

ToroV, a kernel in user-space, or sort of

www.torokernel.io

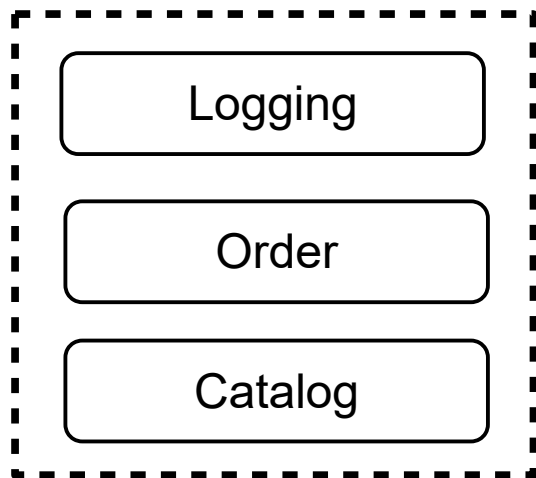
Matias Vara Larsen
matiassevara@torokernel.io

Who am I?

- I enjoy working on operating systems and playing with virtualization
- I worked at Citrix, Tttech, Huawei ...
- <https://github.com/MatiasVारा>

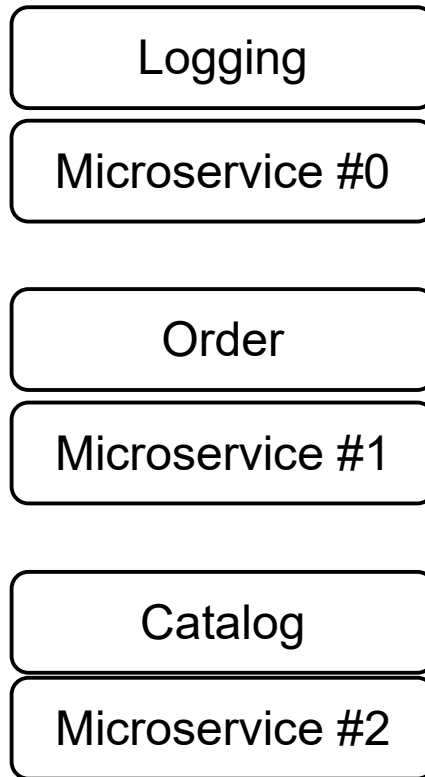


e.g., Amazon website

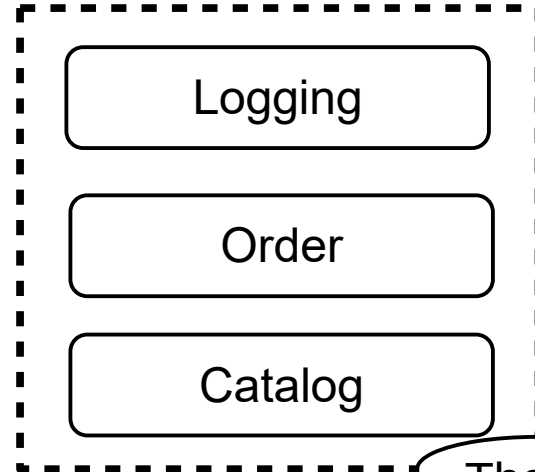


Monolithic Application

Decomposed Application
into Services



e.g., Amazon website

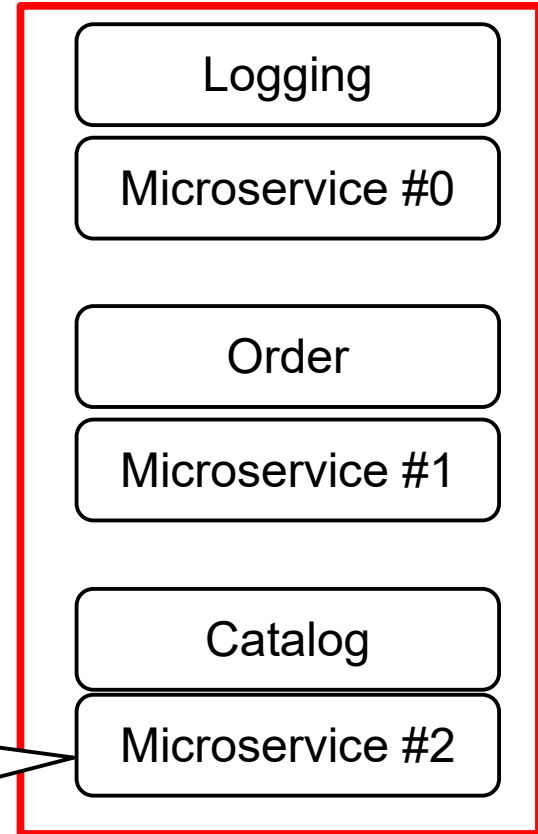


Monolithic Application

Decomposed Application
into Services



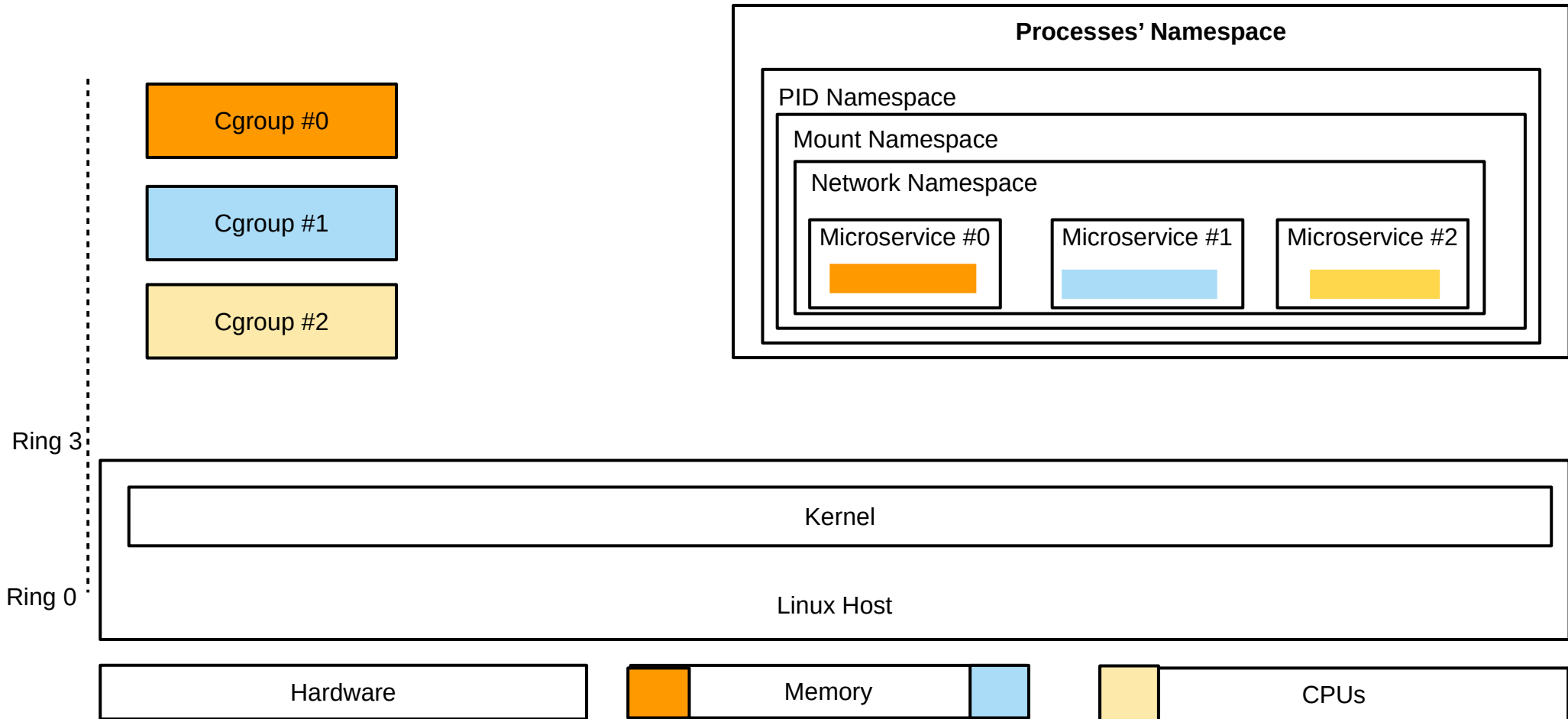
The Cloud provider
deals with deployment,
i.e., performance and isolation

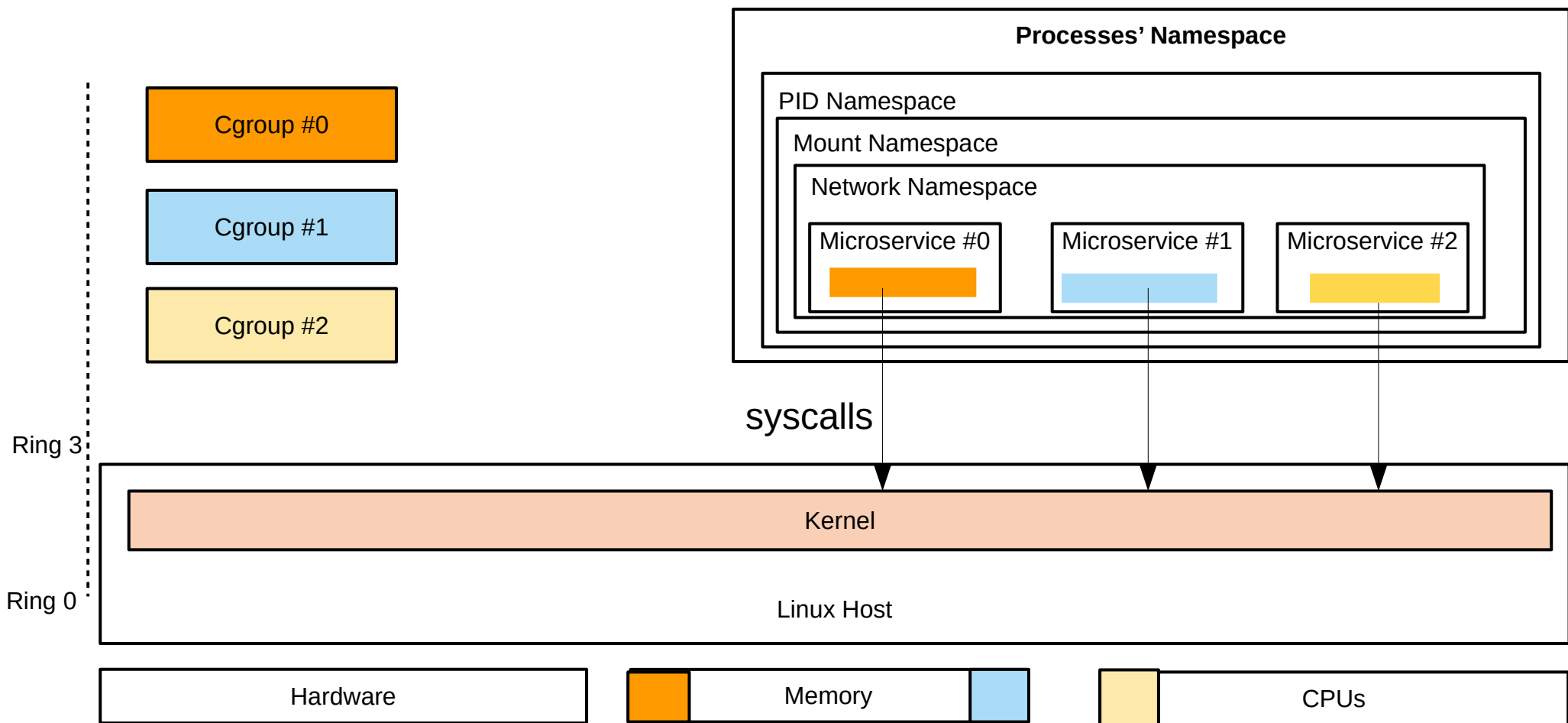


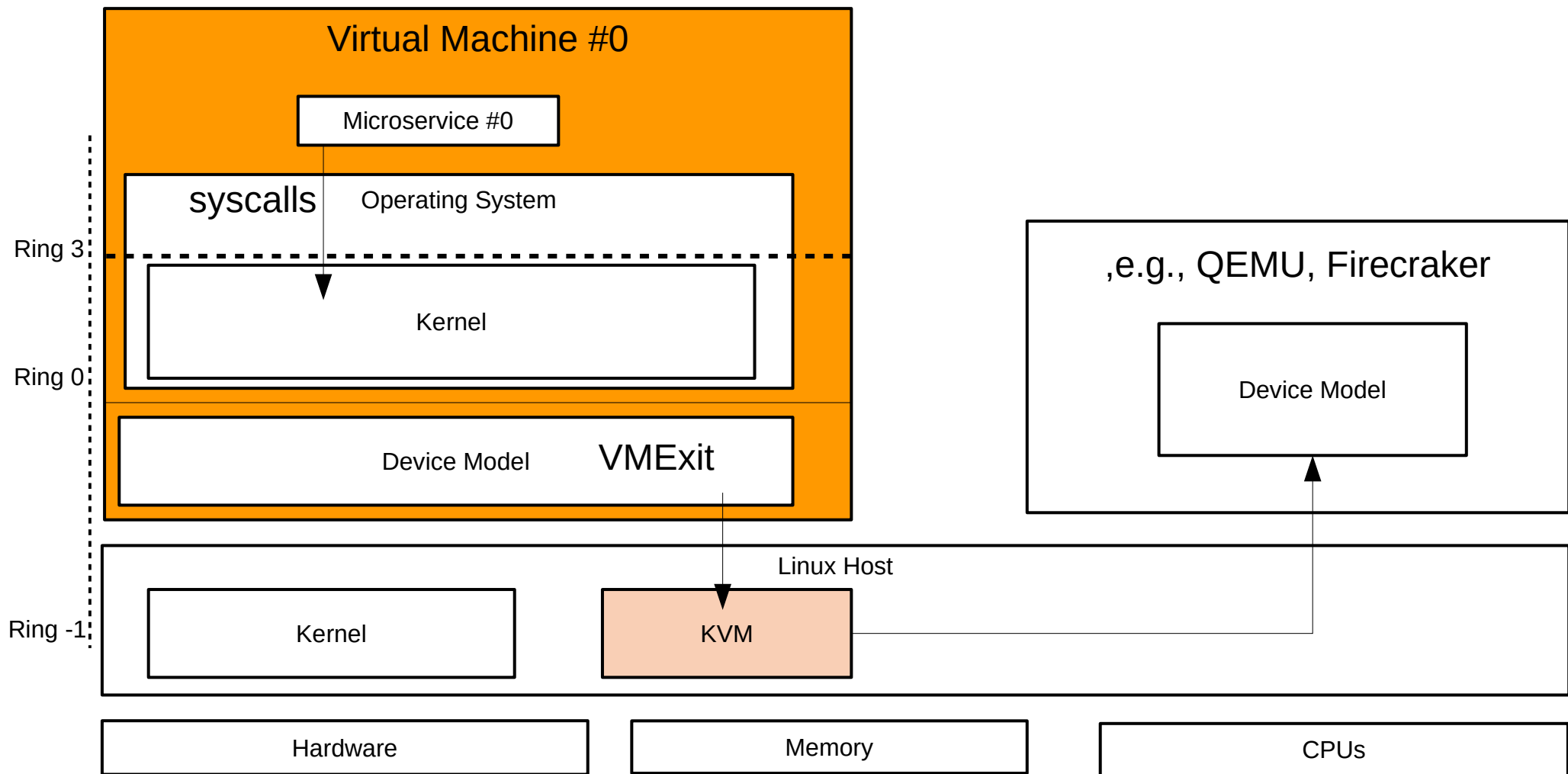
Deployed as server-less

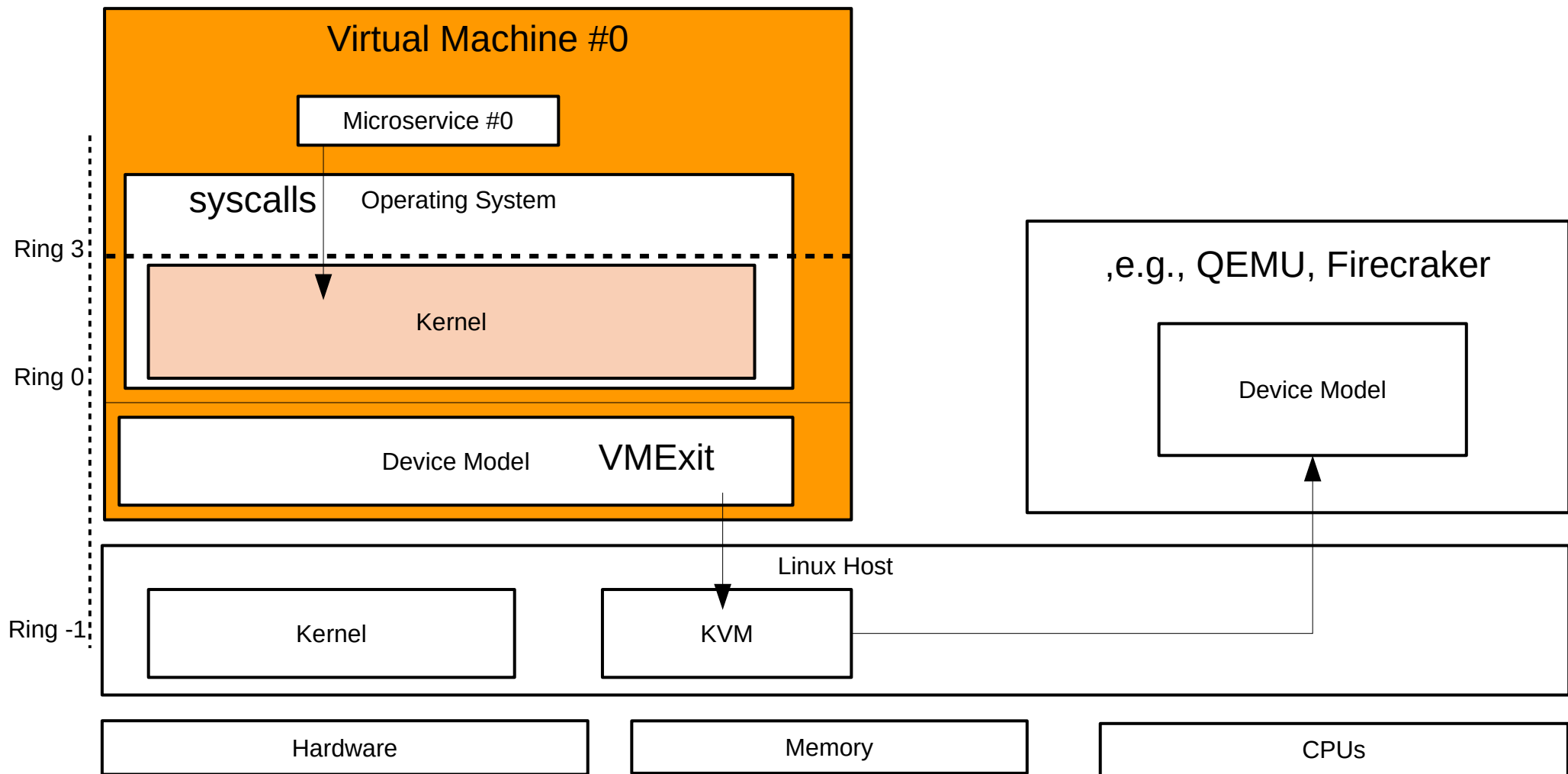
How server-less applications can be deployed?

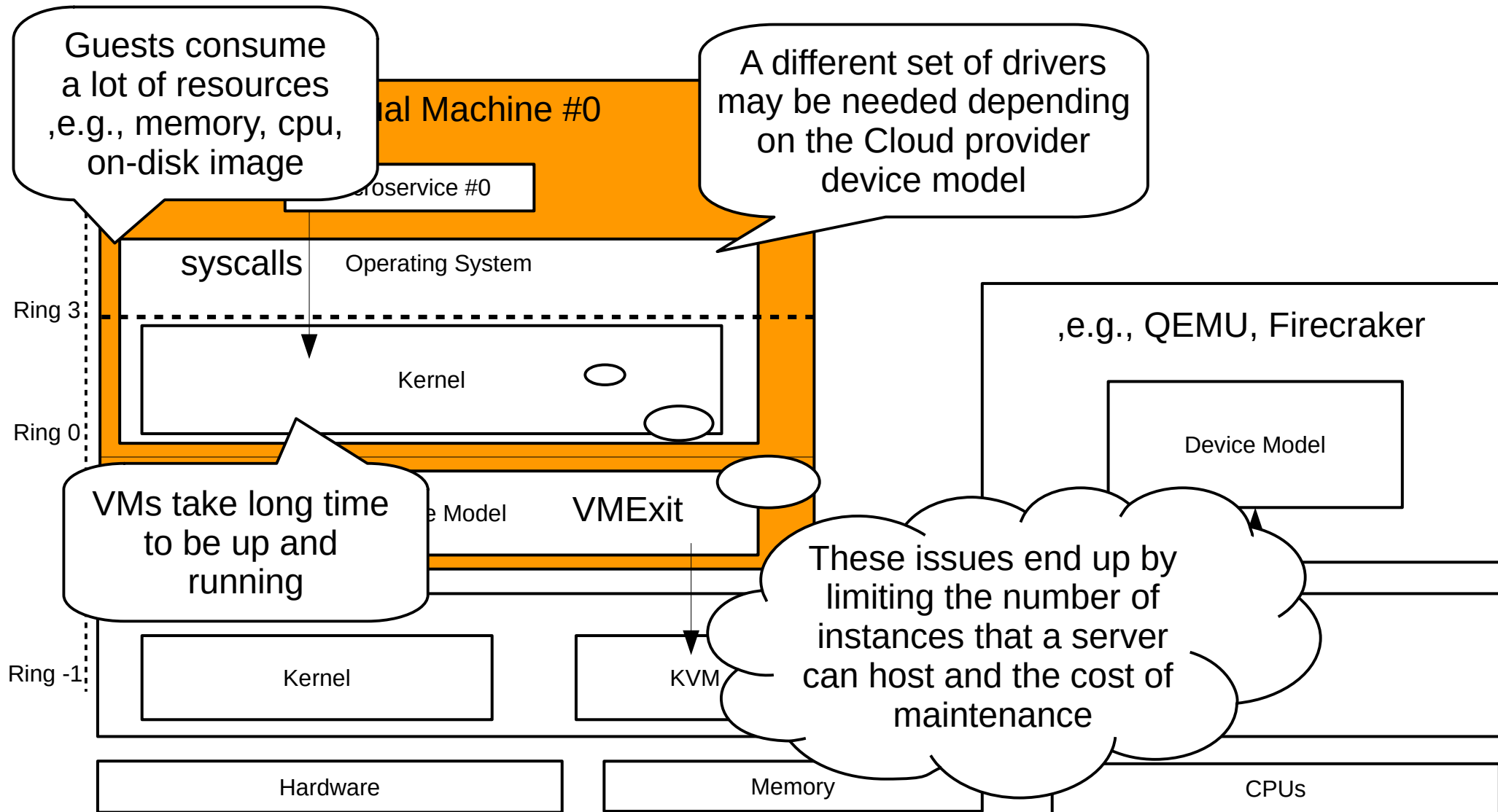
- By using containers, i.e., software-based virtualization
- By using VMs, i.e., hardware-based virtualization
 - General Purpose OS
 - unikernel
- These mechanisms are chosen based on a trade-off between performance and security

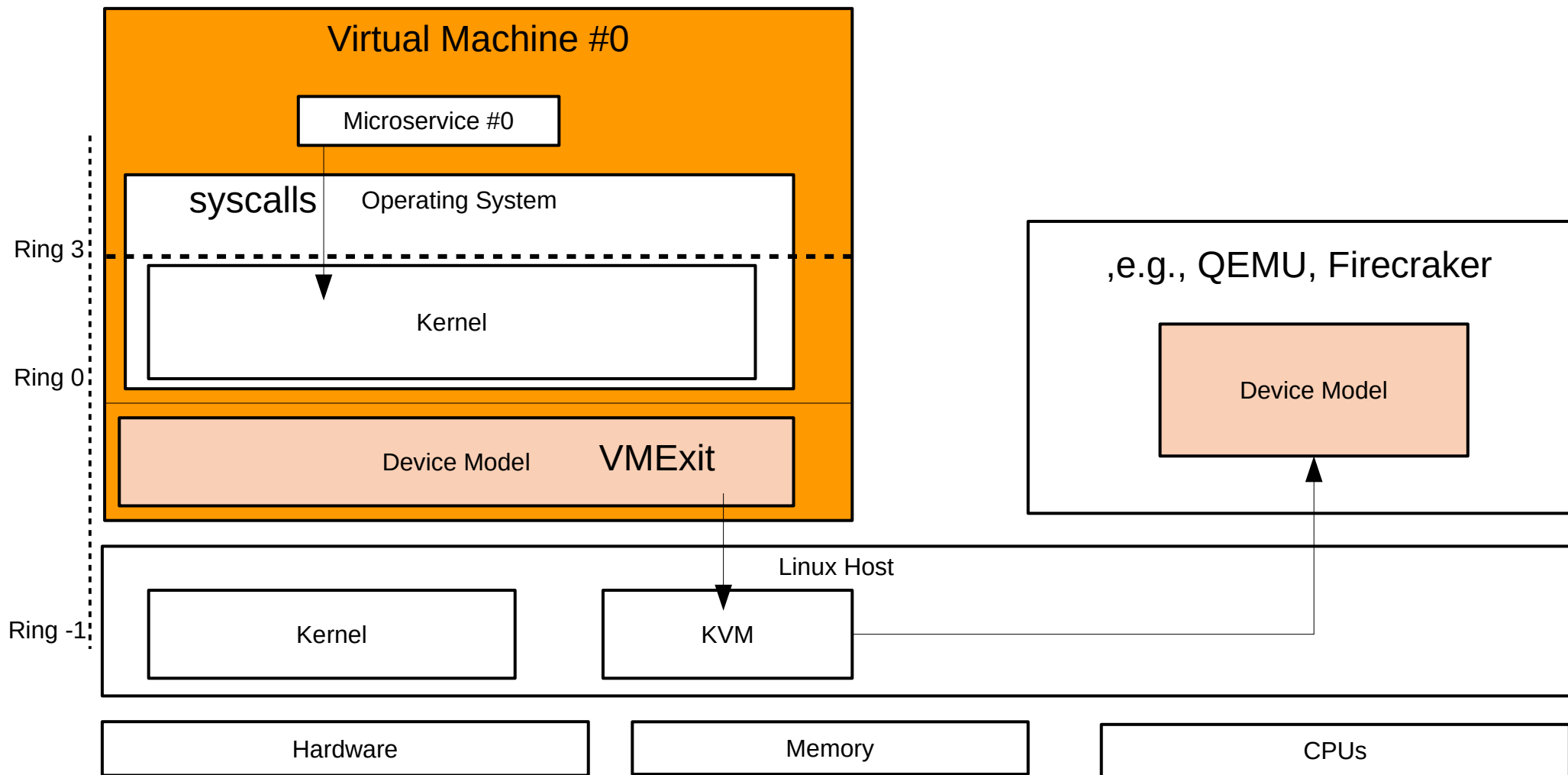


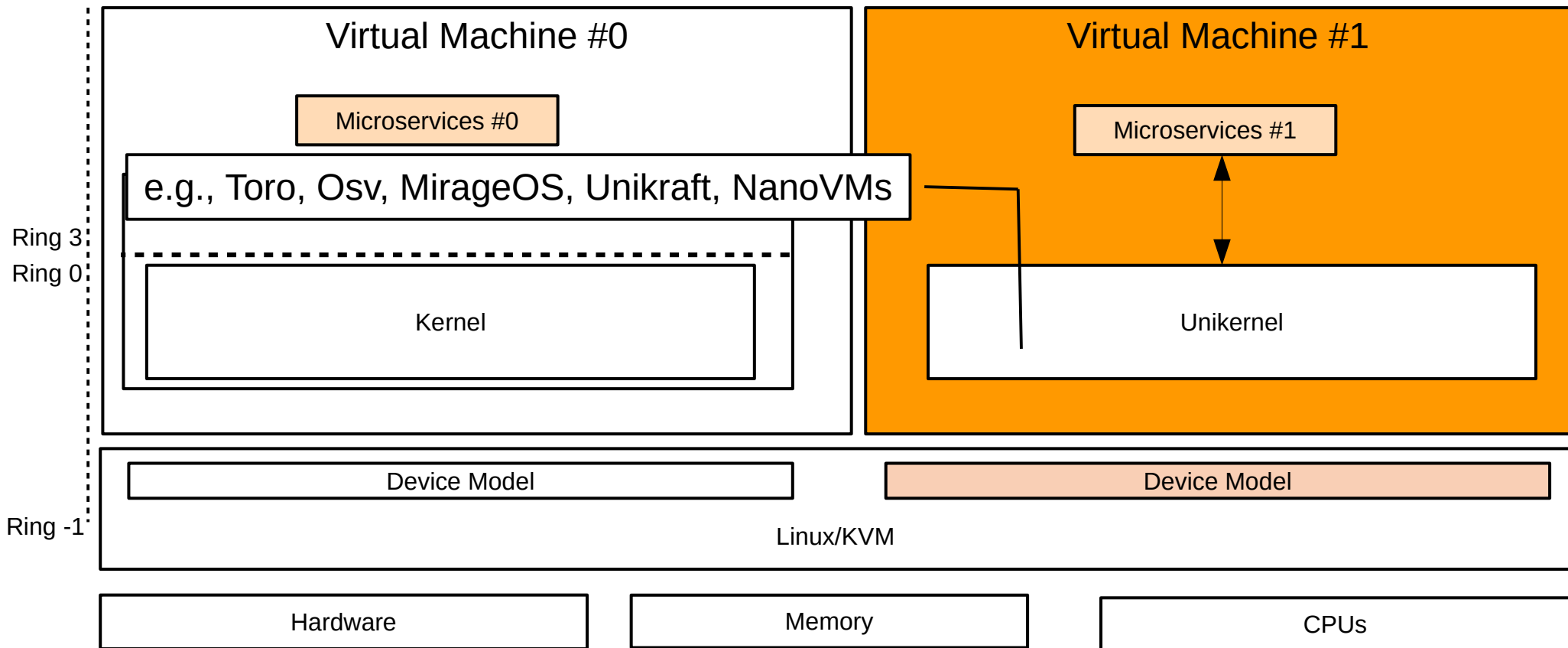






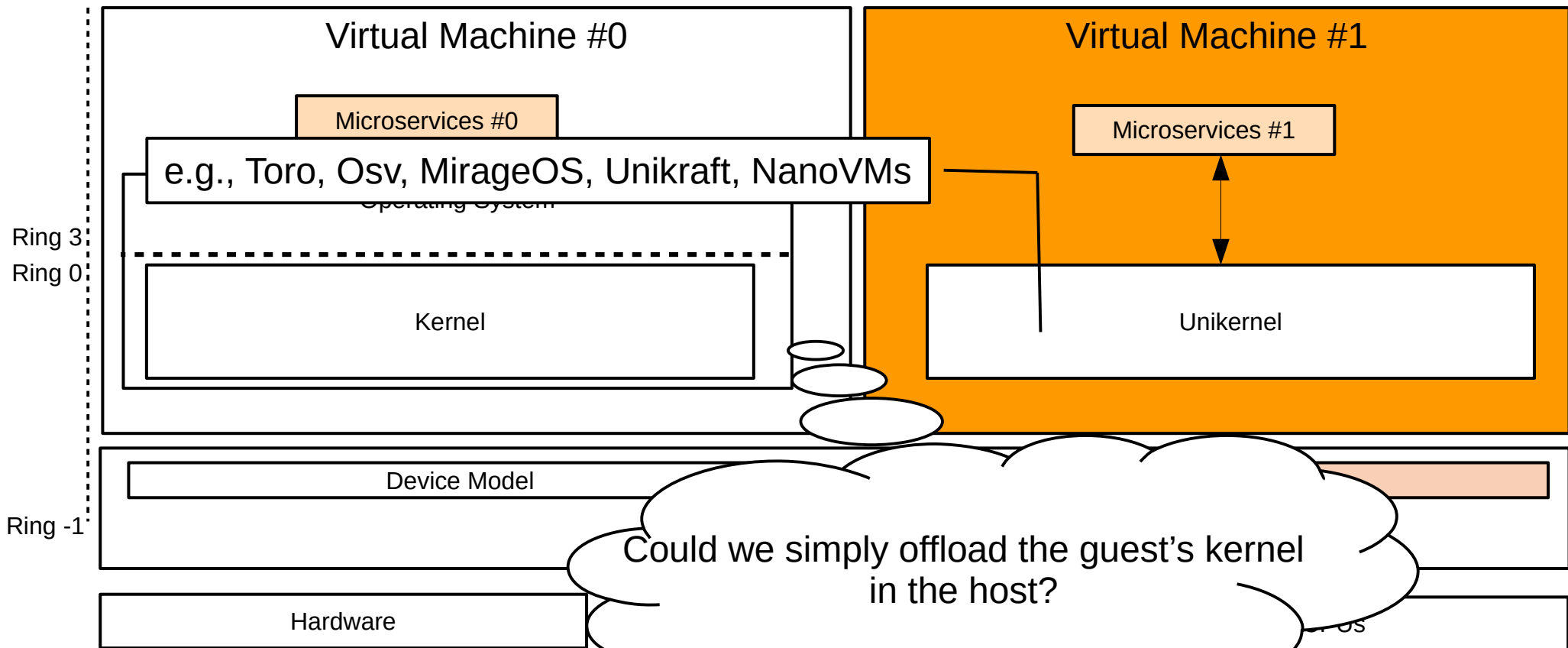






[2] "Unikernels: library operating systems for the cloud", Madhavapeddy et al., 2013

[3] "Unikernels: the next stage of Linux's dominance", Ali Raza et al., 2019

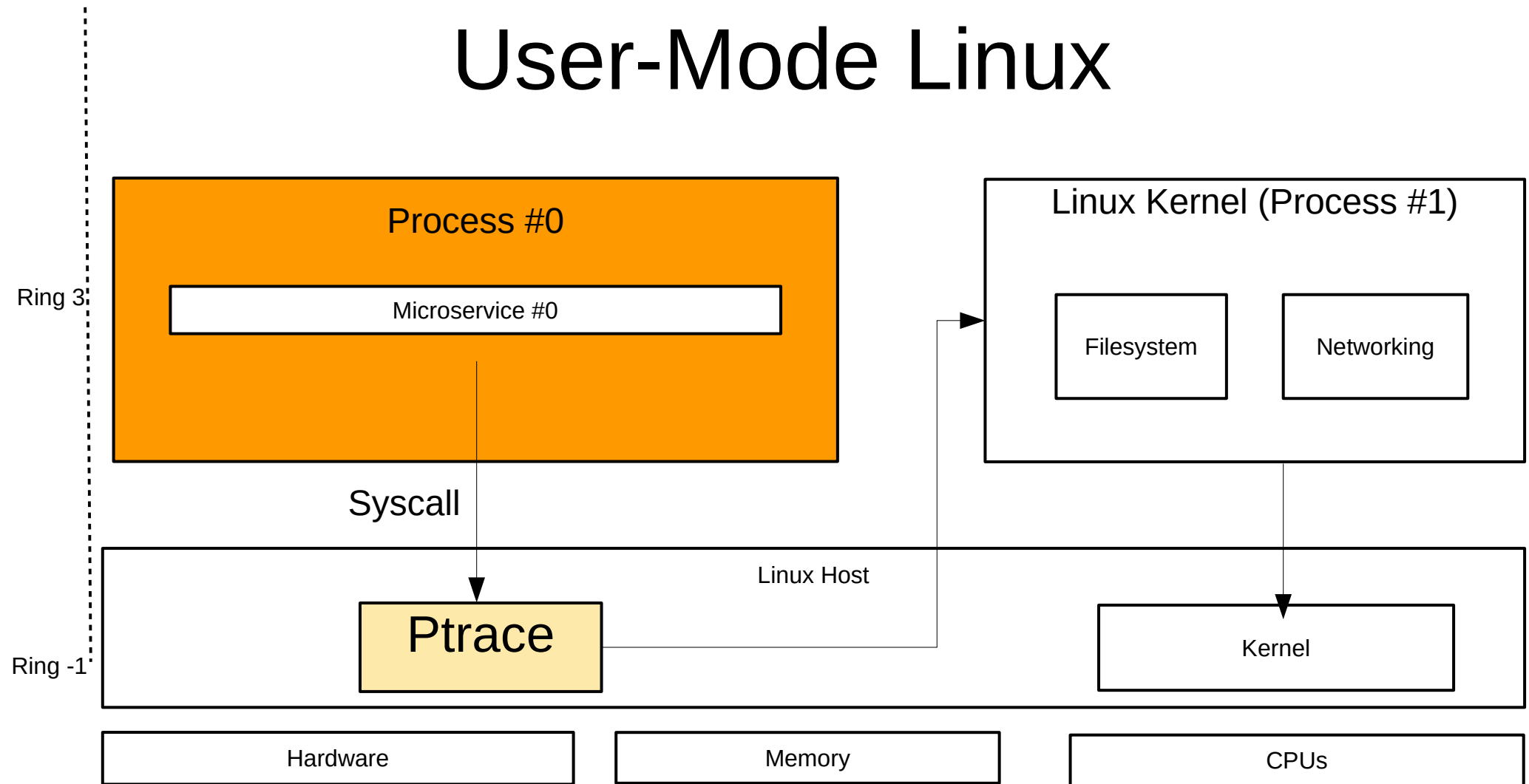


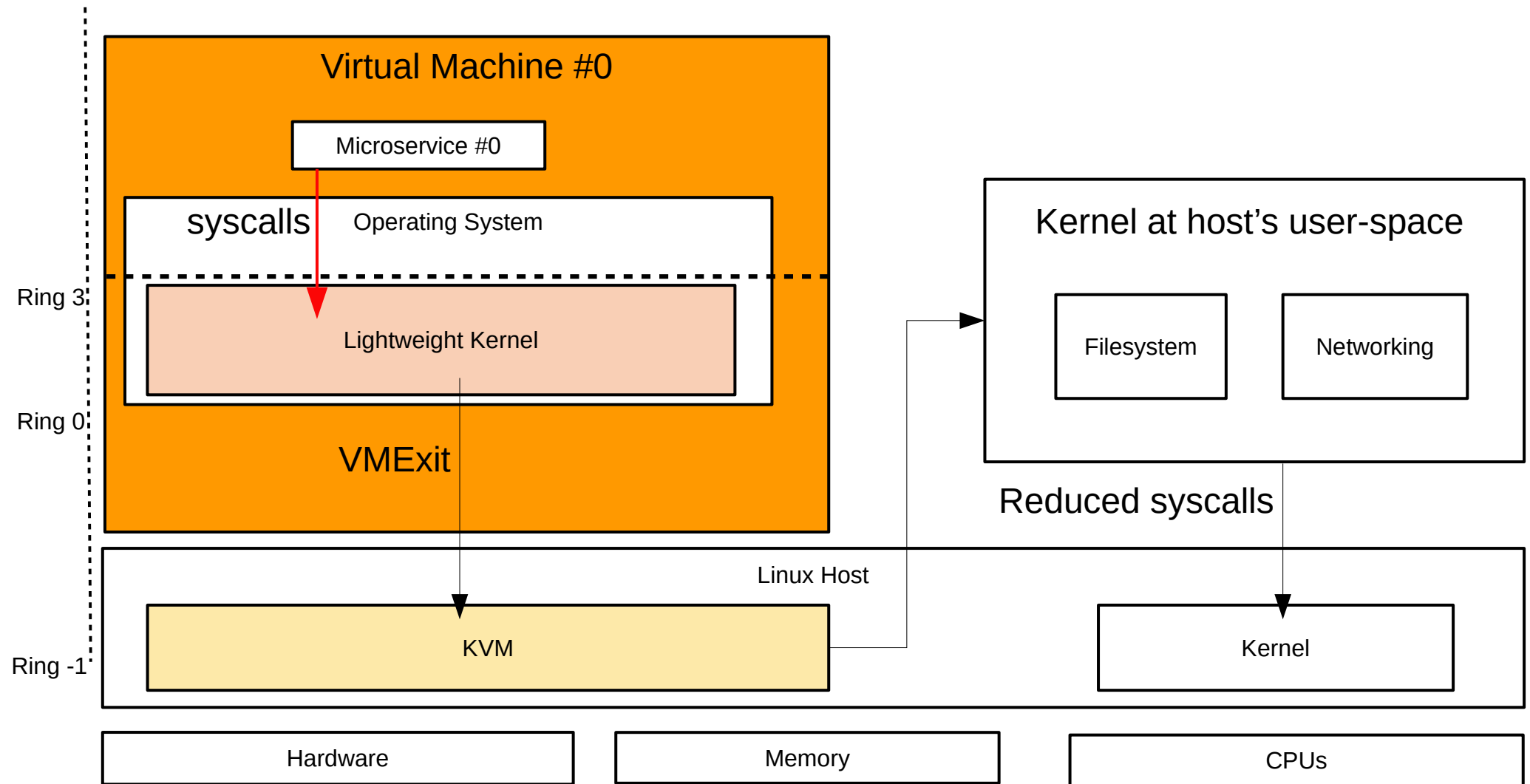
- [2] "Unikernels for the cloud", Madhavapeddy et al., 2013
[3] "Unikernels: the next stage of Linux's dominance", Ali Raza et al., 2019

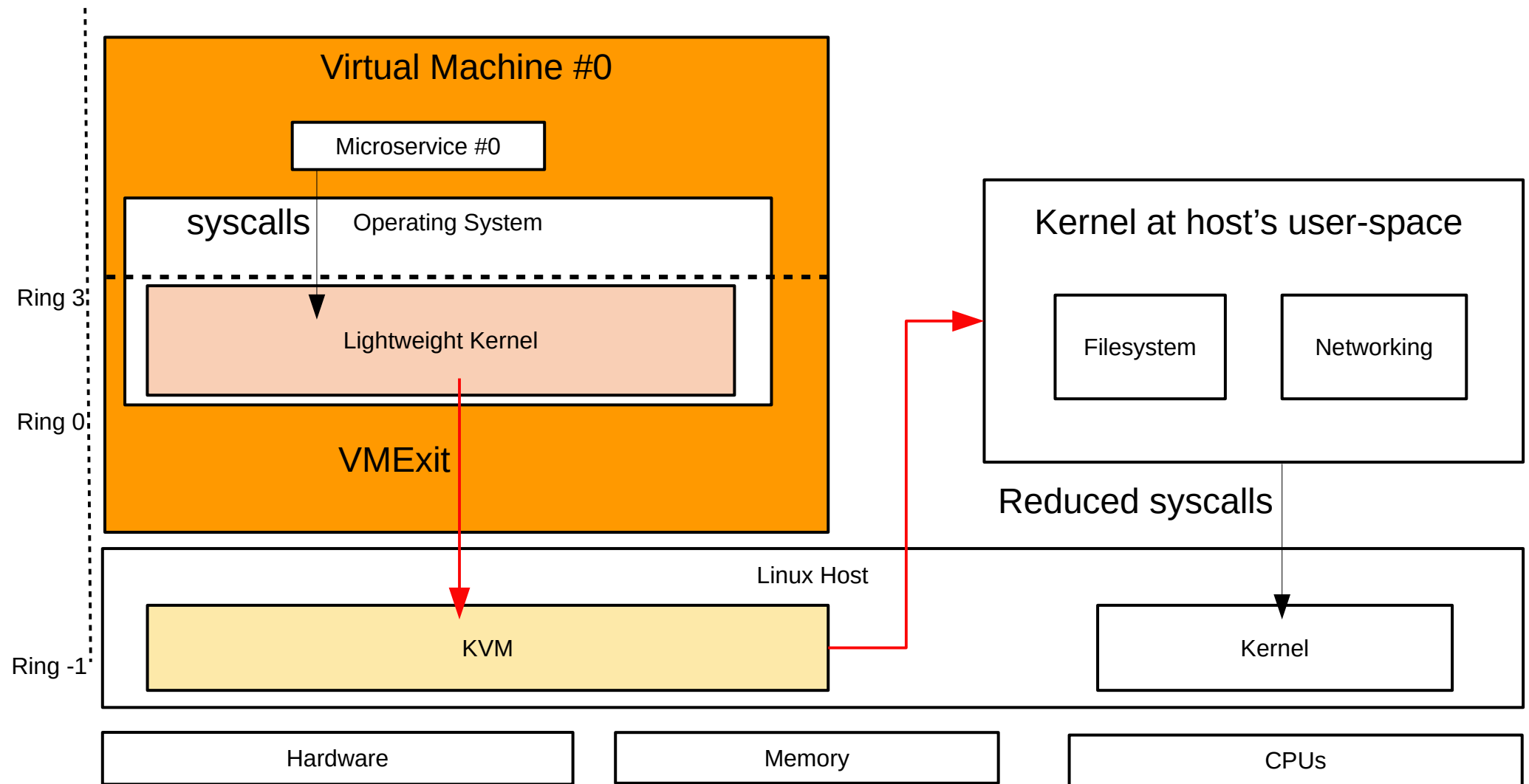
What is a kernel in user's space?

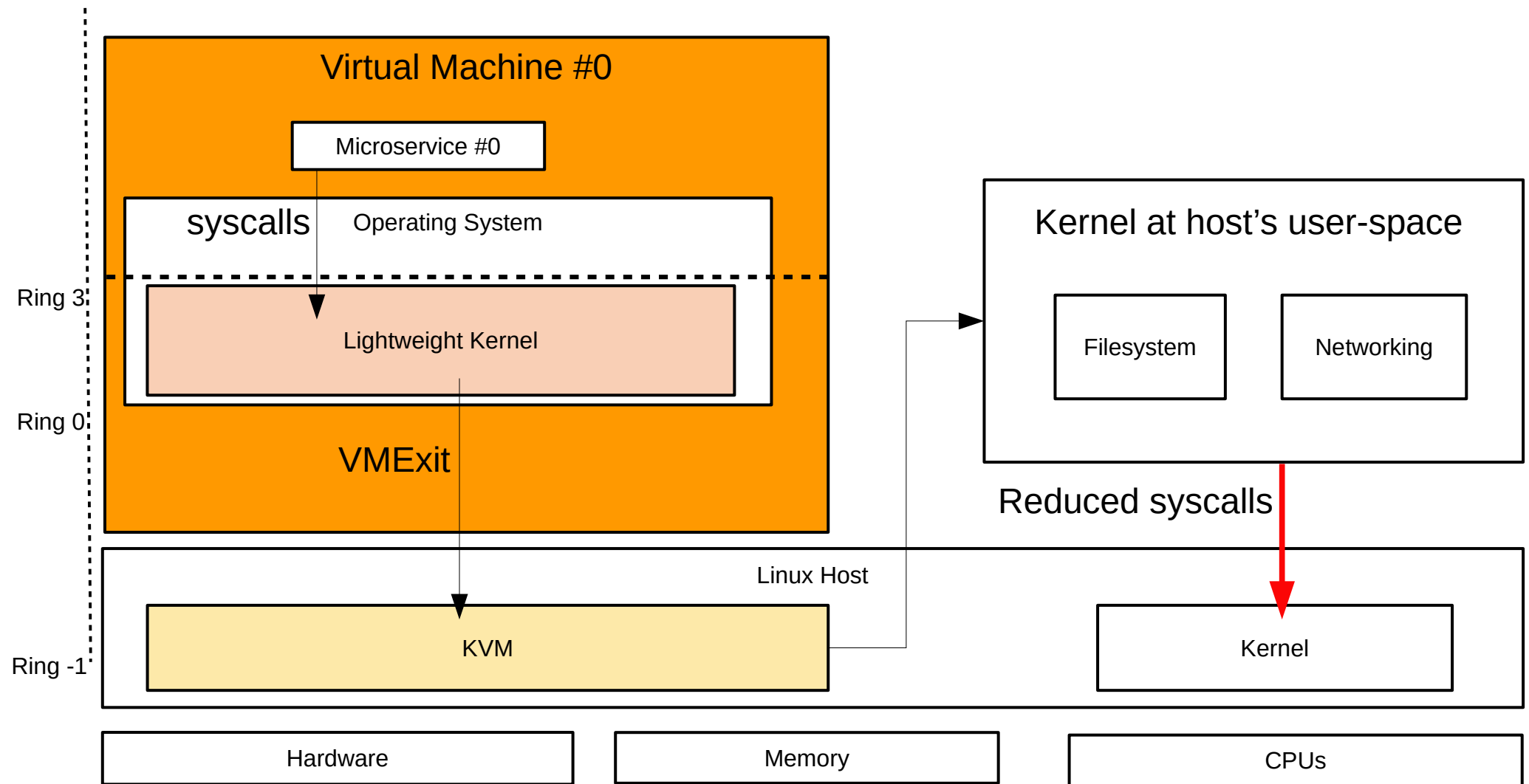
- It is a way to offload guest's kernel in the host
- It allows a guest application to run without a device model thus reducing the attack surface of the host
- It prevents the host to be exposed by emulating kernel services.
- Approaches:
 - User-Mode Linux
 - Gvisor
 - ToroV

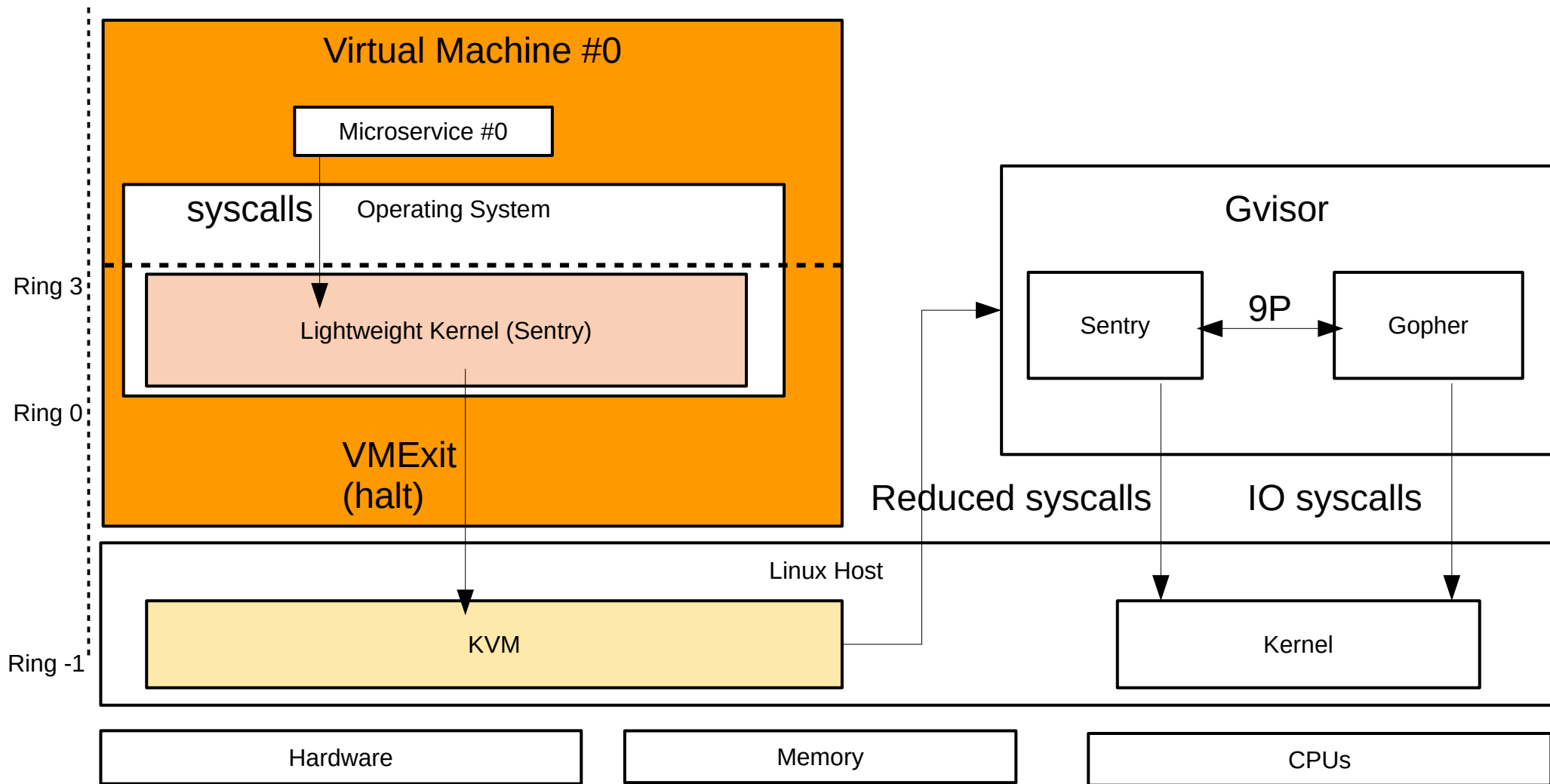
User-Mode Linux

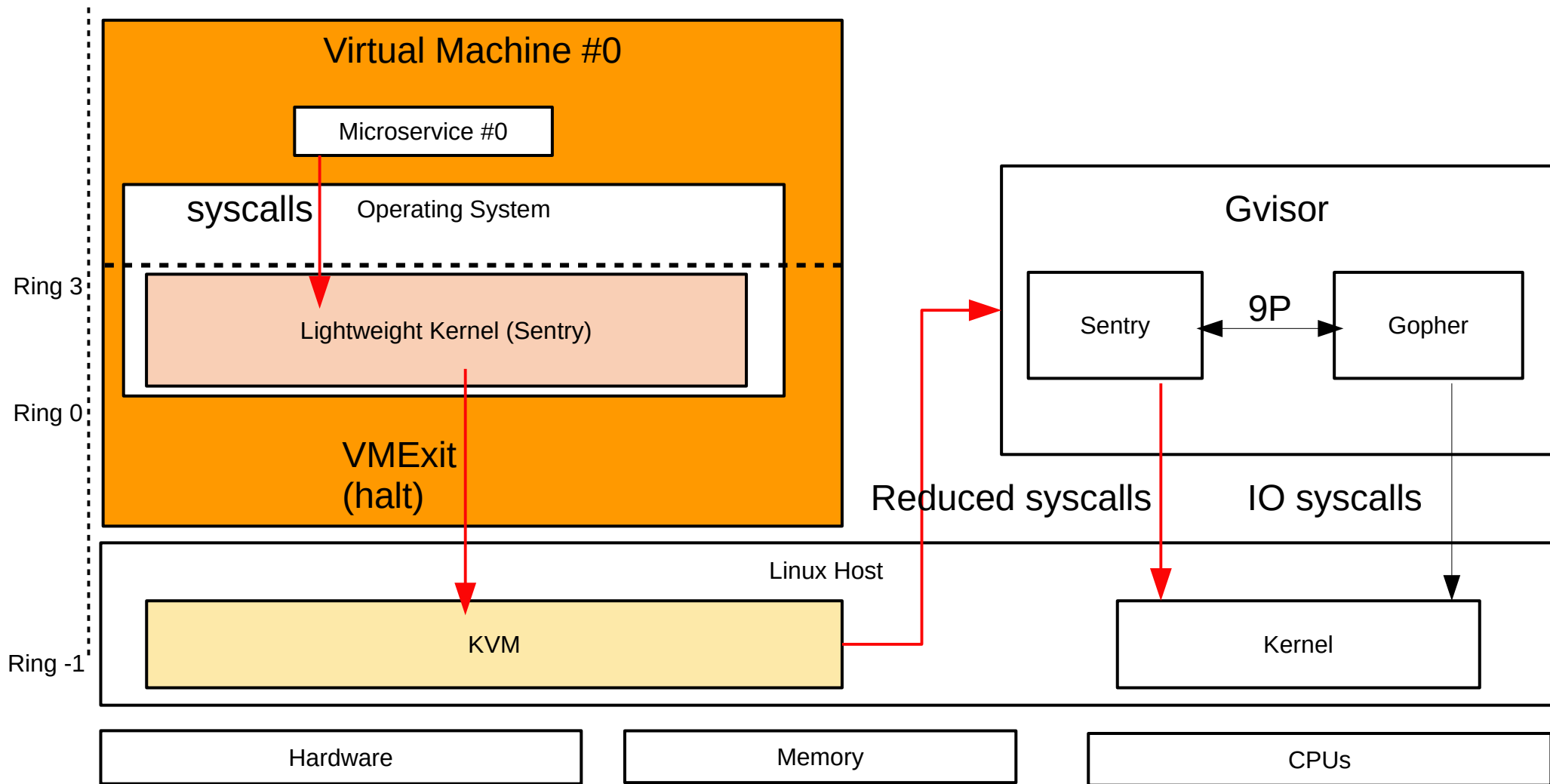






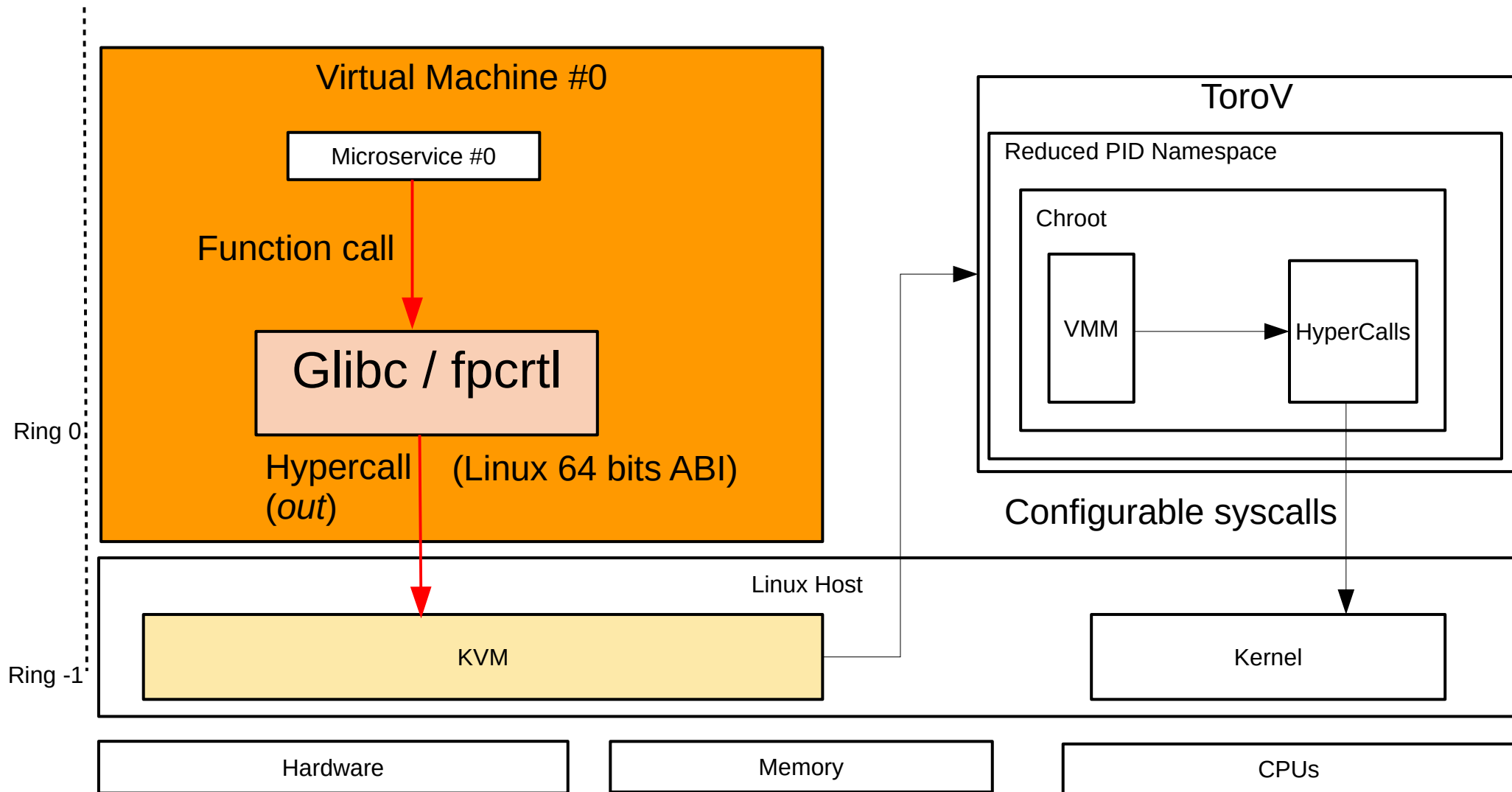


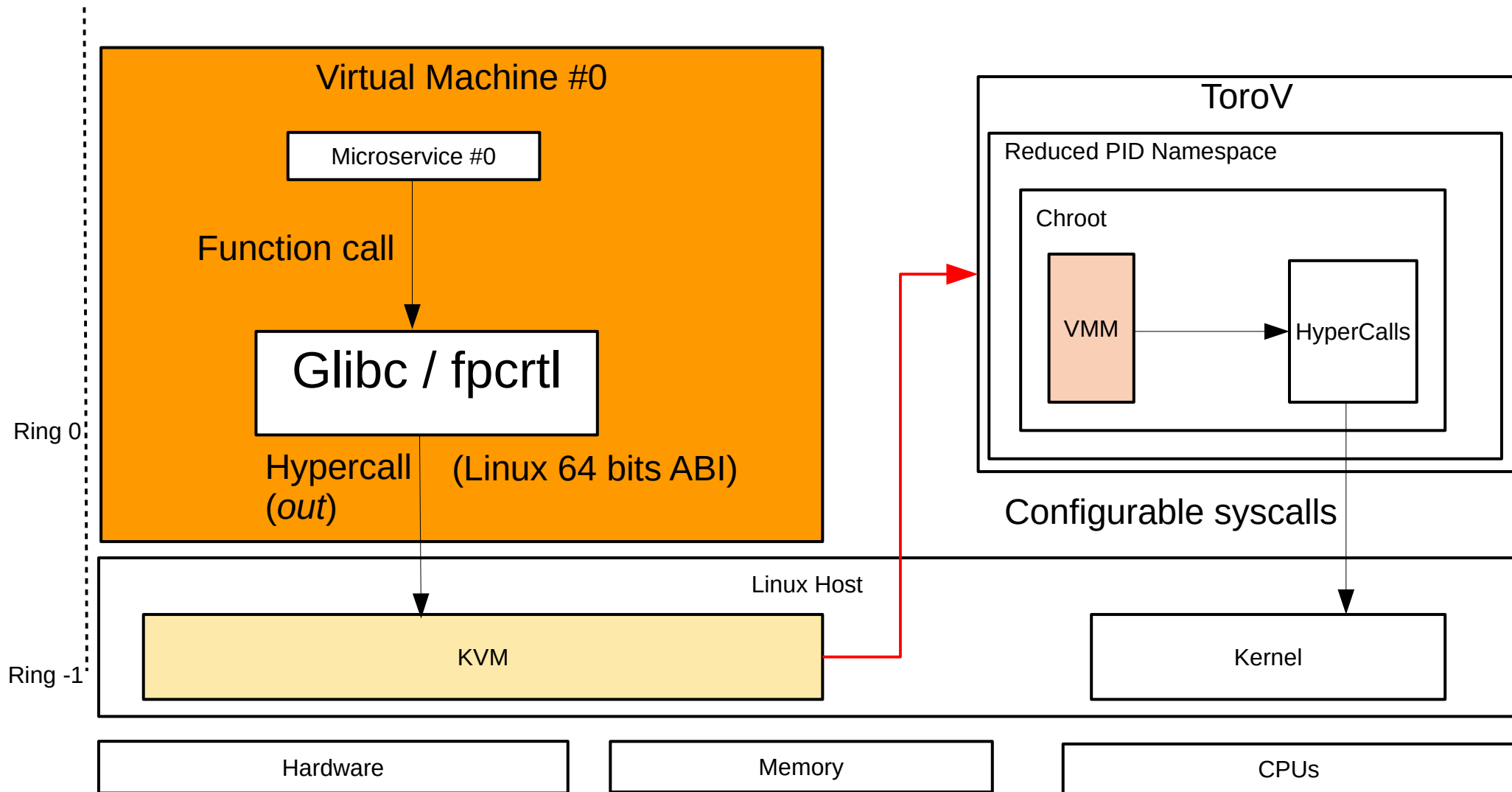


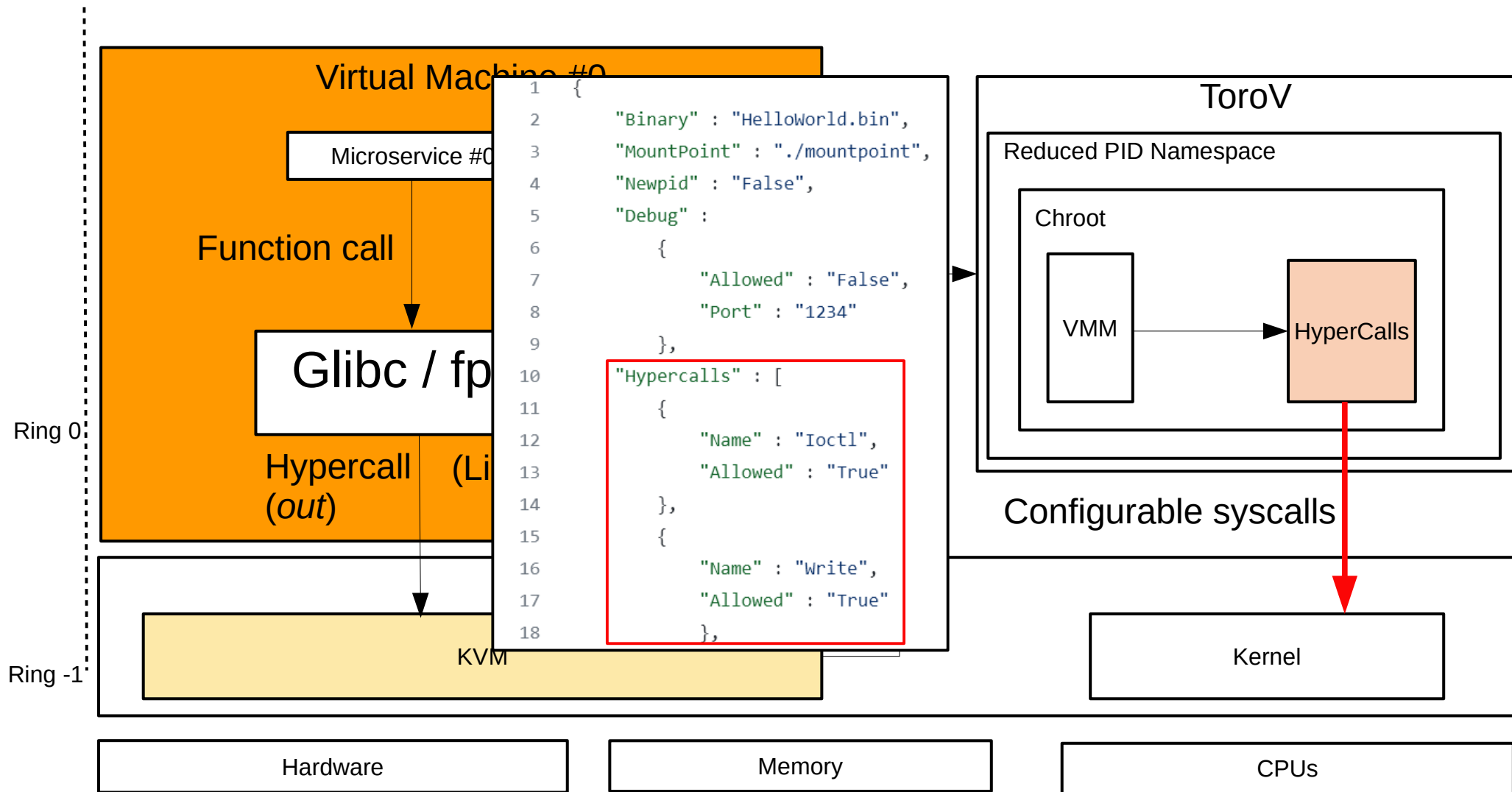


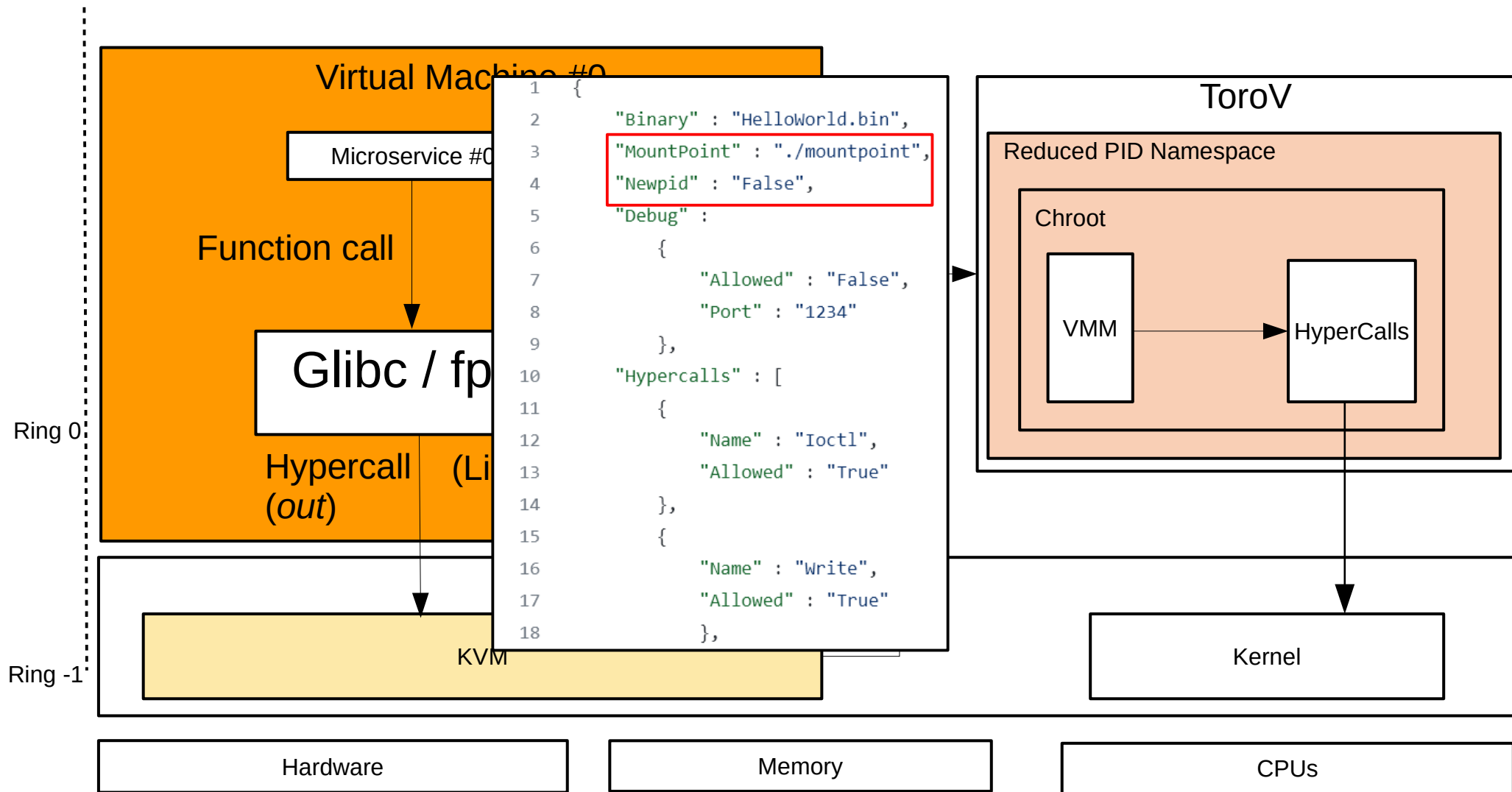
What is ToroV?

- It is a minimalist kernel in user space in which syscalls from the guest are forwarded to the host.
- It allows the user to configure what syscalls are allowed per application.
- It provides a modified stdlib that the user's application must be compiled within.
- It exposes a POSIX API based on hypercalls to the guest.
- It runs as a containerized process to reduce host attack surface.
- It allows the user to debug guest's applications by simply using a GDB.

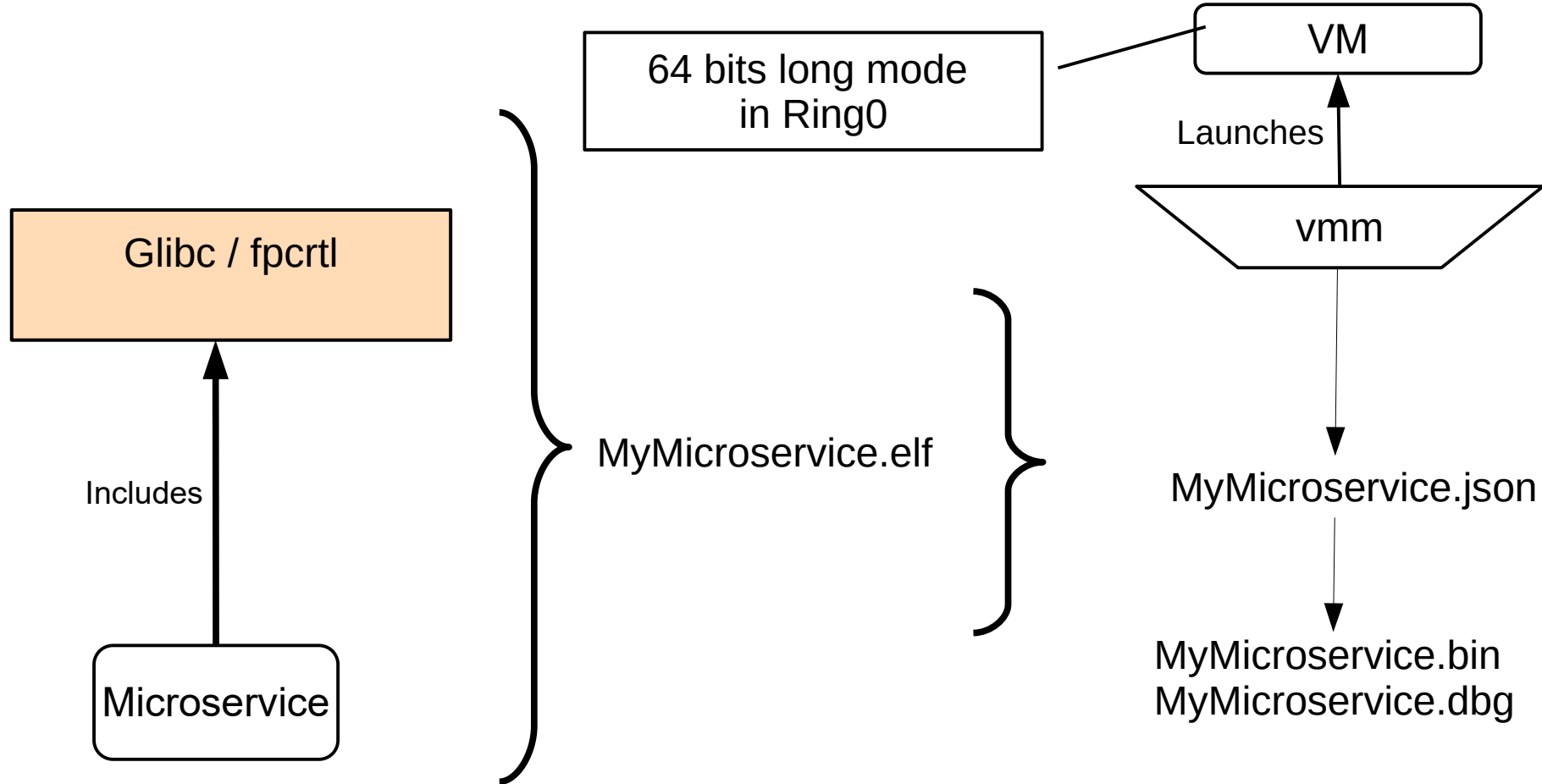




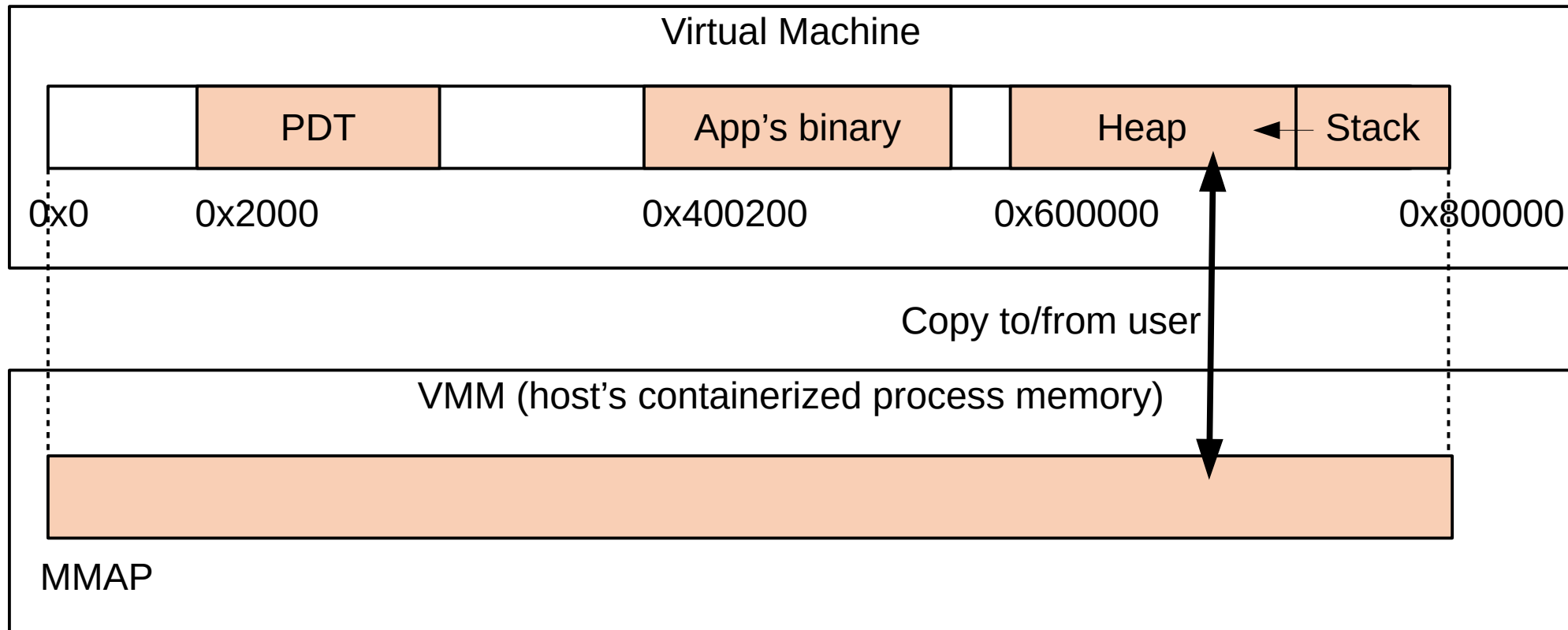




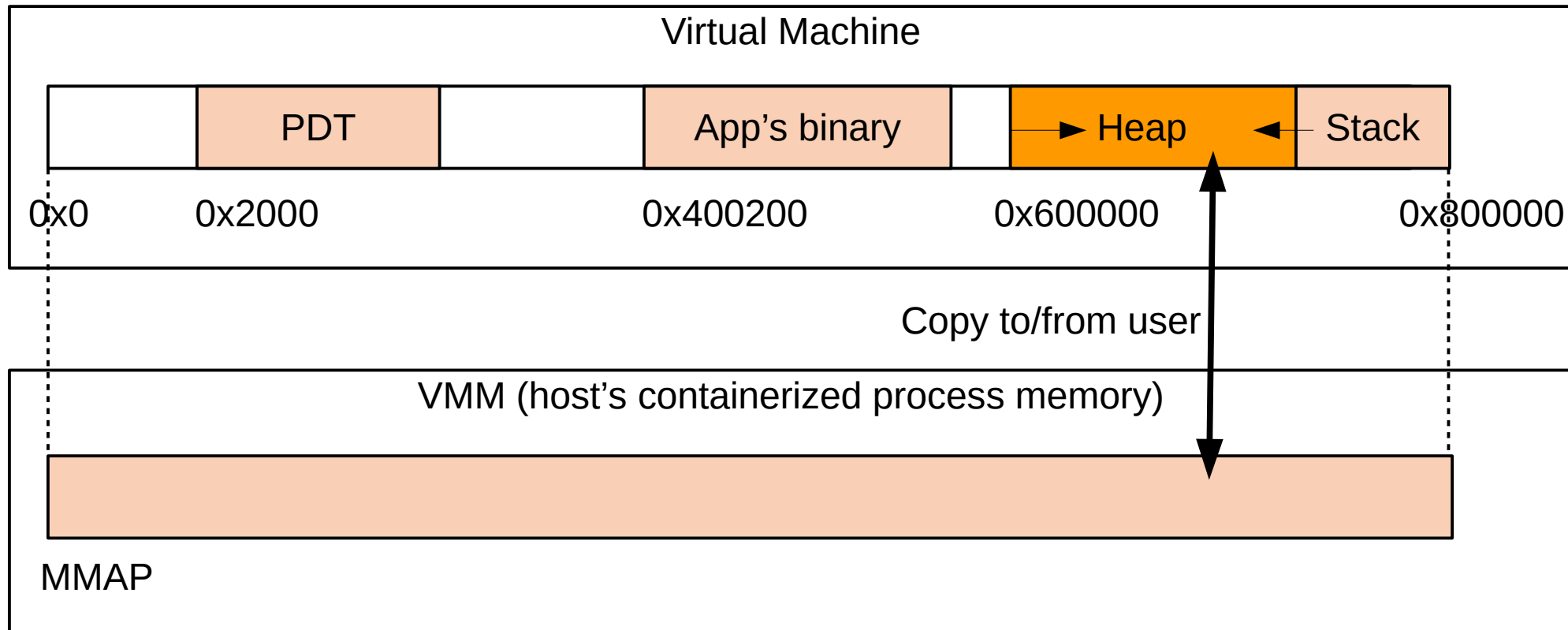
Compilation Process



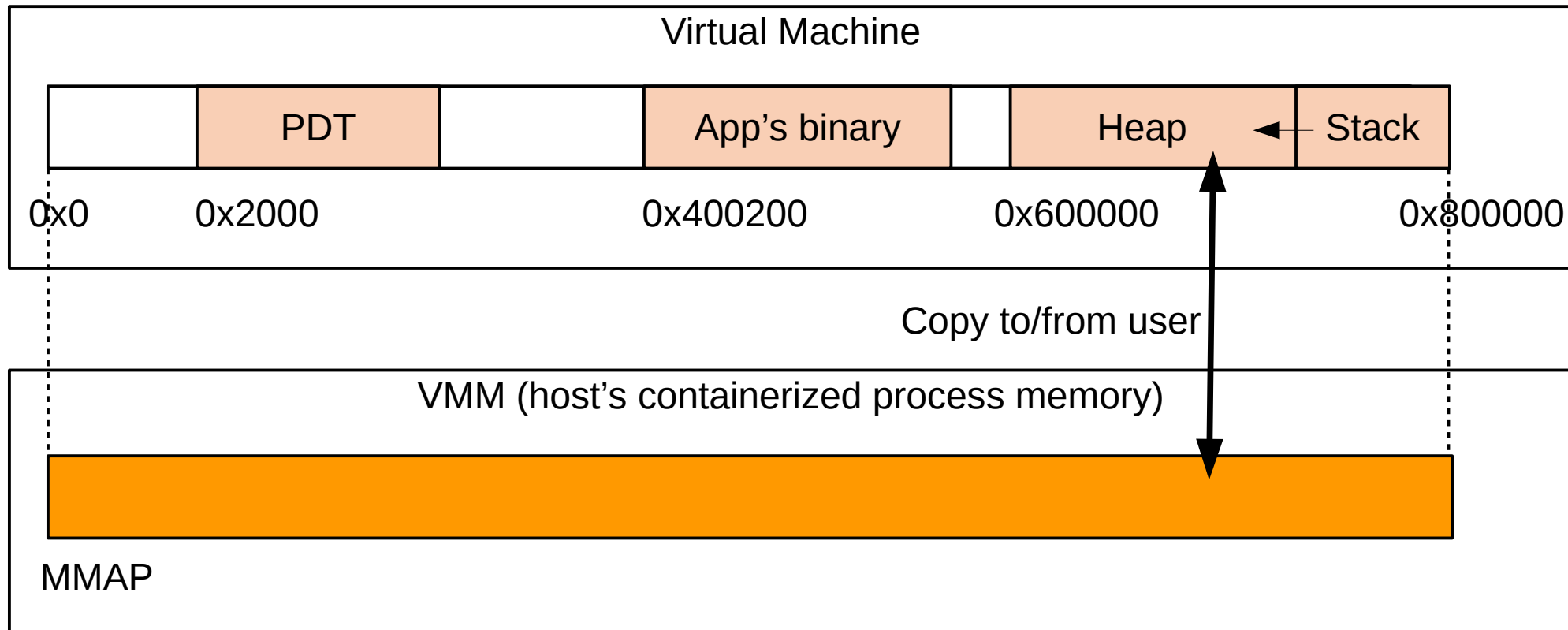
Memory Layout



Memory Layout



Memory Layout



HelloWorld example

```
{
  "Binary": "HelloWorld.bin",
  "MountPoint":
    "./mountpoint",
  "Newpid": "False",
  "Debug":
    {
      "Allowed": "False",
      "Port": "1234"
    },
  "Hypercalls": [
    {
      "Name": "Ioctl",
      "Allowed": "True"
    },
    {
      "Name": "Write",
      "Allowed": "True"
    },
    {
      "Name": "Getrlimit",
      "Allowed": "True"
    }
  ]
}
```

```
strace -f ../../src/vmm/vmm helloworld.json
```

```
ioctl(5, KVM_CREATE_VCPU, 0) = 6
mmap(NULL, 12288, PROT_READ|PROT_WRITE, MAP_SHARED, 6, 0) = 0x7f8908125000
ioctl(6, KVM_GET_SREGS, {cs={base=0xffff0000, limit=65535, selector=61440, type=11, present=1,
dpl=0, db=0, s=1, l=0, g=0, avl=0}, ...}) = 0
ioctl(6, KVM_SET_SREGS, {cs={base=0, limit=4294967295, selector=8, type=11, present=1, dpl=0,
db=0, s=1, l=1, g=1, avl=0}, ...}) = 0
ioctl(6, KVM_SET_REGS, {rax=0, ..., rsp=0x7fff00, rbp=0, ..., rip=0x400200, rflags=0x2}) = 0
clone(child_stack=0x513a30, flags=CLONE_VM|CLONE_FILES|CLONE_NEWNS|SIGCHLD) =
16528
wait4(16528, strace: Process 16528 attached
<unfinished ...>
[pid 16528] chroot("./mountpoint") = -1 ENOENT (No such file or directory)
[pid 16528] chdir("/") = 0
[pid 16528] ioctl(6, KVM_RUN, 0) = 0
[pid 16528] ioctl(6, KVM_GET_REGS, {rax=0x1, ..., rsp=0x7ffe28, rbp=0x7ffe38, ..., rip=0x4005c5,
rflags=0x2}) = 0
[pid 16528] write(2, "Hello World, I am ToroV!\n", 25Hello World, I am ToroV!
) = 25
[pid 16528] ioctl(6, KVM_SET_REGS, {rax=0x19, ..., rsp=0x7ffe28, rbp=0x7ffe38, ..., rip=0x4005c5,
rflags=0x2}) = 0
[pid 16528] ioctl(6, KVM_RUN, 0) = 0
```

HelloWorld example

- ~ 1.5 MB of memory (top)
- ~ 7 ms (median)
- Write() syscall ~ 0.10 ms ~ x10 slower (0.012 ms)

Future Work

- Work on Glibc and other languages like Go or Rust
- Research about how ToroV compares with seccomp
- Enable the use of binaries without recompilation by replacing “syscall” opcode by “out” opcode thus starting the program as-is
- Port the whole project to Rust
- Replace the current syscall mechanism for an asynchronous mechanism, .e.g, virtio device for syscalls.
- Enable that different components handles different syscalls, e.g., SOA

Resources

- GitHub repository at <https://github.com/torokernel/torov>
- “Debugging applications that run as VM by using GDB”, <https://youtu.be/QC8pYtMOWe4>
- “Using ToroV to isolate an app by using virtualization and containerization technologies”, <https://youtu.be/YDpE8jIwVPA>
- “Simple HelloWorld in C in ToroV”, https://youtu.be/E_bQPc64WIM
- “Simple Echo server by relying on POSIX hypercalls”, <https://youtu.be/aJpcmZhDqMw>

Thanks!



www.torokernel.io