

We're making  
this a reality

# #include <os>

From bootloader to REST API with the new C++

Let's put some operating system  
inside your ELF binary

# Unikernels 101

- Coined in 2013 by Madhavapeddy, Mortier et. al.[1]
- Single purpose: lots of servers are anyway
- Library operating system: Link in what you need
- Turns a binary into a bootable disk image
- Mainly targeted at virtual machines and IoT

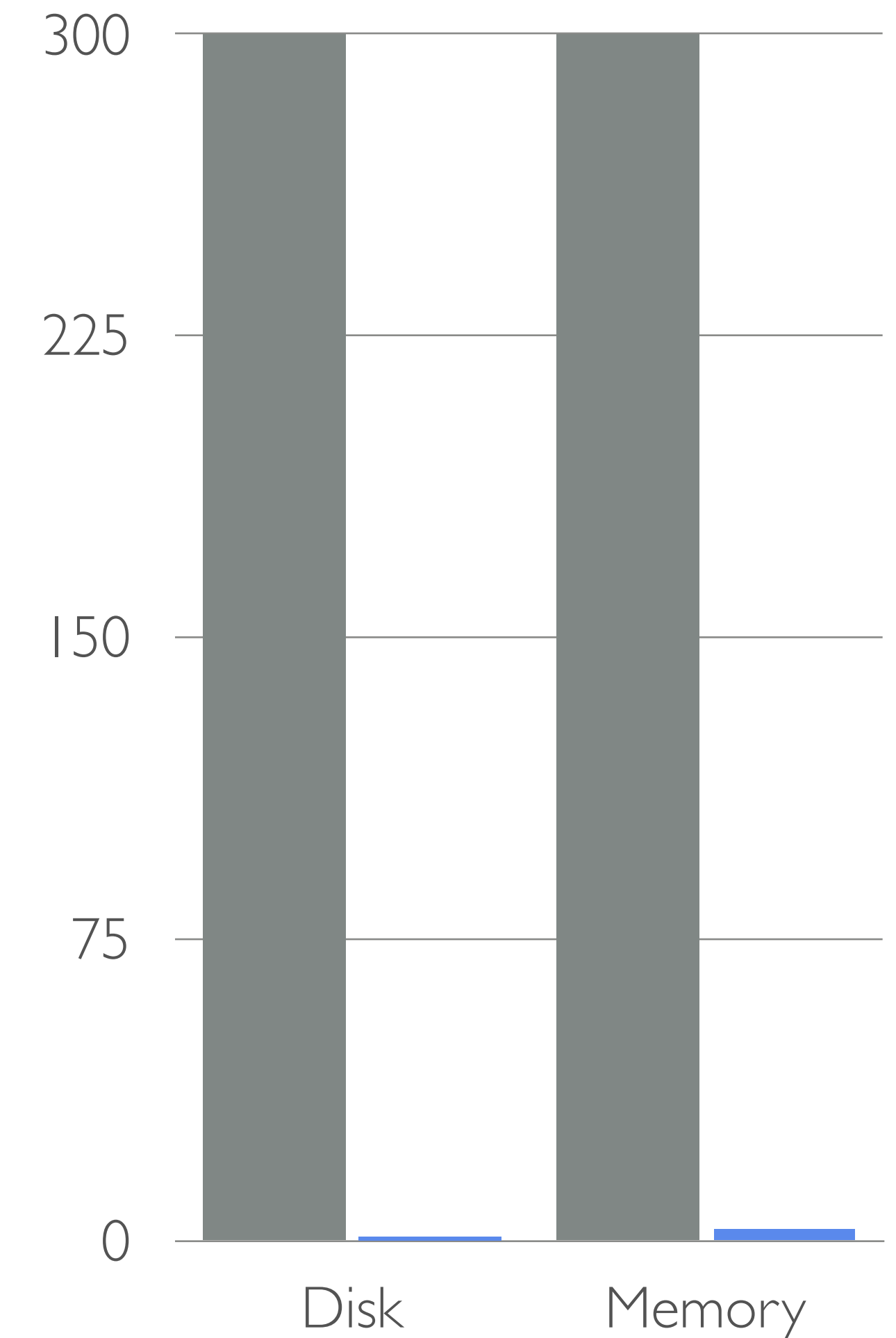
IncludeOS: A minimal, resource efficient implementation  
from scratch in modern C++ [3].



v.s.



- 300x smaller disk image
- 100x less memory usage by default
- That much reduction in attack surface



# includeOS v.s. ubuntu

- IncludeOS is single threaded by default
- Faster, more efficient for single core VM's:
  - 5-20% less CPU cycles spent for the same binary and workload [3]
  - IncludeOS boots in ~300ms. Ubuntu boots in minutes.
- Ubuntu is awesome. As a general purpose OS.

You all know what virtual machines are.  
...let's go through it anyway





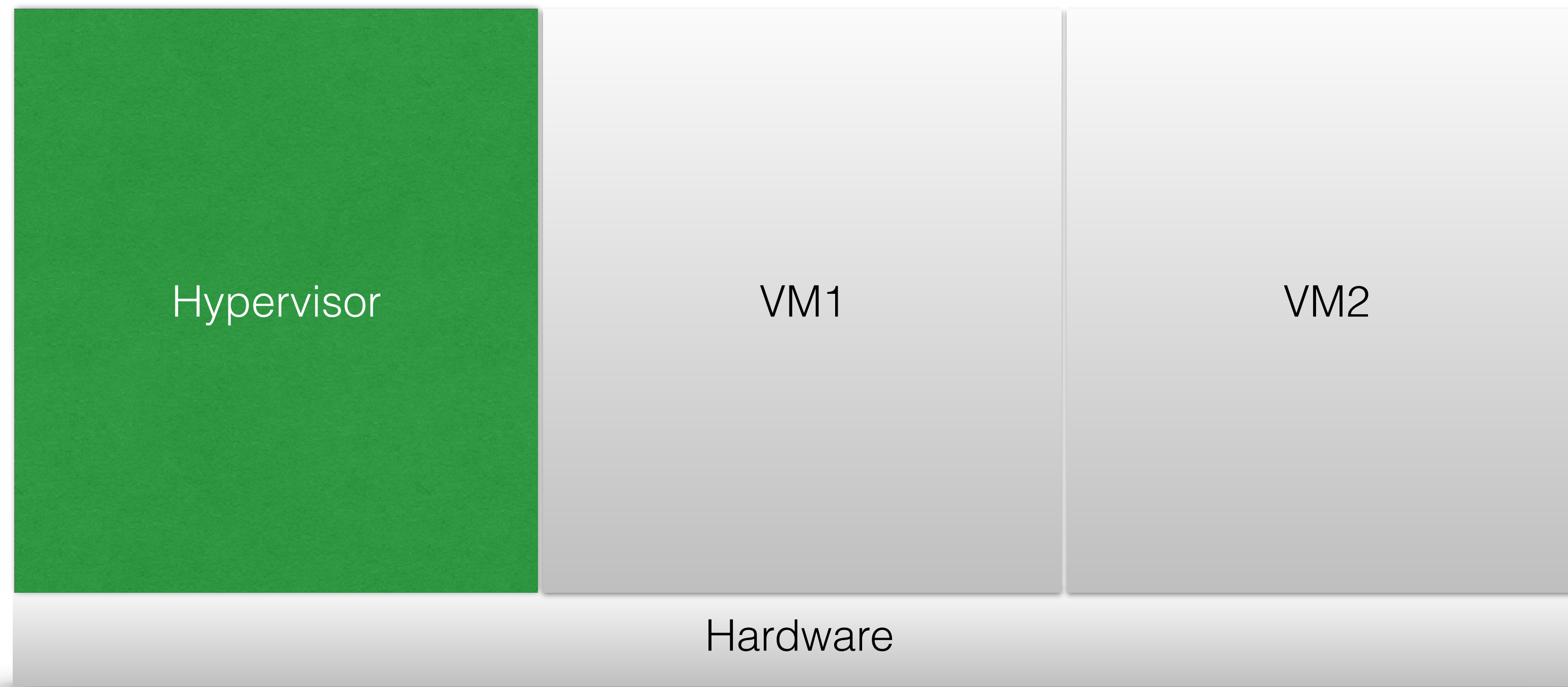
# Popek / Goldberg 1974

Gerald J. Popek, Robert P. Goldberg CACM [4]

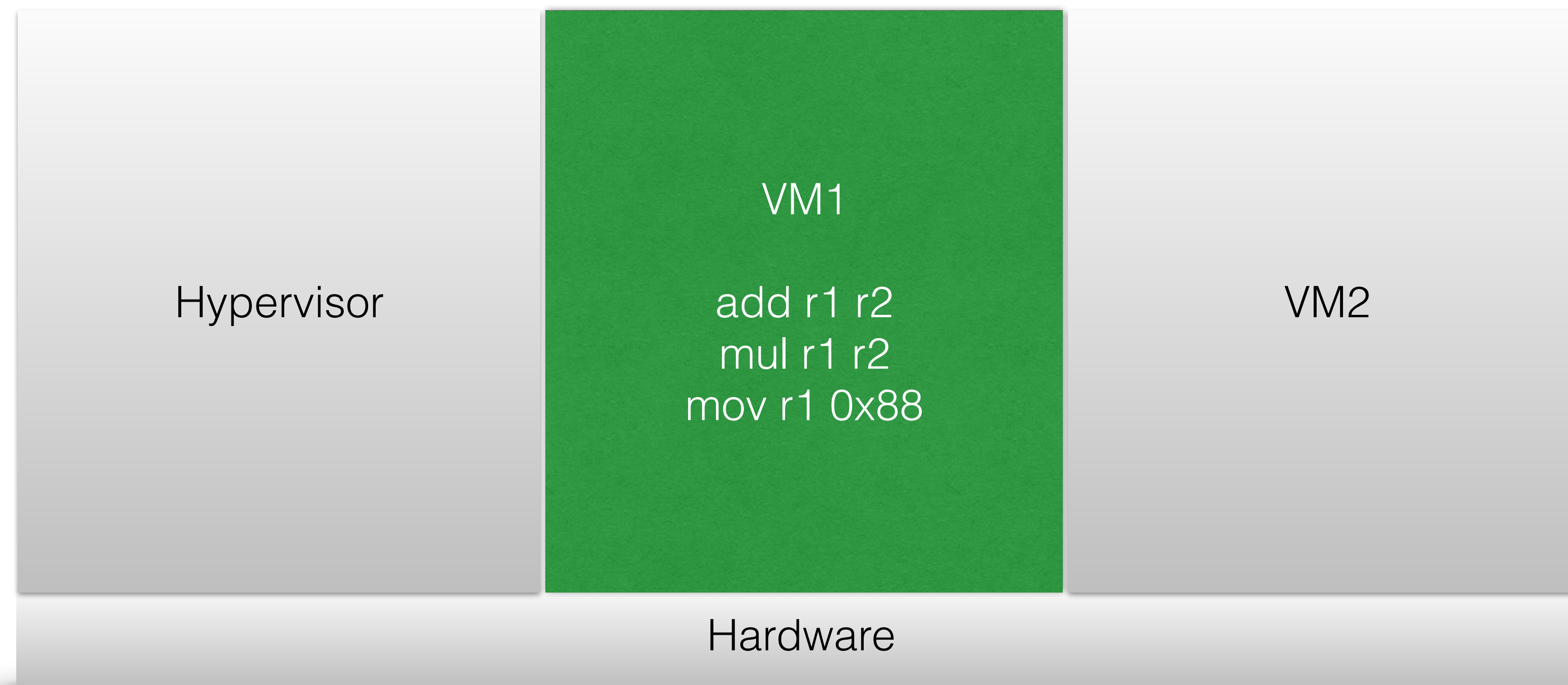
«A virtual machine is taken to be an efficient, isolated duplicate of the real machine»



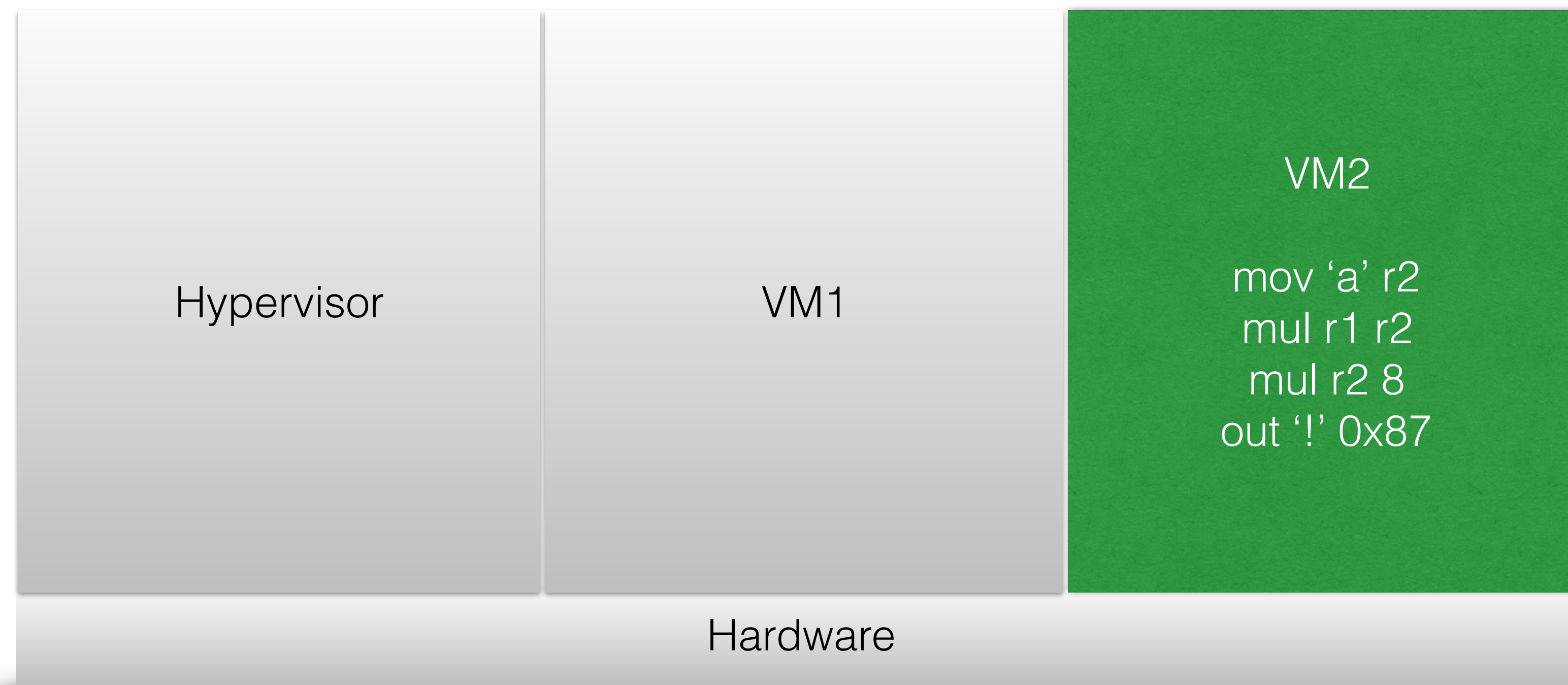
# All sensitive instructions have to trap!



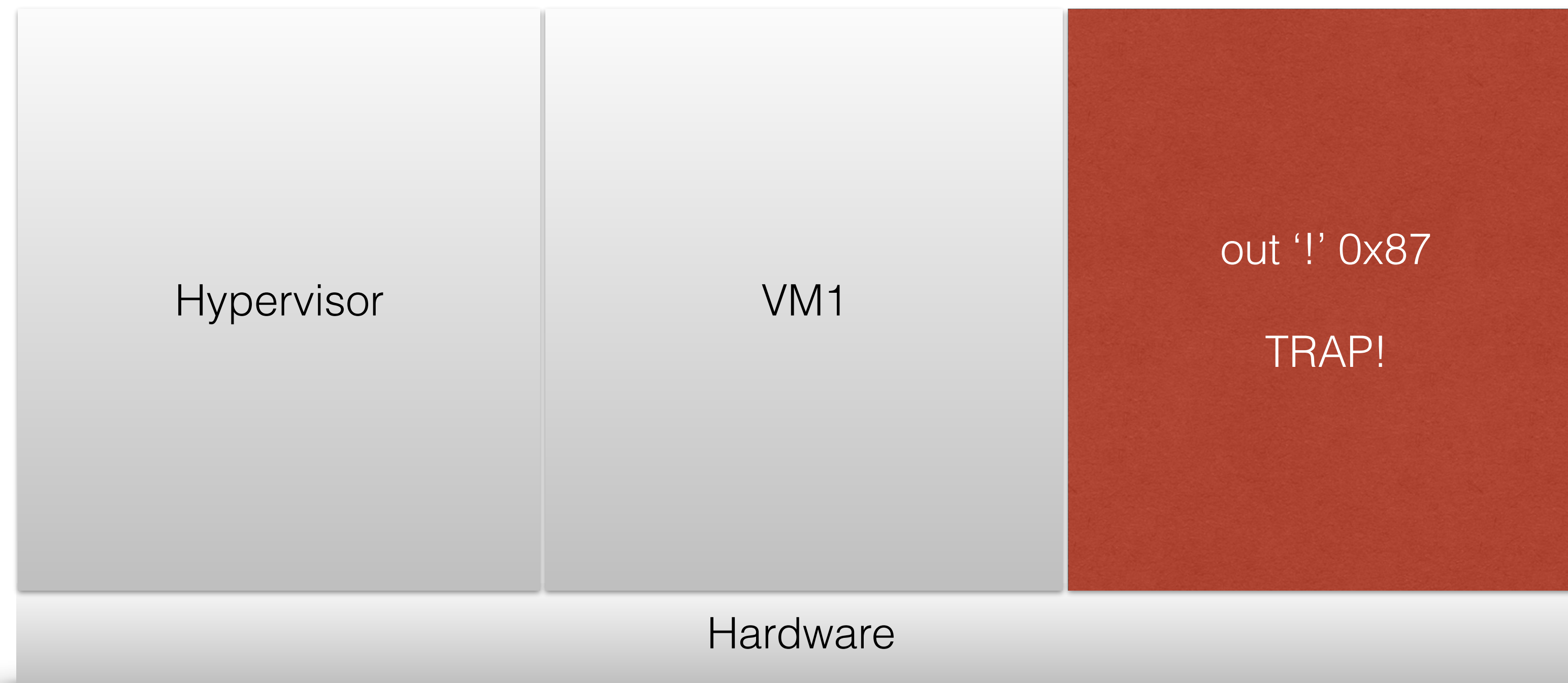
# All sensitive instructions have to trap!



# All sensitive instructions have to trap!

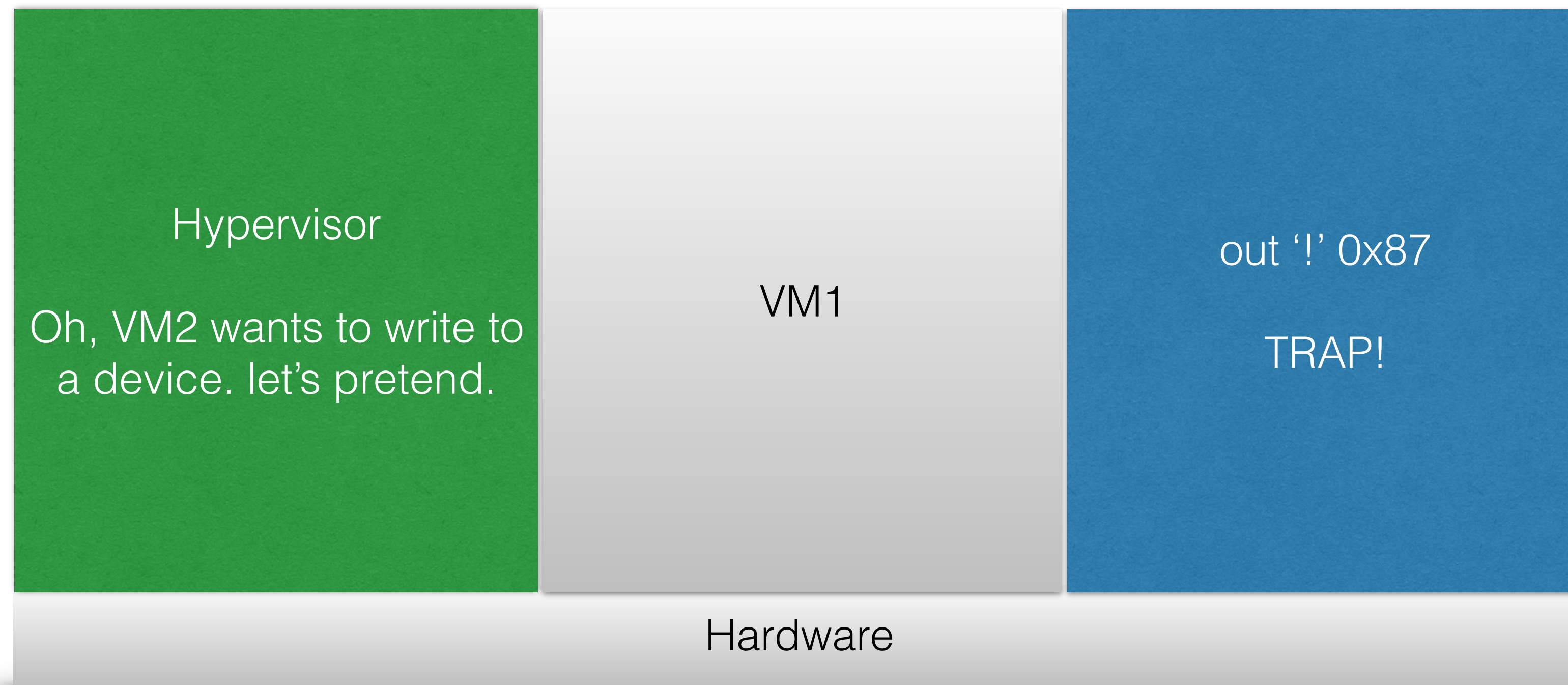


# All sensitive instructions have to trap!

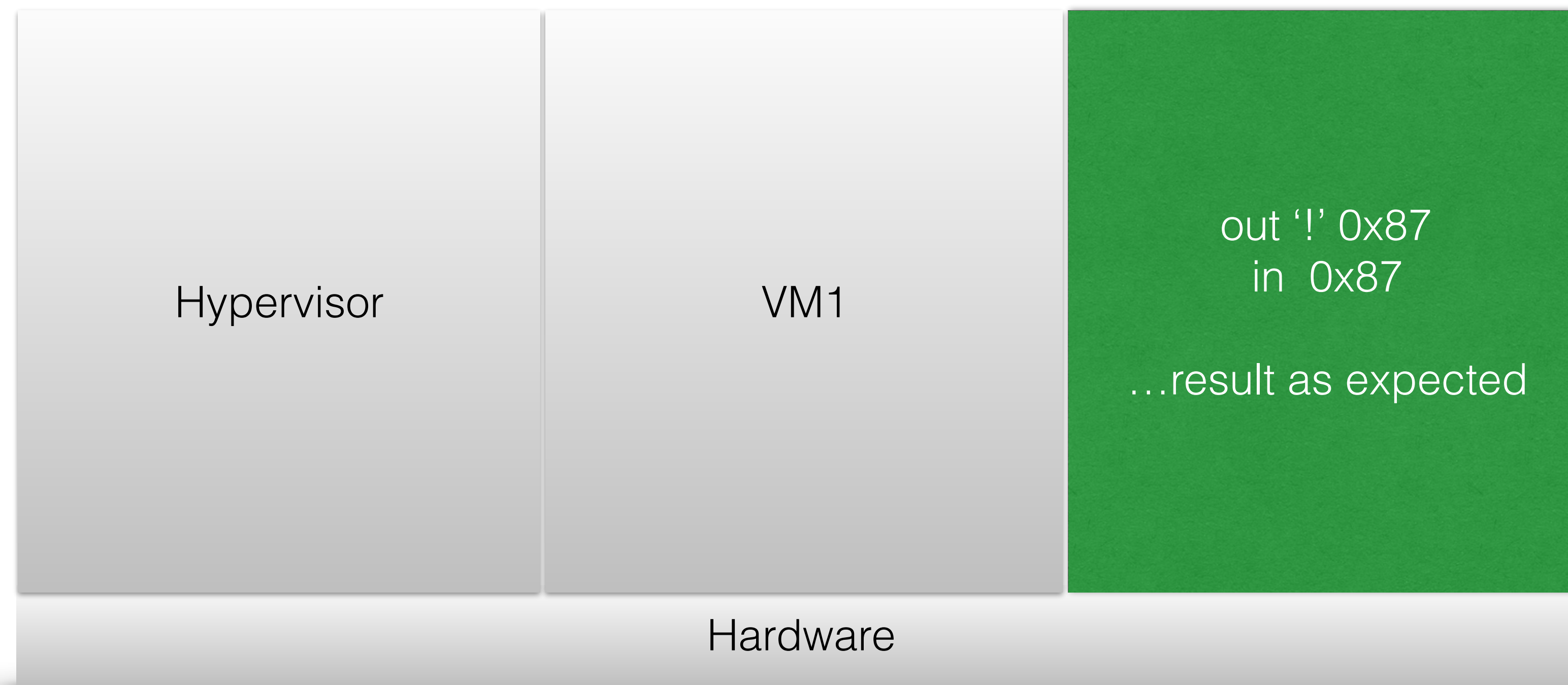




# All sensitive instructions have to trap!

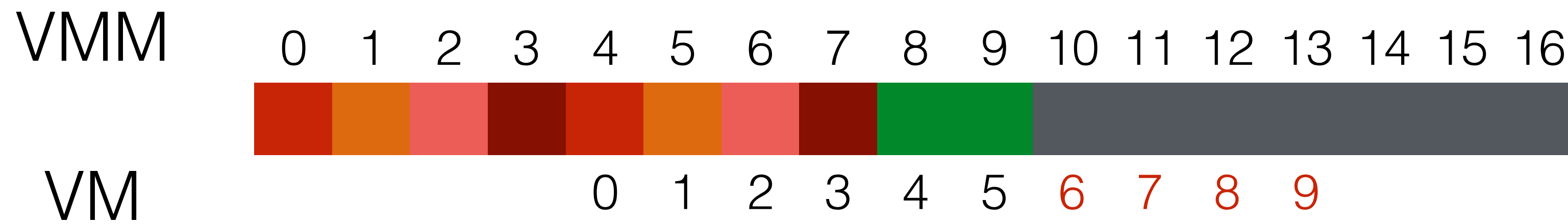


# All sensitive instructions have to trap!



# Virtual memory

If an address translation results in a physical address out of VM bounds: TRAP!



30 years later



Oh, good idea.

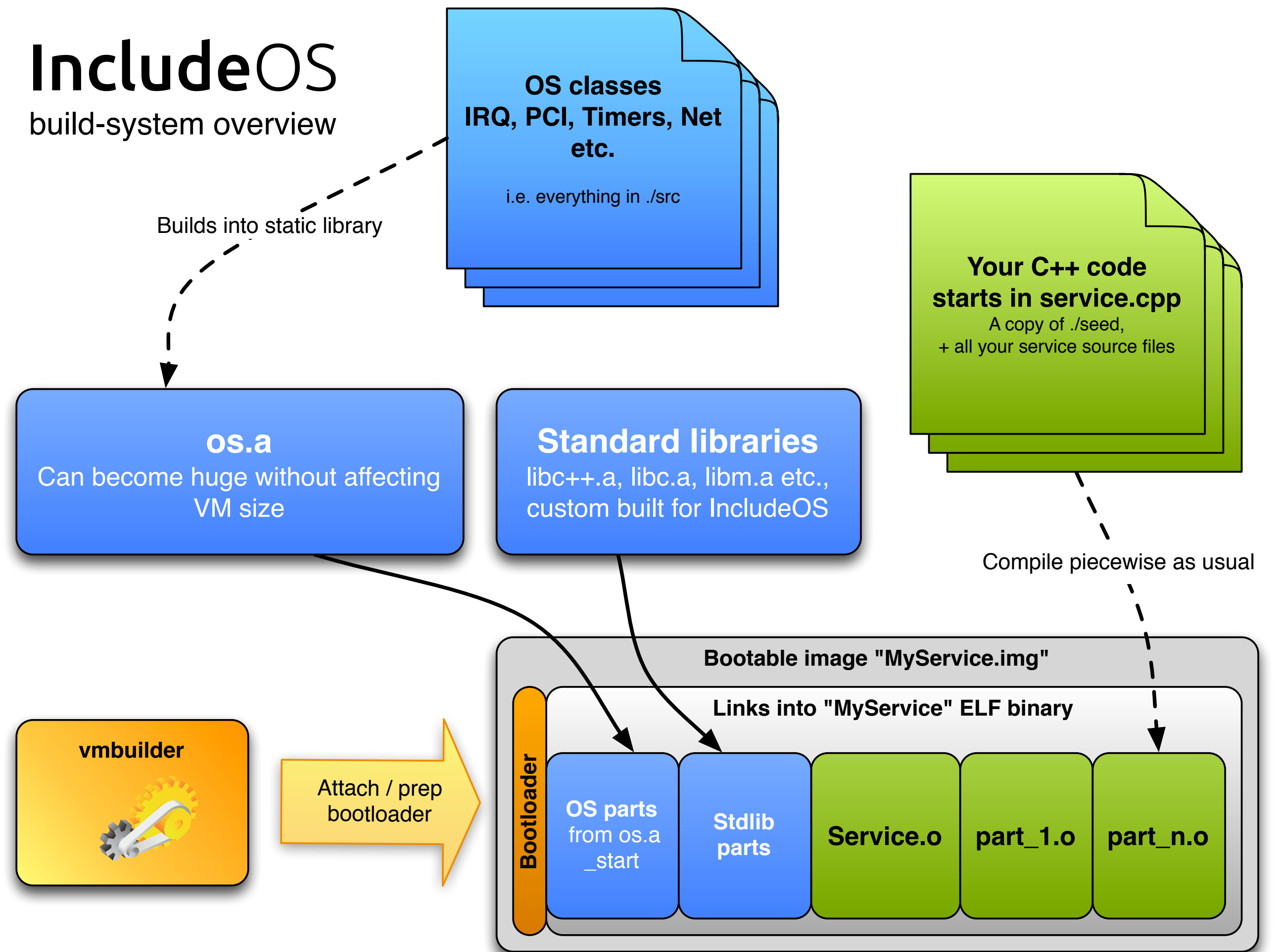


You run directly on hardware.  
The VM-trap is just a glorified context switch.

Since we're 100% self contained  
we can run on Linux, Windows, Mac unchanged

It's like a process.  
It just happens to need some operating system.

- `os.a`: static library containing all OS parts
  - Drivers are now separated
- Linker pulls in what the program needs
- OS library can bloat without affecting image size
- `vmbuild`: Attach and modify bootloader





# Why C++?

(Yes, people really ask that)

# C++ is the most evolved “independent systems language”

- You need to write arbitrary bytes to arbitrary addresses
- You need access to the full instruction set
- You can't write an OS in a 100% type safe language
- But new C++ gets you the best of both worlds

# Hello World!

“kernel” params  
with multiboot

What about  
“main”?

```
#include <os>

void Service::start(const std::string&) {
    std::cout << "Hello world - OS included!\n";
}
```

And where does stdout go?

No parent  
process - where  
do you return?

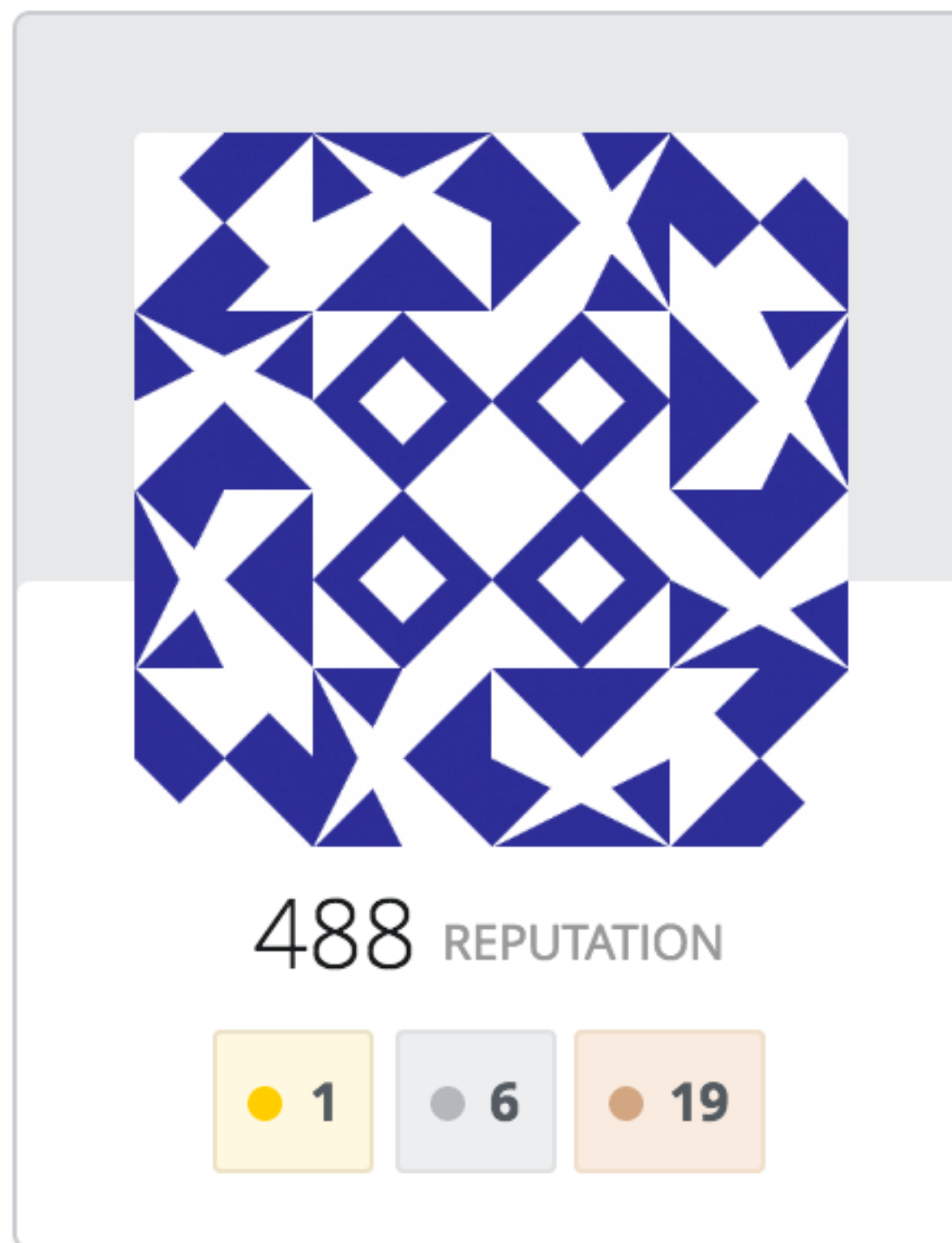


# Delegates everywhere!

A C++11 version of Don Clugstons  
“Fastest possible C++ Delegates” was found  
... on Stack Exchange

<http://codereview.stackexchange.com/questions/14730/impossibly-fast-delegate-in-c11>





**user1095108** **top 25% overall**

Apparently, this user prefers to keep an air of mystery about them.

We'll respect  
your privacy

but we seriously owe  
you a ... beverage of  
your choice

<http://codereview.stackexchange.com/questions/14730/impossibly-fast-delegate-in-c> |

# Routing stdout anywhere

```
// VGA
OS::add_stdout(VGA::get_rsprint_handler());

// Custom
OS::add_stdout([] (const std::string& data) {

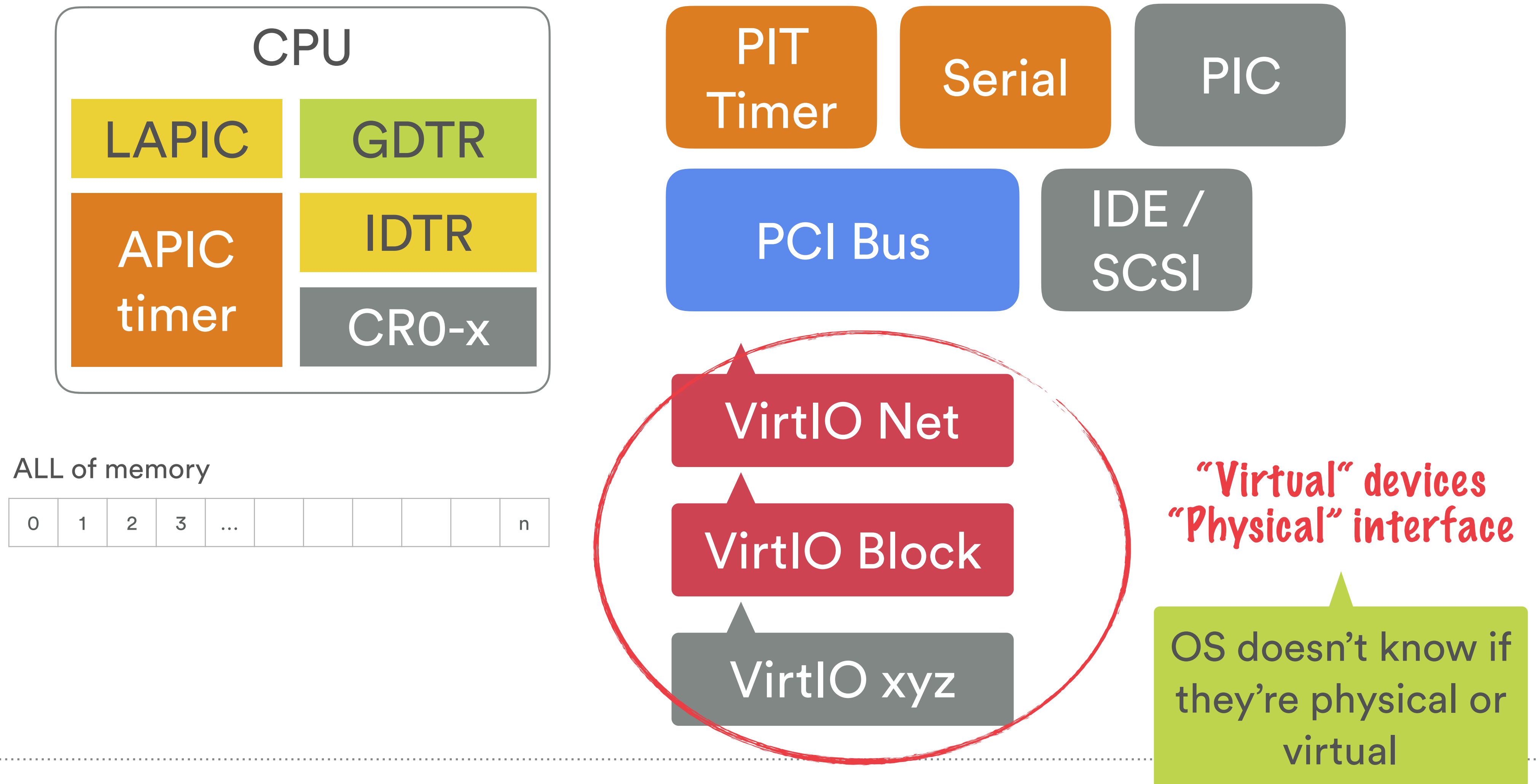
    // Memory-only (webserver 'logger')
    logger->log(timestamp() + data);

});
```

In principle: nothing

What's different  
about an OS  
for virtual machines?

# Hardware overview



How to include only the drivers you need -  
without having to mention them in code?



# Solution: link in self registering drivers

Service Makefile:

DRIVERS = virtionet ...

```
// From IncludeOS v0.9, src/drivers/virtionet.cpp  
// ...VirtioNet implementation
```

```
/** Register VirtioNet's driver factory at the PCI_manager */  
struct Autoreg_virtionet {  
    Autoreg_virtionet() {  
        PCI_manager::register_driver<hw::Nic>(hw::PCI_Device::VENI  
            0x1000, &VirtioNet::new_instance);  
    }  
} autoreg_virtionet;
```

Deployment tools /  
sysadmins can apply  
appropriate drivers

1) Global constructor registers a  
delegate to driver initialization.

2) OS probes PCI bus and applies any  
registered driver

# Kernel overview

```
// 1. Initialize stack, heap, .bss etc.  
// 2. Call global constructors  
// 3. Initialize hardware  
// 4. Call Service::start()
```

```
void OS::event_loop() {
```

```
    while (power_) {  
        IRQ_manager::get().notify();  
    }
```

```
    // Cleanup  
    Service::stop();  
    // ACPI shutdown sequence  
    hw::ACPI::shutdown();  
}  
// From src/kernel/os.cpp
```

Call any delegates  
subscribing to IRQ's,  
then "hlt"

Note: All subscribed  
IRQ's are deferred

# Subscribing to interrupts

// Actual code from src/drivers/virtionet.cpp

```
auto del(delegate<void()>
        ::from<VirtioNet,
        &VirtioNet::irq_handler>(this));
```

Create delegate to  
member function of  
“this”

```
IRQ_manager::get().subscribe(10, del);
```

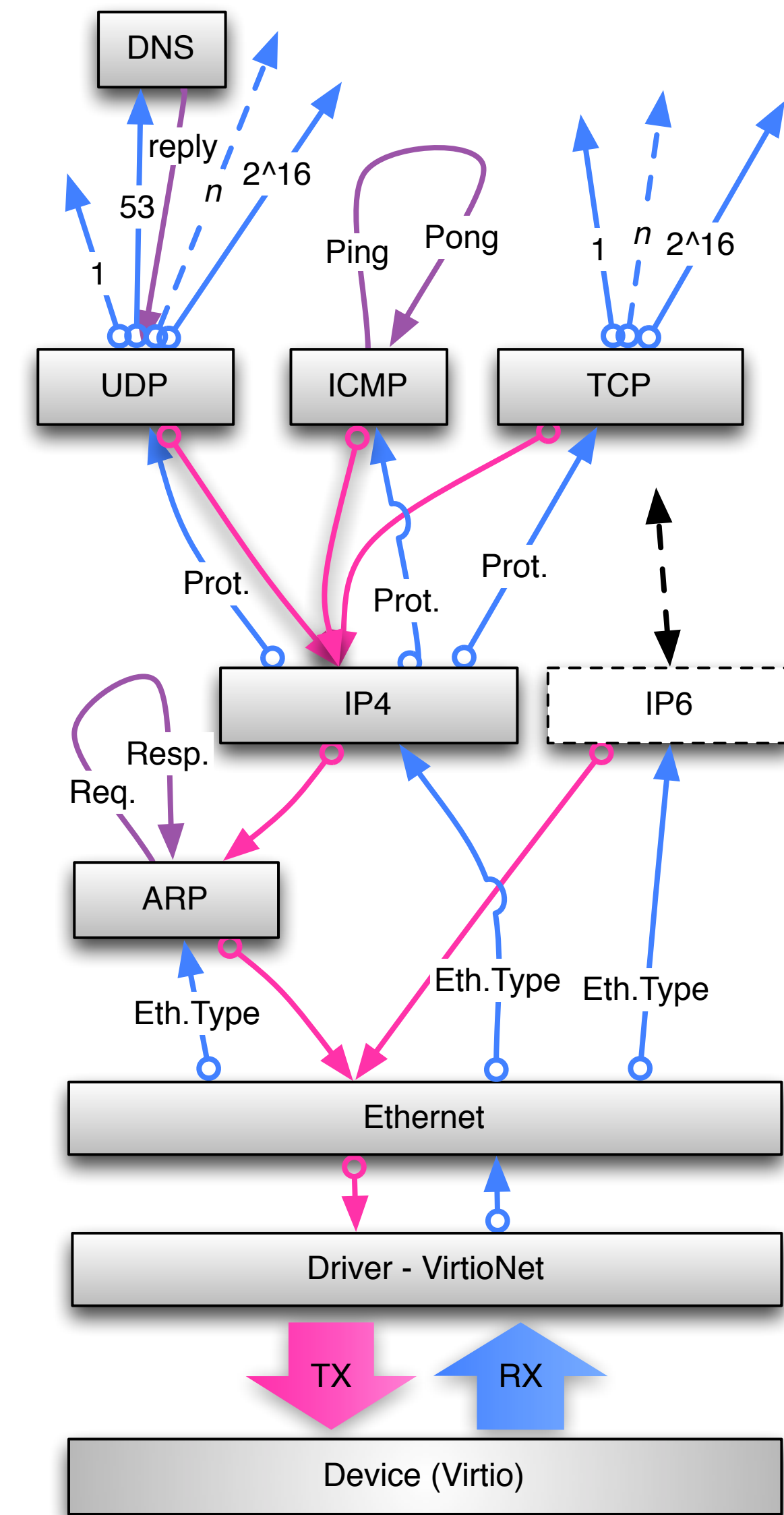


For most use cases don't worry  
about IRQ (this is driver code).  
Abstractions coming up

Subscribe to a given  
IRQ number  
called by event loop

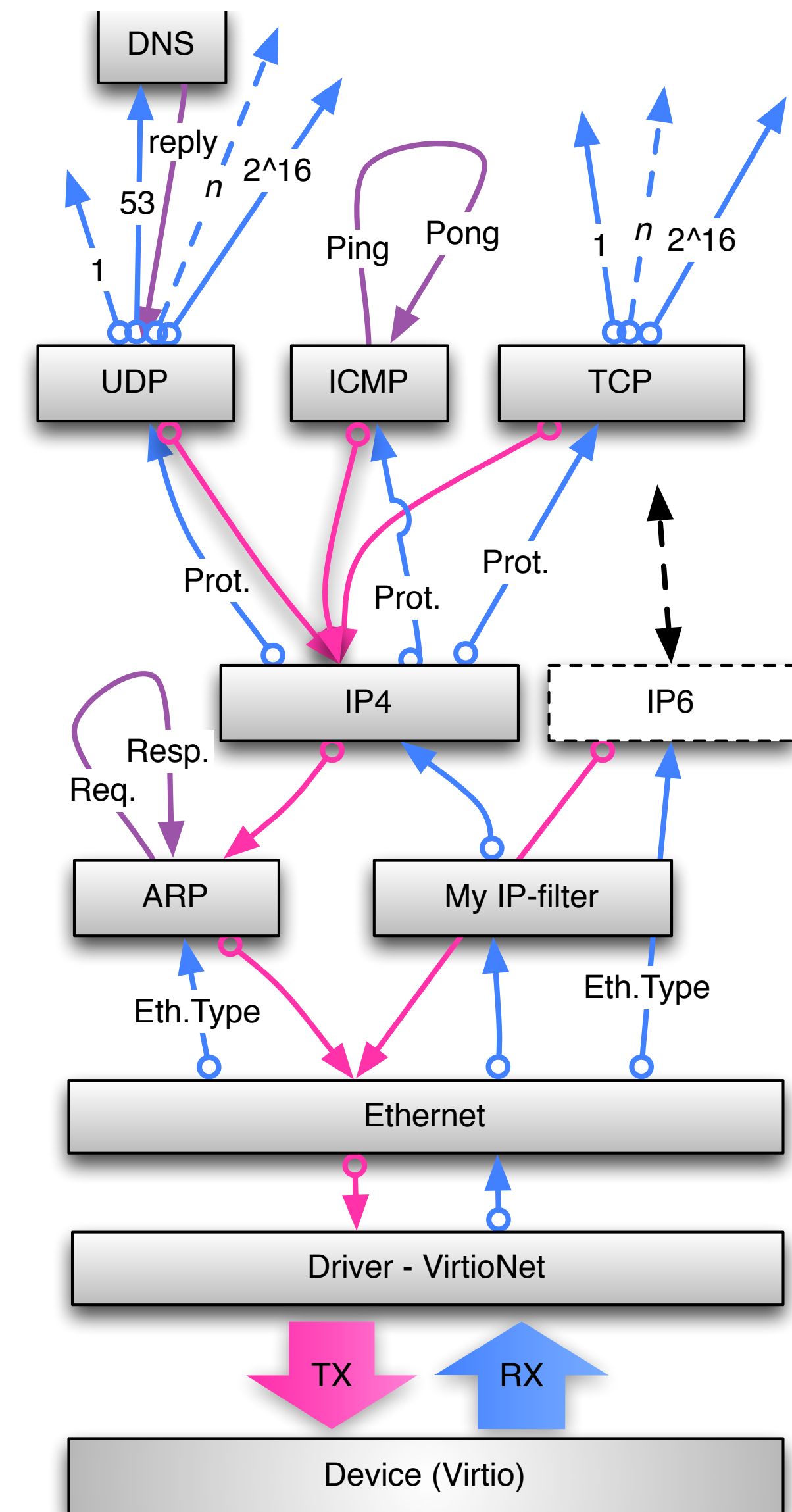
# Wiring the IP stack

- Each “box” is a class instance
- All connections are delegates
- The wiring happens in the “Inet” wrapper class
- Rewiring during runtime is easy
  - e.g. for activating and deactivating an IP-filter
  - Can be done from the application



# Wiring the IP stack

- Each “box” is a class instance
- All connections are delegates
- The wiring happens in the “Inet” wrapper class
- Rewiring during runtime is easy
  - e.g. for activating and deactivating an IP-filter
  - Can be done from the application





# Unikernels:

## The default should be minimal

- Single threaded by default: “no overhead principle”
  - In a VM, threading is more expensive [5]
- 100% async: cheapest way to get performance
- No blocking POSIX calls (yet)
  - You’ll get it. It’s your “constitutional right...”

# Writing a TCP server

```
#include <os>
#include <net/inet4>

void Service::start(const std::string&) {

    auto& server = net::Inet4::stack().tcp().bind(80);

    server.on_connect([] (auto conn) {
        conn->on_read(1024, [conn] (auto buf, size_t n) {
            conn->write("My first unikernel!\n"s);
        });
    });
}
```

Attach any delegate  
to TCP events

# How do you break in?

- We don't know
- If you get in, you're in a hardware protected sandbox
- No shell or other “ad-hoc” code access
- Minimal attack surface: no tools left for the thieves
- Everything inside can be randomized



You made it!

“Leave no room for a lower-level language below C++  
(except assembler)”

- Bjarne Stroustrup

# But do we need another language above C++?

The new C++ is awesome for high-level abstractions



# Introducing “Acorn”

IncludeOS web server  
with C++ framework for REST APIs

<https://github.com/includeos/acorn>



Static content served from FAT file system

Can serve from memdisk inside ELF or VirtioBlock device

IncludeOS supports REST.. Look, a squirrel!

	Name	Age	Occupation	Captured
	Rico	28	Mad Scientist	2016-09-18 12:00
	Martin	16	Build Master	2016-09-18 12:00
	Ingve	24	Inter...ion Master	2016-09-18 12:00

Demo of POST'ing data over REST

Replace with your own content

/static/books/

Type	Name	Size	Modified
↑	..	-	-
📄	<a href="#">poetics.txt</a>	125.42 KB	N/A
📄	<a href="#">borkman.txt</a>	172.42 KB	N/A
📄	<a href="#">fables.txt</a>	241.10 KB	N/A



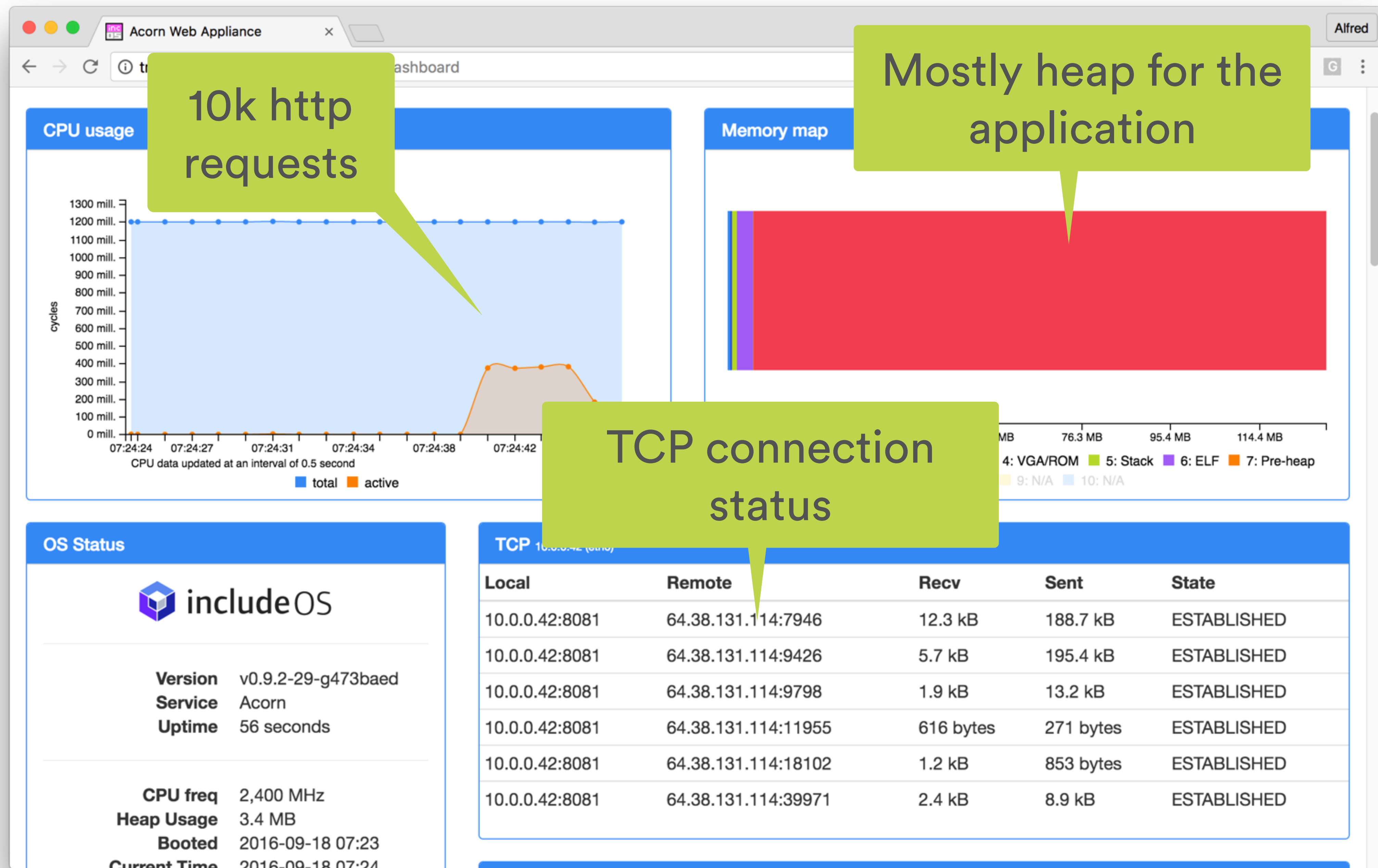
“C++ is not the language for putting  
up a simple web page”

– Bjarne Stroustrup earlier today

How about a fancy web page?



Let's do a dashboard with live CPU stats and live profiling.  
It's surprisingly easy with C++





statman: Centrally located stats for everything, including your service

now

timers	
oneshot_started	70499
oneshot_stopped	70498
periodic_started	3
periodic_stopped	0
syscalls	
sbrk	51
cpu0	
cycles_hlt	333906609728
cycles_total	343925702020

Timer IRQ-based live stack sampling

### Stack Sampler

Active / Asleep:  
1.33 / 98.67 %

Percent	Total	Addr	Function
24.36	456	0x2370e0	Virtio::Queue::kick()
6.20	116	0x239750	net::TCP::checksum(std::__1::shared_ptr<net::tcp::Packet>)
4.70	88	0x28ca40	_malloc_r
2.83	53	0x23d8d0	net::tcp::Connection::rtx_start()
2.78	52	0x28ba80	_free_r
2.67	50	0x2b4c20	bool std::__1::basic_regex<char, std::__1::regex_traits<char>>::__match_at_start_ecma<std::__1::allocator<std::__1::sub_match<char const*> > >(char const*, char const*, std::__1::match_results<char const*, std::__1::allocator<std::__1::sub_match<char const*> > const, std::__1::regex_constants::match_flag_type, bool) const
2.40	45	0x293130	_svfprintf_r
2.30	43	0x28d69c	
1.71	32	0x224f20	VirtioNet::msix_recv_handler()
1.28	24	0x2ffc80	std::__1::__shared_weak_count::__release_shared()
1.18	22	0x22dc60	RTC::now()
1.12	21	0x233400	hw::Serial::write(char)

Acorn Web Appliance

trident3.vlab.cs.hioa.no:8081/app/dashboard

irq

126	0
0	164273
64	30719
65	30303
66	0
67	30212
120	328

eth0

packets_rx	100757
packets_tx	70976

ethernet

packets_rx	100758
packets_tx	70976
packets_dropped	0

arp

Logger

```

[2016-09-18.14:23:53]
[2016-09-18.14:23:53]
[2016-09-18.14:23:53]      +--> nome.js
[2016-09-18.14:23:53]      +--> squirrel.js
[2016-09-18.14:23:53]
[2016-09-18.14:23:53]      +--> stats.js
[2016-09-18.14:23:53]      +--> app.js
[2016-09-18.14:23:53]
[2016-09-18.14:23:53]      +--> index.html
[2016-09-18.14:23:53]
[2016-09-18.14:23:53]      +--[ images ]
[2016-09-18.14:23:53]      +   .
[2016-09-18.14:23:53]      +   ..
[2016-09-18.14:23:53]      +--> prettypicture.jpg
[2016-09-18.14:23:53]
[2016-09-18.14:23:53]      +--[ css ]
[2016-09-18.14:23:53]      +   .
[2016-09-18.14:23:53]      +   ..
[2016-09-18.14:23:53]      +--> style.css
[2016-09-18.14:23:53]
[2016-09-18.14:23:53]      +--> index.html
[2016-09-18.14:23:53]      +--> README.md
[2016-09-18.14:23:53]
[2016-09-18.14:23:53] =====
[2016-09-18.14:23:53]      [ Router ] Registered routes:
[2016-09-18.14:23:53] GET      /api/squirrels/
[2016-09-18.14:23:53] GET      /api/users/
[2016-09-18.14:23:53] GET      /api/users/:id(\d+)/:name/something/:something([a-z]+)
[2016-09-18.14:23:53] GET      /api/dashboard/
[2016-09-18.14:23:53] GET      /api/dashboard/memmap
[2016-09-18.14:23:53] GET      /api/dashboard/stack_sampler
[2016-09-18.14:23:53] GET      /api/dashboard/status
[2016-09-18.14:23:53] GET      /api/dashboard/statman
[2016-09-18.14:23:53] GET      /api/dashboard/tcp

```

Logs: stdout routed to in-memory ring-buffer



Individual REST  
endpoints for each  
component.  
  
100% done in  
application.

The OS exposes  
nothing  
except an API to  
access STL data  
structures from App



# Routes!

RESTful endpoints for anything in IncludeOS

# “Node/Express.js”-style routes

```
Router route http://my_domain/api/users/42
// Assume matches regex and captures “id”, “42”

// GET /api/users/:id where 'id' is a numeral
router.on_get("/api/users/:id(\\d+)", [users] (auto req, auto res) {
    try {
        // Try to retrieve param "id"
        auto& params = req->params();
        std::string id = params.get("id");
        // Look for the user inside the bucket
        auto& user = users->pick_up(std::stoi(id));

    } catch (const std::exception& e) {
        res->send_code(http::Not_Found);
    }
});
```

Regex-based mini-language for expressing routes.  
ported to C++ from [express.js](#)

The “:id” part is made accessible to the app as key to a value.

# Default route: serve from FAT fs

```
Router router;

// Serve index.html from root URL
router.on_get("/", [](auto, auto res) {
    // Stat disk for index.html
    disk->fs().cstat("/public/index.html", [res] (auto err, const auto& entry) {
        if(err)
            res->send_code(http::Not_Found);
        else
            // Serve index.html
            res->send_file({disk, entry});
    });
});
```



# Middleware

Applying a stack of functions to all requests

# “Node/Express.js”-style middleware

```
server.use([](auto req, auto res, auto next) {  
    if(has_content_type_json(req)) {  
        auto json = std::make_shared<Json_doc>();  
        json->doc().Parse(req->get_body().c_str());  
        req->set_attribute(json);  
    }  
    (*next)();  
});
```

Call next middleware  
(or exit to hang up)

Check for JSON  
content type

Parse and attach to  
request

```
router.on_post("/users", [](auto req, auto res) {  
    auto json = req->get_attribute<Json_doc>();  
    if(json) {  
        auto& doc = json->doc();  
        // Register user based on parsed JSON data  
    }  
});
```

Expect parsed JSON



# Readymade middleware

```
// Director: Provide an auto-generated HTML directory listing
auto director = std::make_shared<middleware::Director>(disk, "/public/static");
server.use("/static", director);

// Cookie parser: parse all cookies on incoming requests
auto cookie_parser = std::make_shared<middleware::CookieParser>();
server.use(cookie_parser);

// Butler: Serve static files on root route with disk root on /public
auto opt = {"index.html", "index.htm"};
auto butler = std::make_shared<middleware::Butler>(disk, "/public", opt);
server.use(butler);
```

# Performance:

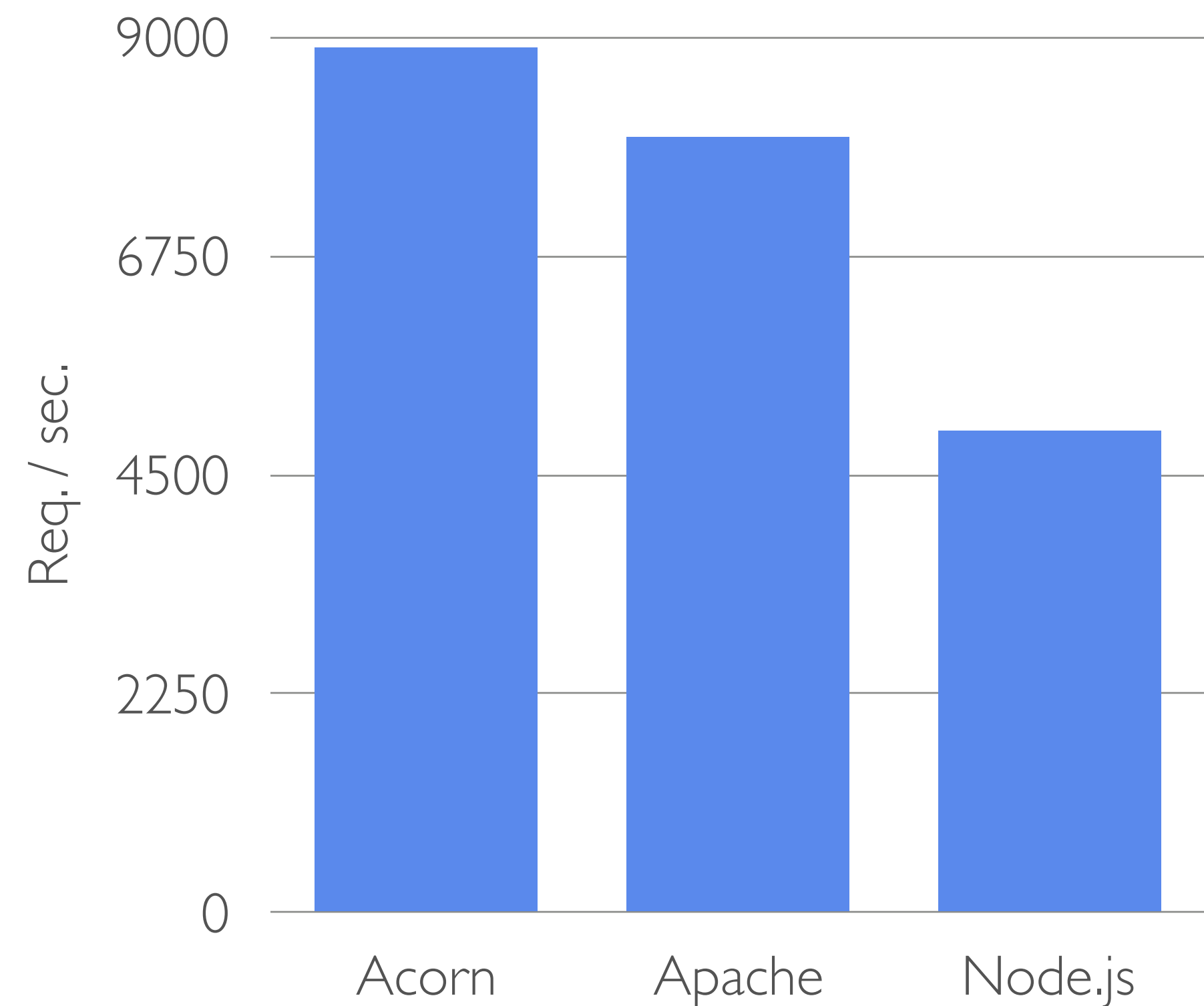
## Virtual machines aren't faster than bare metal



The goal is to have thousands of self-contained, secure micro services

# “Acorn” preliminary performance data

- All services run in single core VM
- Apache / Node on Ubuntu 16.10
- “Acorn” on IncludeOS
  - 11% faster than Apache
  - 79% faster than Node.js
- 10k http req. with httpperf
  - Average over 66 samples
  - Intel i7 desktop PC w. Ubuntu 16.04



# Deploy in the cloud!



# Deploy in OpenStack

Literally three  
commands

```
$ make  
$ glance --image-create --file <My_image.img> ...  
$ nova boot --image ...
```

Get an IP,  
Fire up a browser!

We'll be adding support for  
AWS, Google compute,  
Azure, you name it.

# Project status

- 1.0 “Stable API” depends on your feedback!
- Production ready core by years end
- “Acorn”: open source alpha today (!)
- TLS, POSIX and multi-arch are priorities
- Tooling improvements:
  - project unik for deployment / management
  - Package manager (3rd party)
  - CMake builds
  - IncludeOS specific CLI tool



# Recap

- Virtual machines aren't heavy - OS'es are
- IncludeOS is faster, more efficient, more secure for single purpose
- C++ is great at both high- and low level.
  - Definitely the best language for the job



# Please contribute!

 Follow @includeos

- Pull Requests
- Ideas for cool use cases
- Expert advice!

alfred@includeos.org

## Demo!

Open Content session  
showing the implementation  
of an IRC server on  
IncludeOS on Thursday

Fork:



Chat:



# Questions?



# References

1. “Unikernels: Library operating systems for the cloud”, Anil Madhavapeddy, Richard Mortier et. al, ASPLOS 2013
2. “Unikernels: Rise of the Virtual Library Operating System”, Anil Madhavapeddy, Dave Scott, ACM Queue, 2013
3. “IncludeOS: A minimal, Resource efficient Unikernel for Cloud Systems, A. Bratterud, A. A. Walla et. al, IEEE CloudCom 2015
4. “Formal requirements for Virtualizable Third Generation Architectures”, G. Popek, R.P. Goldberg, Communications of the ACM Vol. 17, 1974
5. “APLE: Addressing Lock Holder Preemption Problem with High Efficiency”, J. Shan, X. Ding and N. Gehani, IEEE CloudCom 2015