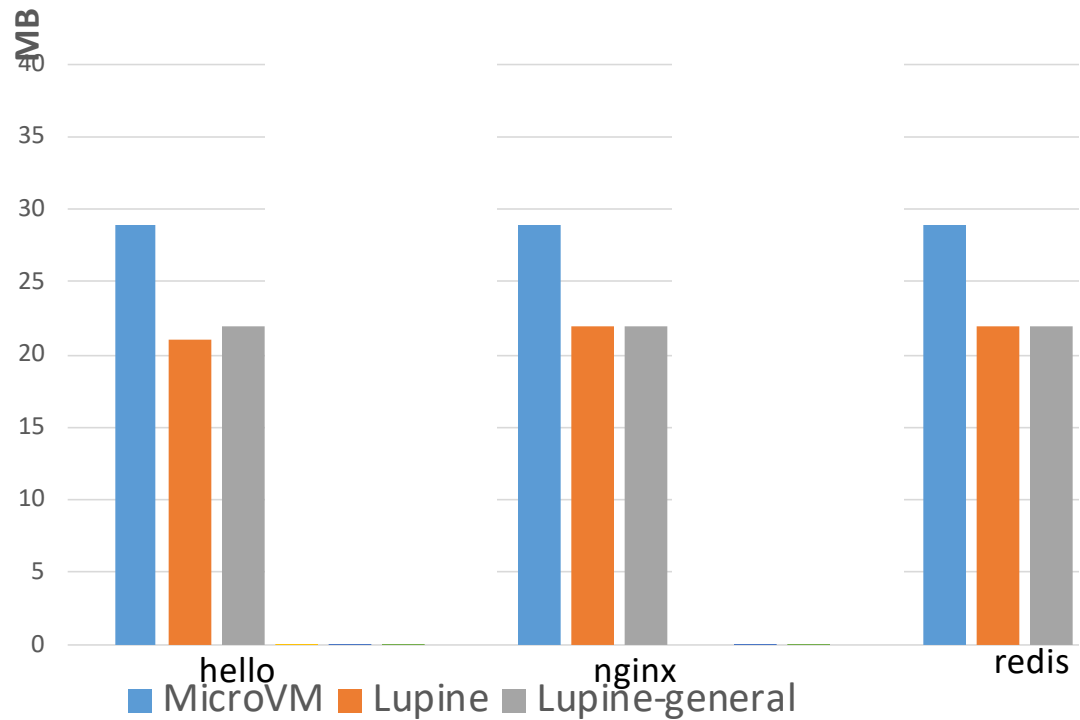
# Boot time

- Measured via I/O port write from guest
- Way better than MicroVM! (59%)
- Even *Lupine-general* boots faster than Hermitux, OSv



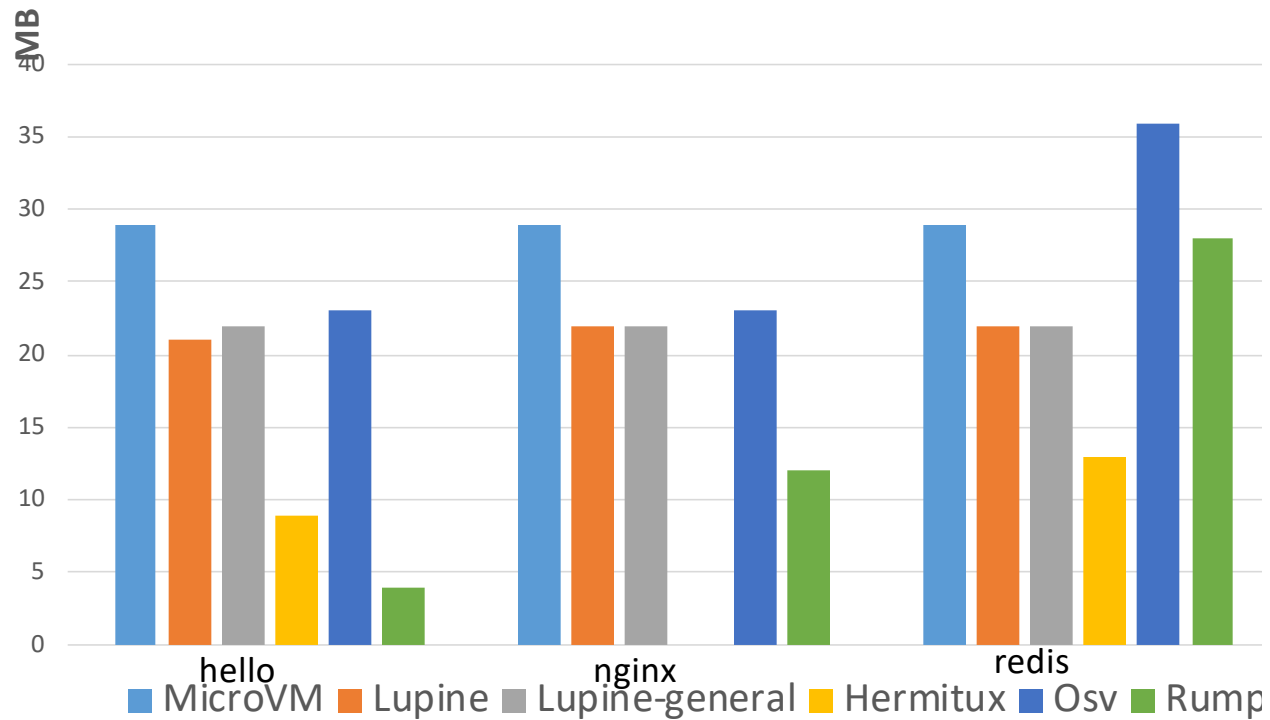
# Memory Footprint

- Repeatedly tested app with decreasing memory allotment
- Better than MicroVM(28%)



# Memory Footprint

- Repeatedly tested app with decreasing memory allotment
- Better than MicroVM(28%)





# Application performance

- Throughput normalized to MicroVM
- Lupine outperforms MicroVM by up to 29%

Name	redis-get	redis-set	nginx-conn	nginx-sess
MicroVM	1.00	1.00	1.00	1.00
<b>Lupine</b>	1.20	1.21	1.29	1.16
Lupine-general	1.19	1.20	1.29	1.15
Hermitux	.66	.67		
OSv			.87	.53
Rump	.99	.99	1.25	.53

**Table 4.** Application performance normalized to MicroVM (Note: higher value is better).

# Related work

- Unikernel-like work that leverages Linux
  - LightVM (TinyX): VMs can be as light as containers
  - X-Containers: Xen paravirt for Linux to be a libOS
  - UKL: modify Linux build to include kernel call to application main
- Linux configuration studies
  - Alharthi et al.: 89% of 1530 studied vulnerabilities nullified via config specialization
  - Kurmus et al.: 50-85% of attack surface reduction via configuration

Segue back to Dan for open challenges...

# Takeaways

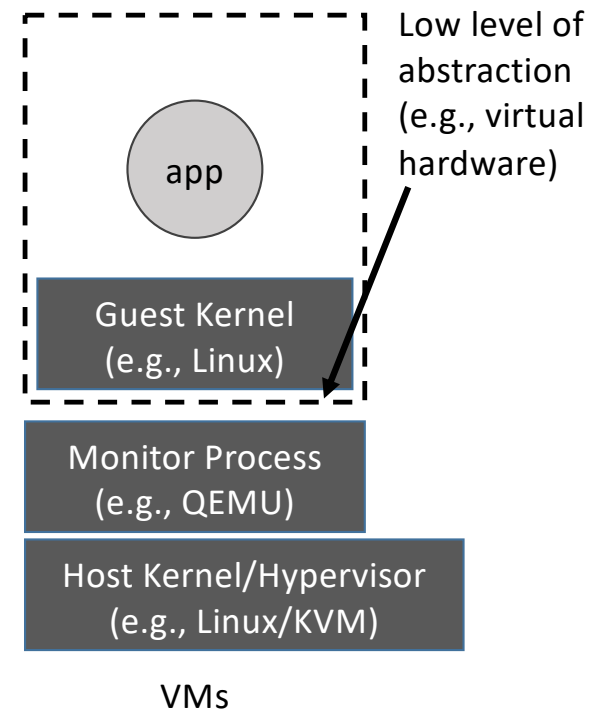
- **Specialization is important:**
  - 73% smaller image size, 59% faster boot time, 28% lower memory footprint and 33% higher throughput than the state-of-the-art microVM
- **Specialization per application may not be:**
  - 19 options (lupine-general) cover at least 83% of downloaded apps with at most 4% reduction in performance

# Getting Lupine benefits into community

- Most benefits are achieved through specialized config
  - But *[lupine-general.config](#)* can run top 20 Docker containers
- Challenges/risks
  - How do we know lupine-general is general enough?
    - Research needed: discovery vs. fallback?
  - Tension with container ecosystem (kata agent --> more general kernel config?)
    - Research needed: bloat-aware agent design?

# Continuing challenges with virtualization-enabled containers

- Sharing for container-like performance
- E.g., volume sharing
  - Virtiofs
- How to ensure safety?



# Thank you!

- EuroSys 20 Paper: <https://dl.acm.org/doi/10.1145/3342195.3387526>
- <https://github.com/hckuo/Lupine-Linux>
  
- [djwillia@us.ibm.com](mailto:djwillia@us.ibm.com)
- [hckuo2@illinois.edu](mailto:hckuo2@illinois.edu)