**USENIX ATC'22 Best Paper Award!**

# Co-opting Linux Processes for High-Performance Network Simulation

Rob Jansen, U.S. Naval Research Laboratory
Jim Newsome, Tor Project
Ryan Wails, Georgetown University & U.S. Naval Research Laboratory
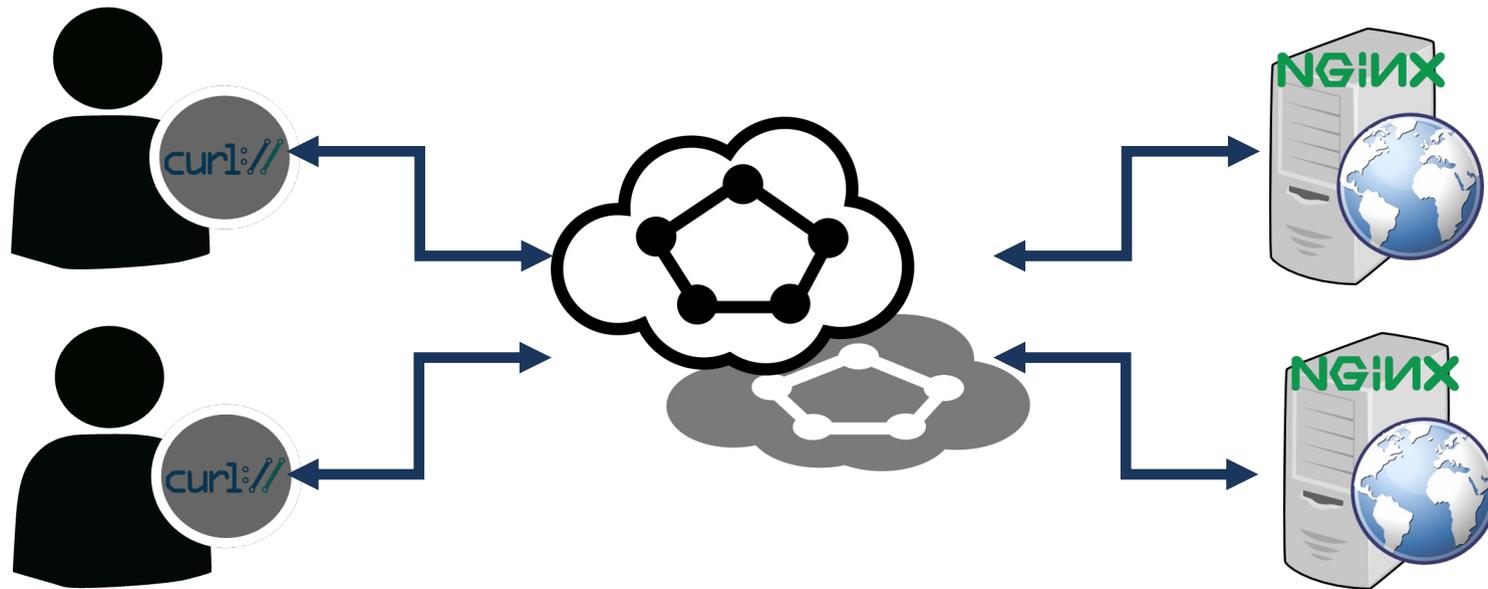
**Rob Jansen, Ph.D.**
Computer Security Research Scientist
Center for High Assurance Computer Systems
U.S. Naval Research Laboratory

USENIX Annual Technical Conference
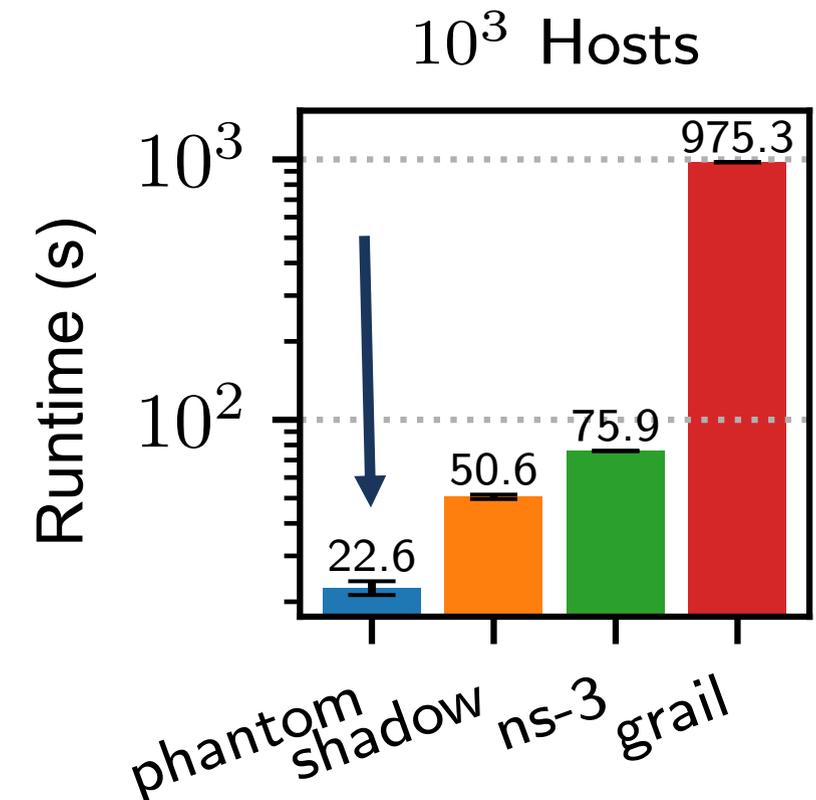Carlsbad, CA, USA
July 11th, 2022

## Designed a new, hybrid network simulator/emulator

- co-opts Linux processes into a discrete-event network simulation that emulates kernel functionality
- enables large-scale, distributed system experiments

- **2.3x** faster than Shadow v1
- **3.4x** faster than NS-3
- **43x** faster than gRaIL [ToN'19]



$10^3$ Hosts

- Merged into the open-source Shadow project and synonymous with **Shadow v2**
- Artifacts: https://netsim-atc2022.github.io

# Outline

**motivation**
design
evaluation

- Important properties of test networks:
  - **Controllable**: isolate important factors
  - **Replicable**: identically replicate experiments (determinism)

- Requirements for large distributed systems:
  - **Accurate**: directly execute system software (not an abstraction)
  - **Scalable**: decouple from time, computational constraints of host

# Problems with Traditional Approaches

- Simulation (e.g., ns-3)
  - Not realistic: runs abstractions instead of real applications
  - Hard to maintain and can lead to invalid results
- Emulation (e.g., mininet)
  - Not controllable: results will not be identical
  - Not scalable: CPU overload → time distortion

# Problems with Traditional Approaches

- Simulation (e.g., ns-3)
  - Not realistic: runs abstractions instead of real applications
  - Hard to maintain and can lead to invalid results
- Emulation (e.g., mininet)
  - Not controllable: results will not be identical
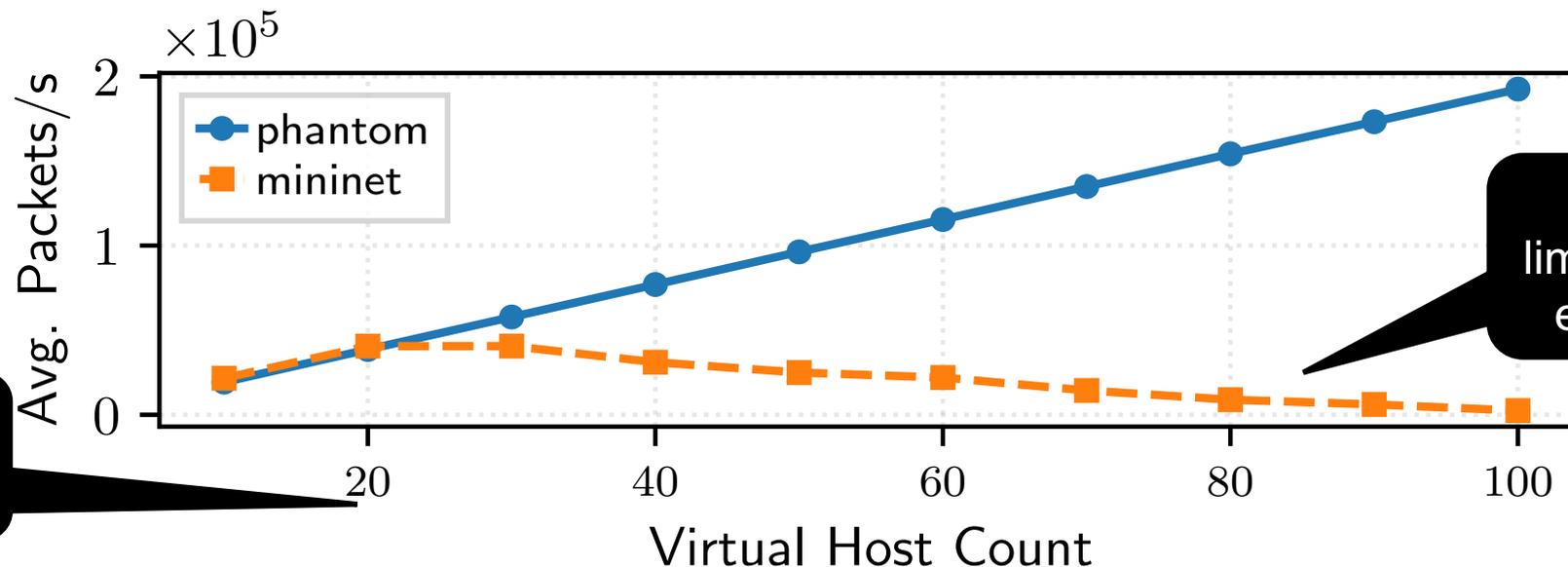  - Not scalable: CPU overload → time distortion

**...ıllns-3**

**Mininet** > sudo mn



CPU overload with >20 hosts

Forwarding capacity limited: fewer packets than expected are forwarded

- Hybrid architecture
  - Network simulation, but directly execute application code
  - Enjoys advantages of both simulation and emulation
  - Best opportunity to scale to large-scale distributed systems

| Architecture | Example Tool | Scalability | Realism | Control |
|---|---|---|---|---|
| Emulation | Mininet | ✗ | ✓ | ✗ |
| Simulation | NS-3 | ✓ | ✗ | ✓ |
| Hybrid | This work | ✓ | ✓ | ✓ |

Executing application code via
plugin (link-map) namespaces

- appid = dlmopen(app.so)
- func = dlsym(appid, "main")
- func()

NS-3-DCE, Shadow

Limitations

- Compatibility (must build PIC/PIE)
- Correctness (intercept libcalls only)
- Maintainability (custom ld, threading)

## Executing application code via plugin (link-map) namespaces

- appid = dlmopen(app.so)
- func = dlsym(appid, "main")
- func()

## NS-3-DCE, Shadow

## Limitations

- Compatibility (must build PIC/PIE)
- Correctness (intercept libcalls only)
- Maintainability (custom ld, threading)

**Table 2:** Application Properties Supported in Hybrid Simulators

| Application Property | Shadow | Phantom |
|---|:---:|:---:|
| Multiple threads (e.g., support for pthreads) | ● | ● |
| Multiple processes (e.g., support for fork) | ◑ | ◑ |
| Not position-independent (i.e., PIC or PIE) | ○ | ● |
| Not dynamically linked to libc | ○ | ● |
| Symbols not exported to dynamic symbol table | ○ | ● |
| System calls made in statically linked code | ○ | ● |
| System calls made in assembly (i.e., avoiding libc) | ○ | ● |
| 100% statically linked (e.g., some go programs) | ○ | ◑ |

○ Does not work in tool or architecture ● Works in tool & architecture
◑ Not implemented in tool (as of writing) but supported by architecture

## Executing code via Linux processes

- fork() + exec() $\rightarrow$ ptrace()
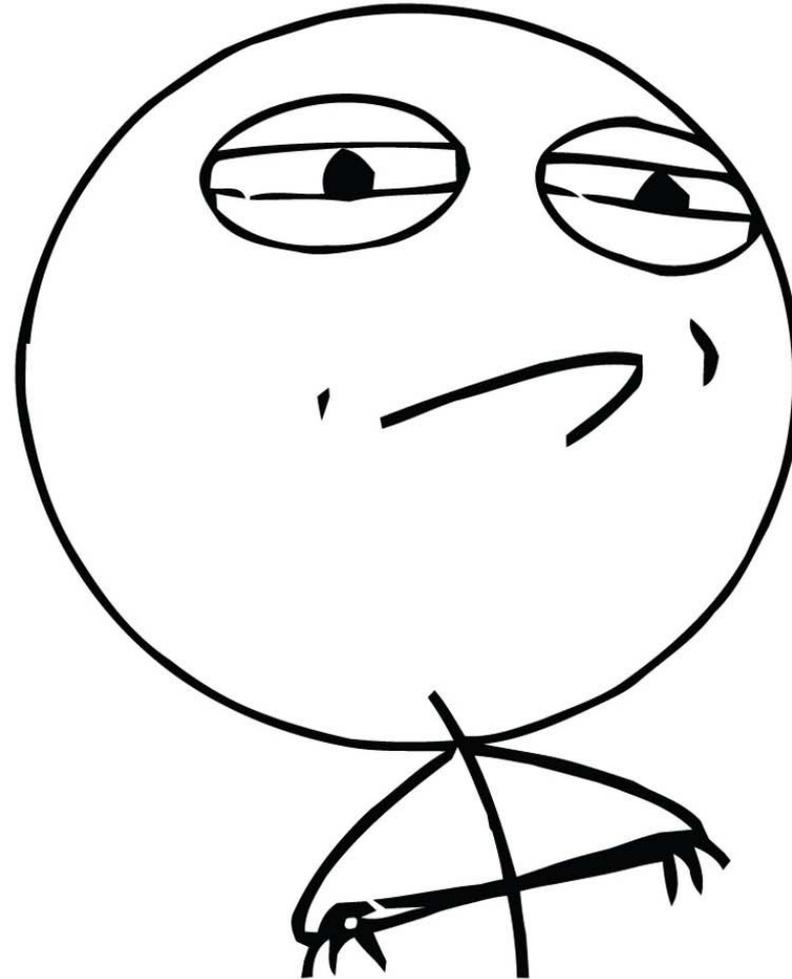
## gRaIL [ToN'19]

## Limitations: ptrace is slow!

- Process creation overhead **quadratic** in total number of processes
- Synchronization: 4 context switches for **every syscall**
- Data transfer: extra syscall needed to **copy for every word of memory**

$10^3$ Hosts

Runtime (s)

$10^3$

$10^2$

22.6   50.6   75.9   975.3

phantom   shadow   ns-3   grail

gRaIL (on ns-3) is **13x** slower than ns-3 alone

Can we design a tool with the
**performance benefits** of a uni-process
plugin-based architecture
AND
the improved **modularity** and **isolation**
of a mutli-process architecture?
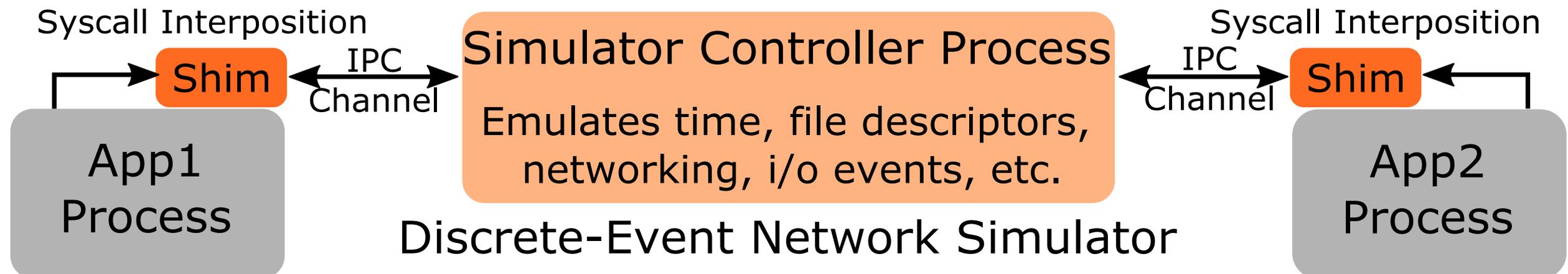
# Outline

motivation

**design**

evaluation

# Design Overview

1. parallel workers
2. direct execution
3. syscall interposition
4. syscall emulation
5. IPC
6. process control
7. CPU affinity

- Discrete-event packet-level network simulator
- Directly executes apps as standard Linux processes
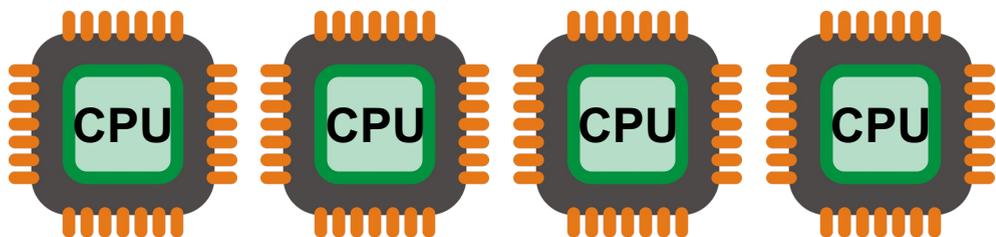- Intercepts all system calls made by apps and emulates them

- Simulates system call behavior and networking
  - File descriptors (files, sockets, pipes)
  - Event notification (poll, epoll, select)
  - Networking (buffers, protocols, ifaces)
  - DNS and routing (latency, bandwidth)

Syscall Interposition | Syscall Interposition

**Shim** ⟷ IPC Channel ⟷ **Simulator Controller Process** Emulates time, file descriptors, networking, i/o events, etc. ⟷ IPC Channel ⟷ **Shim**

App1 Process | | App2 Process

**Discrete-Event Network Simulator**

1. **parallel workers**
2. direct execution
3. syscall interposition
4. syscall emulation
5. IPC
6. process control
7. CPU affinity

# Parallel Worker Threads

## Goal: efficiently parallelize simulation workload

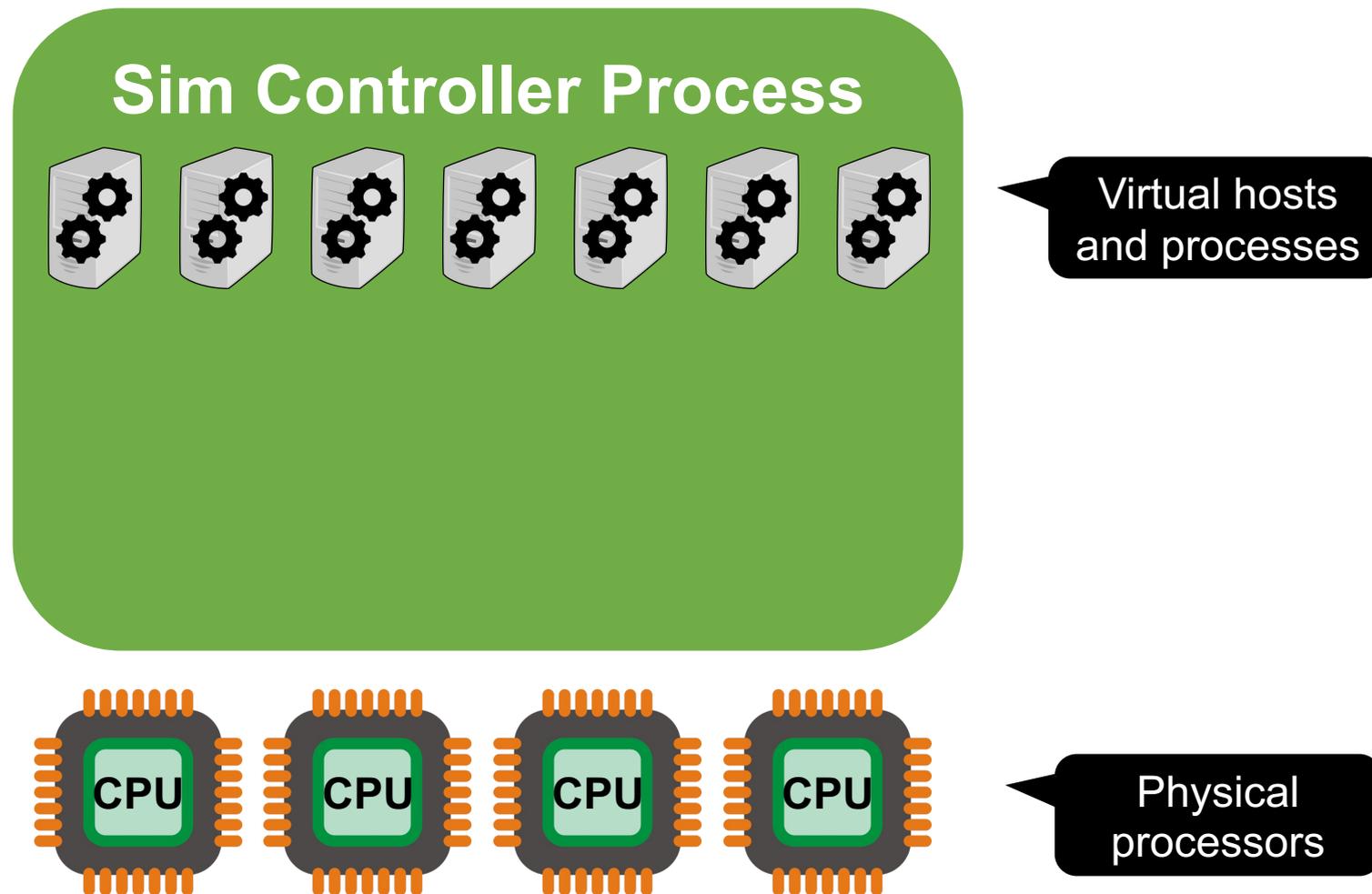

Sim Controller Process

CPU   CPU   CPU   CPU

Physical processors

1. **parallel workers**
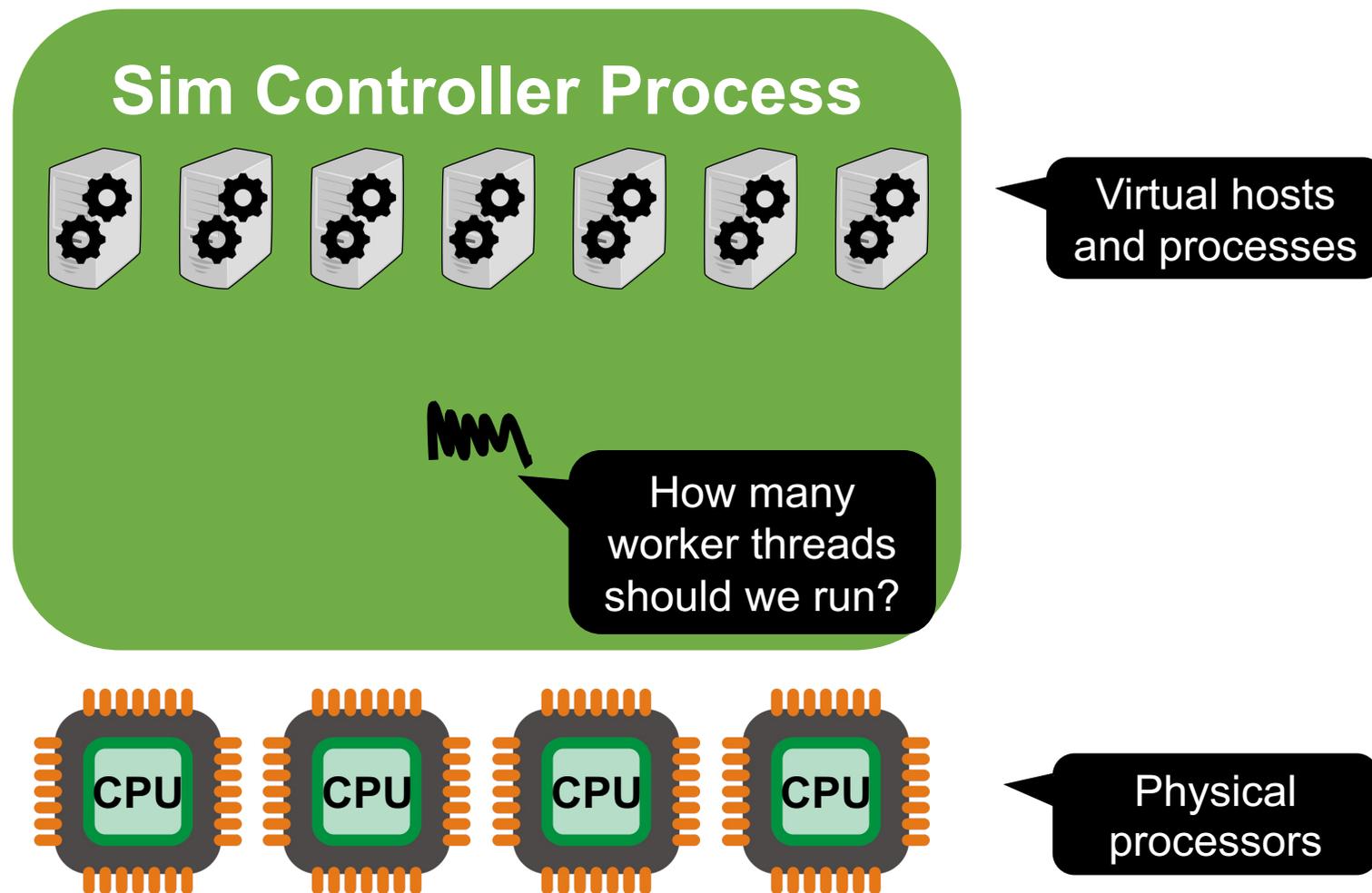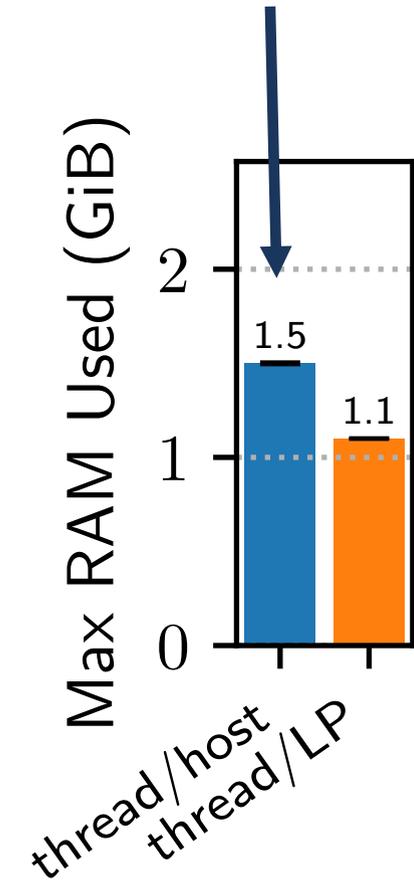2. direct execution
3. syscall interposition
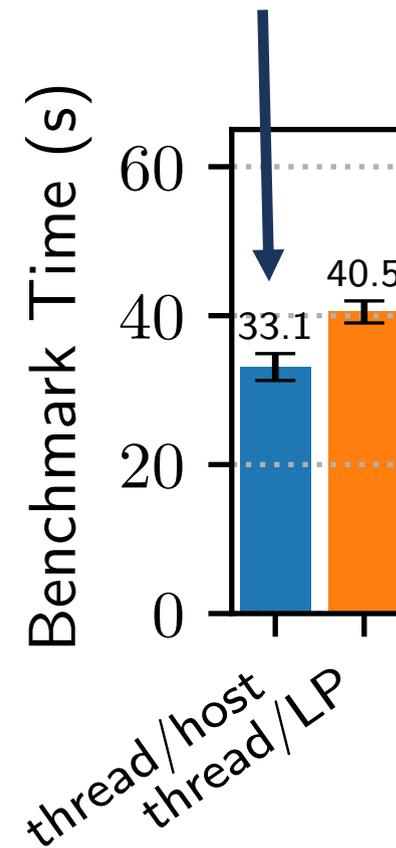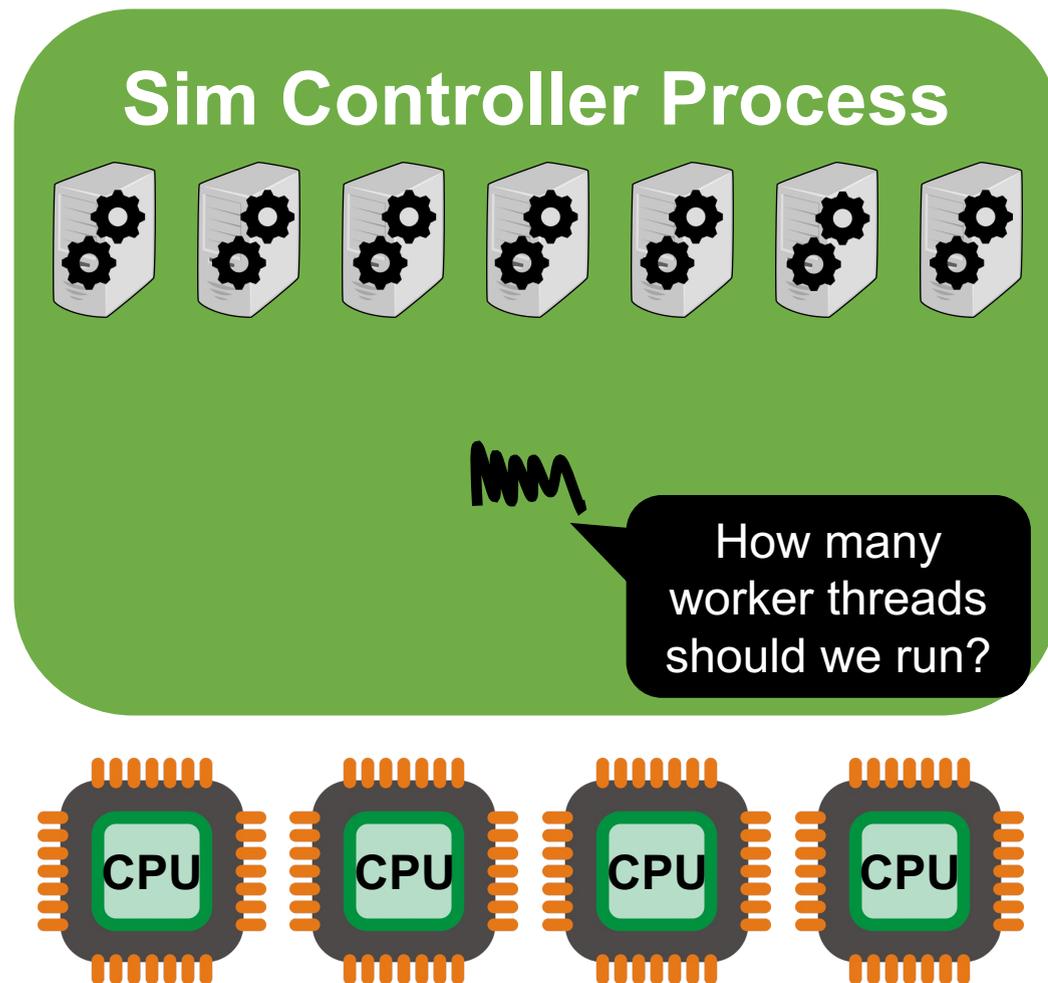4. syscall emulation
5. IPC
6. process control
7. CPU affinity

# Parallel Worker Threads

## Goal: efficiently parallelize simulation workload



Sim Controller Process

Virtual hosts and processes

Physical processors

1. **parallel workers**
2. direct execution
3. syscall interposition
4. syscall emulation
5. IPC
6. process control
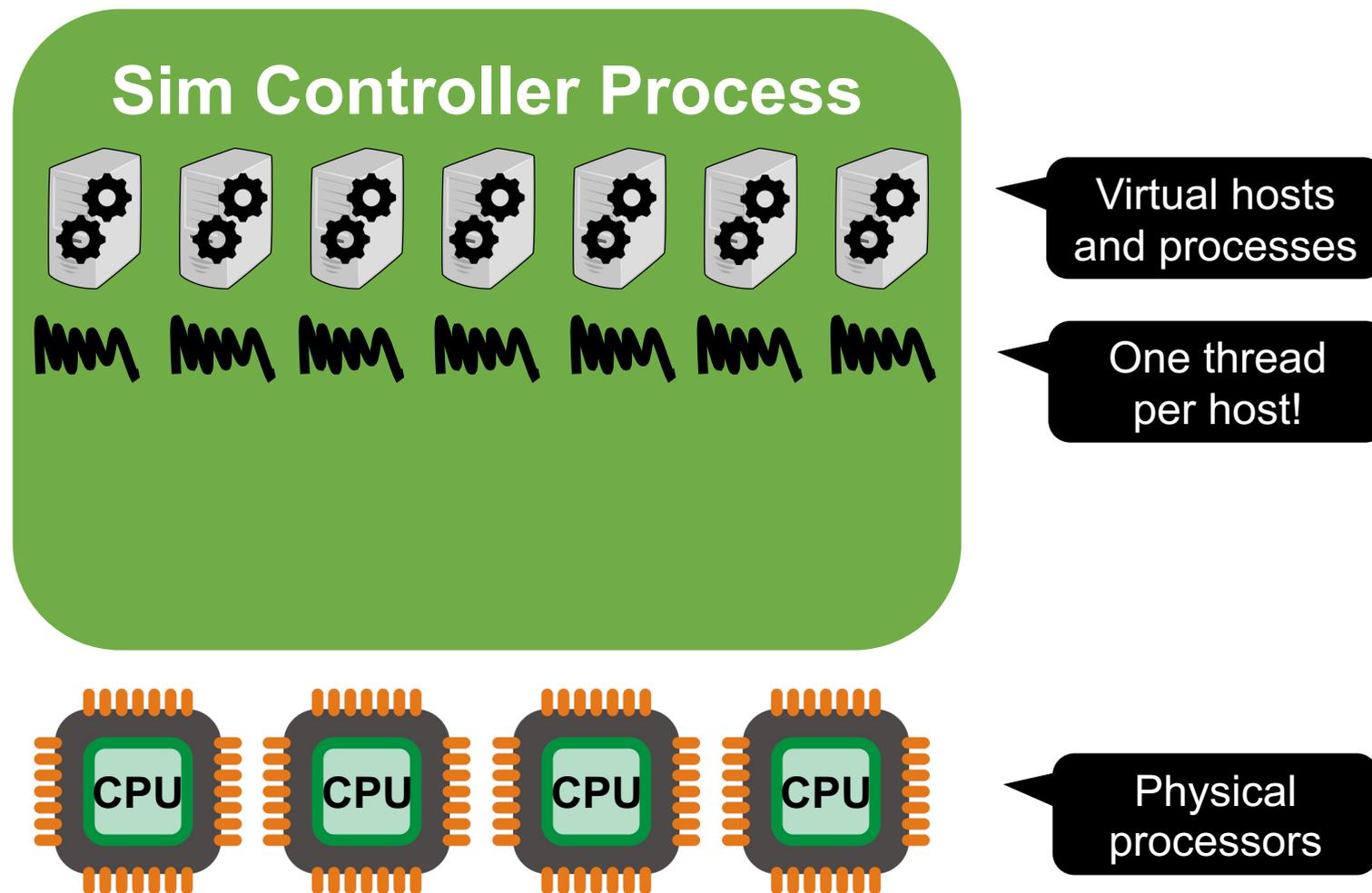7. CPU affinity

# Parallel Worker Threads

## Goal: efficiently parallelize simulation workload

1. **parallel workers**
2. direct execution
3. syscall interposition
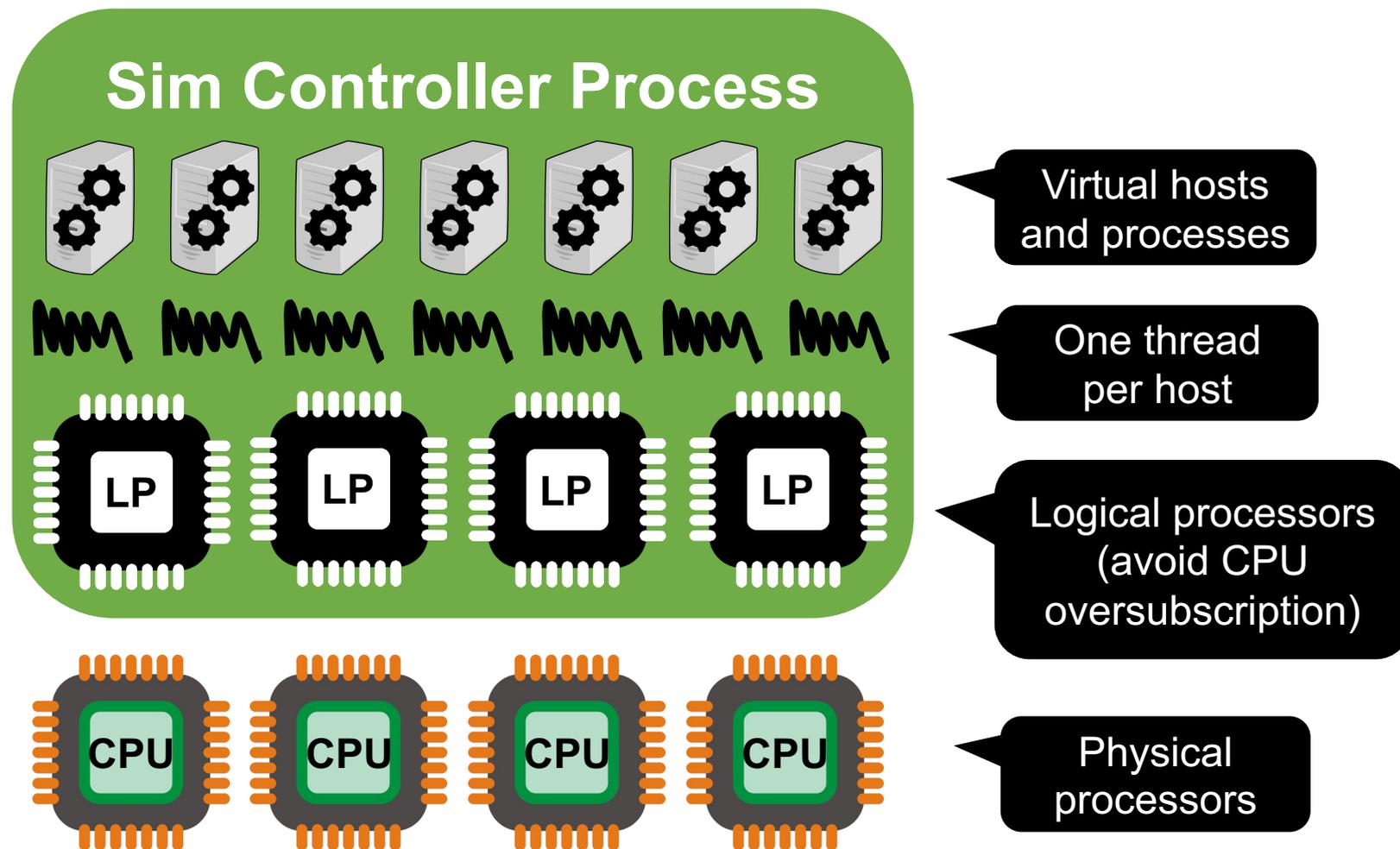4. syscall emulation
5. IPC
6. process control
7. CPU affinity

# Parallel Worker Threads

## Goal: efficiently parallelize simulation workload

1. **parallel workers**
2. direct execution
3. syscall interposition
4. syscall emulation
5. IPC
6. process control
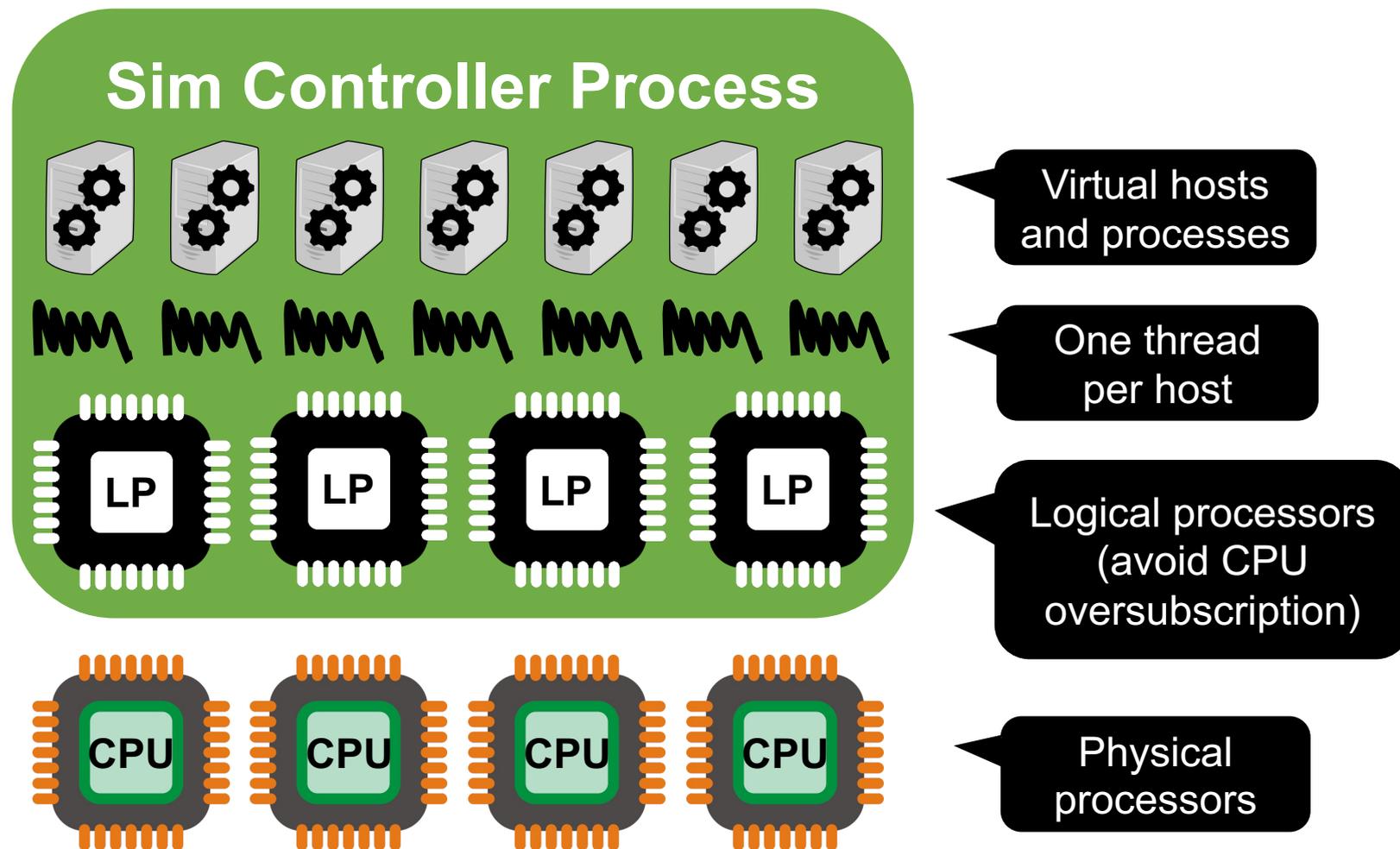7. CPU affinity

# Parallel Worker Threads

## Goal: efficiently parallelize simulation workload

1. **parallel workers**
2. direct execution
3. syscall interposition
4. syscall emulation
5. IPC
6. process control
7. CPU affinity

# Parallel Worker Threads

## Goal: efficiently parallelize simulation workload



Virtual hosts and processes

One thread per host

Logical processors (avoid CPU oversubscription)

Physical processors

## Thread scheduling:

- Work stealing

- Each LP starts a thread
  1. Runs all assigned events in current round (1 ms)
  2. Set thread to *waiting*
  3. Starts next *waiting* thread (if any)

- When all threads *waiting*
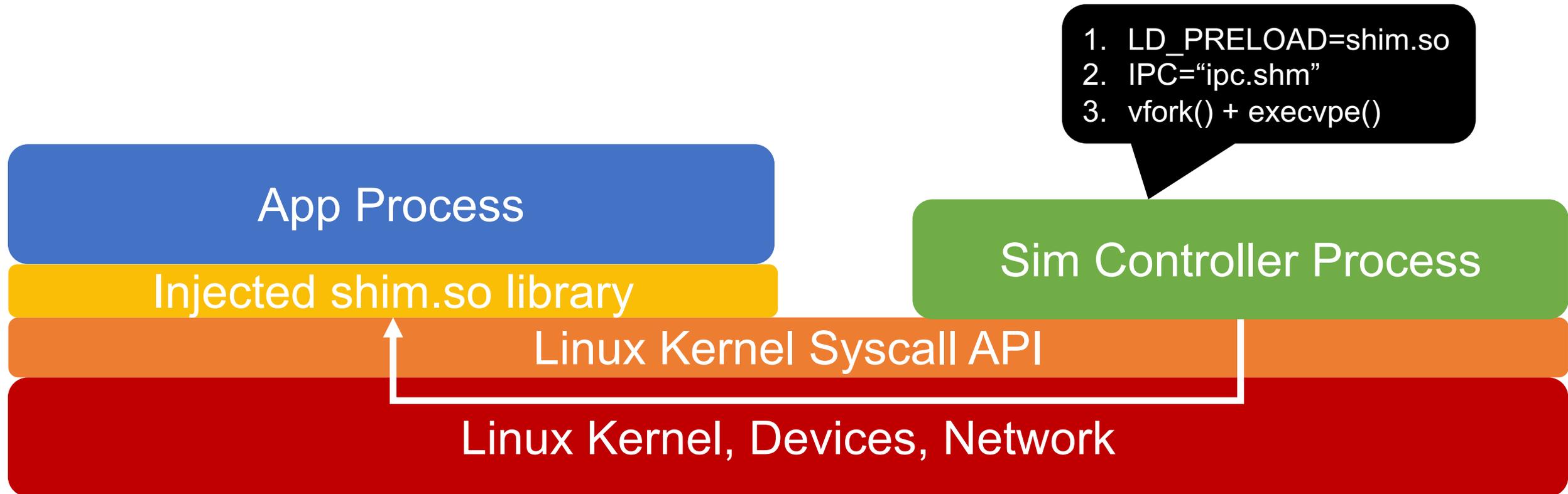  - Advance round clock
  - Repeat

1. parallel workers
2. **direct execution**
3. syscall interposition
4. syscall emulation
5. IPC
6. process control
7. CPU affinity

# Direct Execution
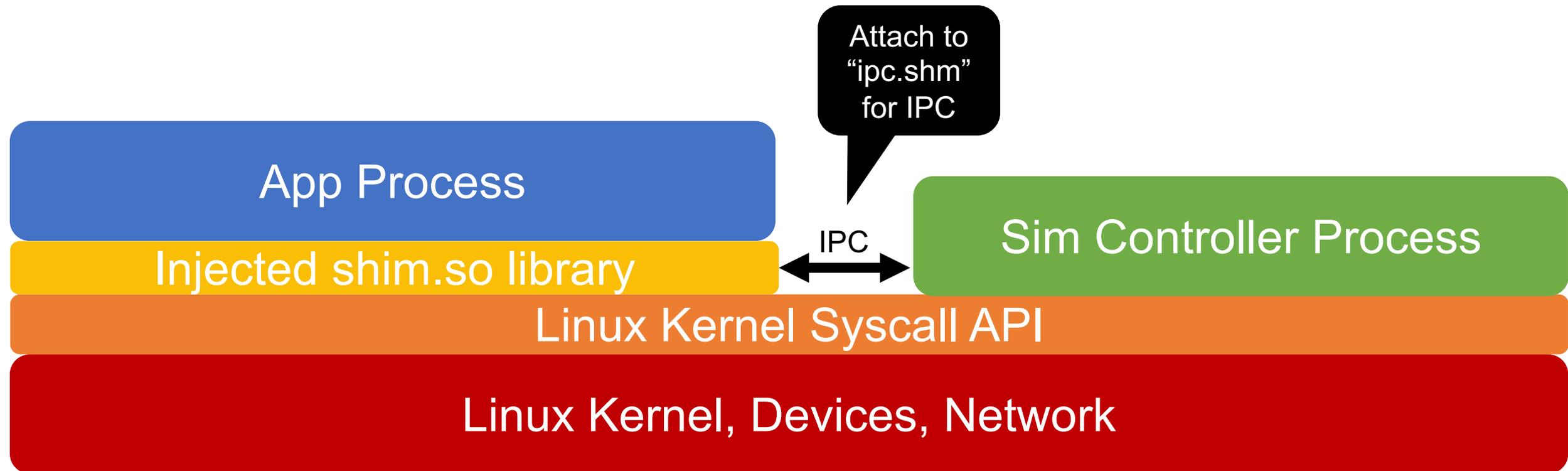
Sim Controller Process

Linux Kernel Syscall API

Linux Kernel, Devices, Network

1. parallel workers
2. direct execution
3. **syscall interposition**
4. syscall emulation
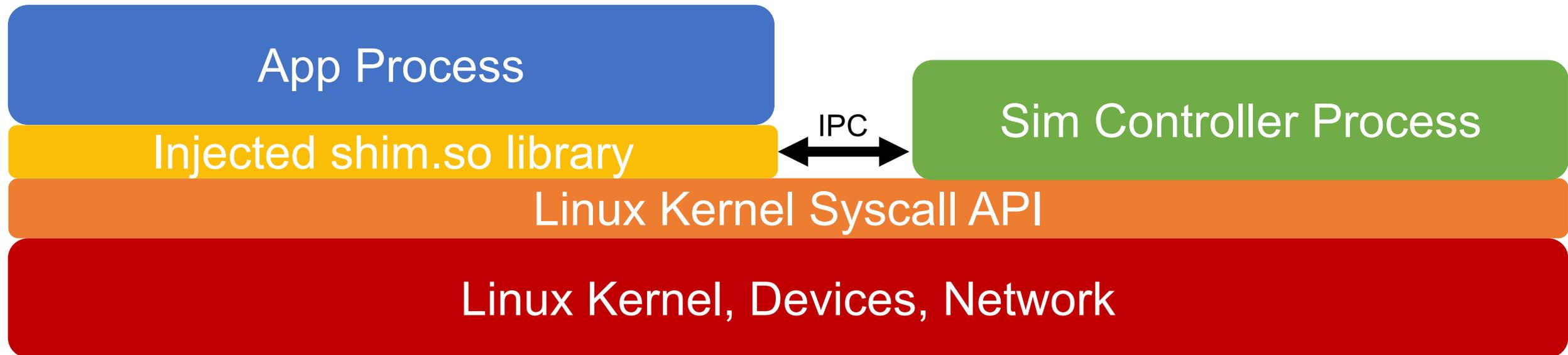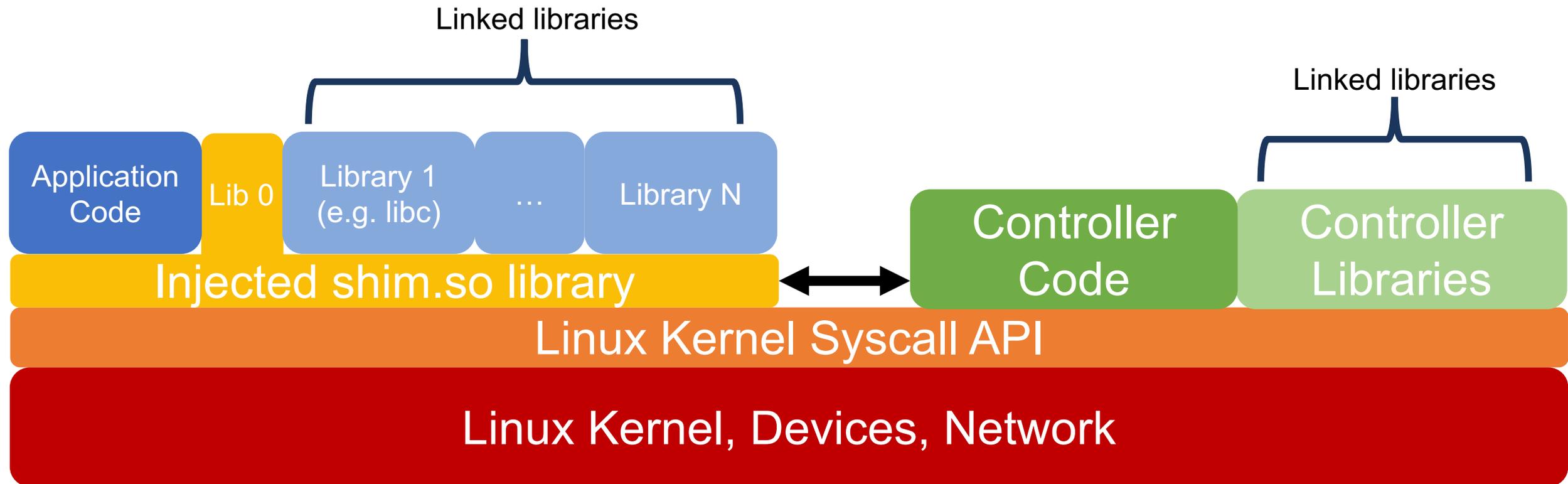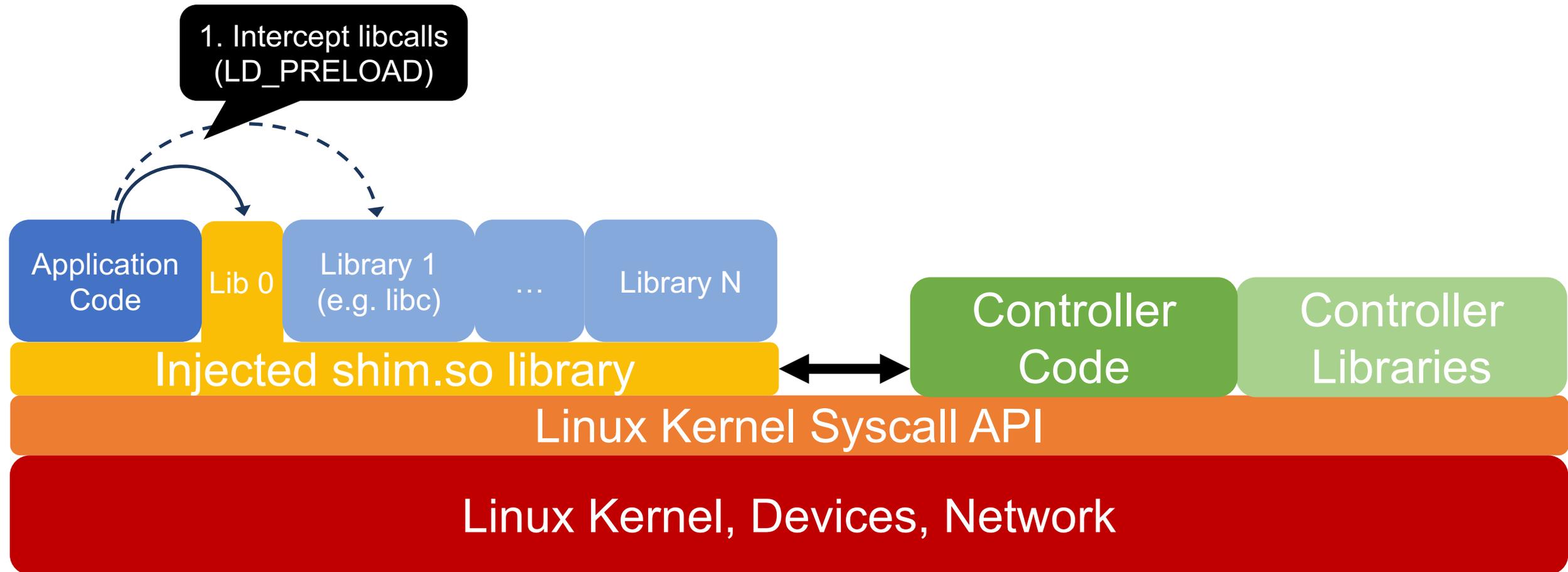
5. IPC
6. process control
7. CPU affinity

# Syscall Interposition

1. parallel workers
2. direct execution
3. **syscall interposition**
4. syscall emulation
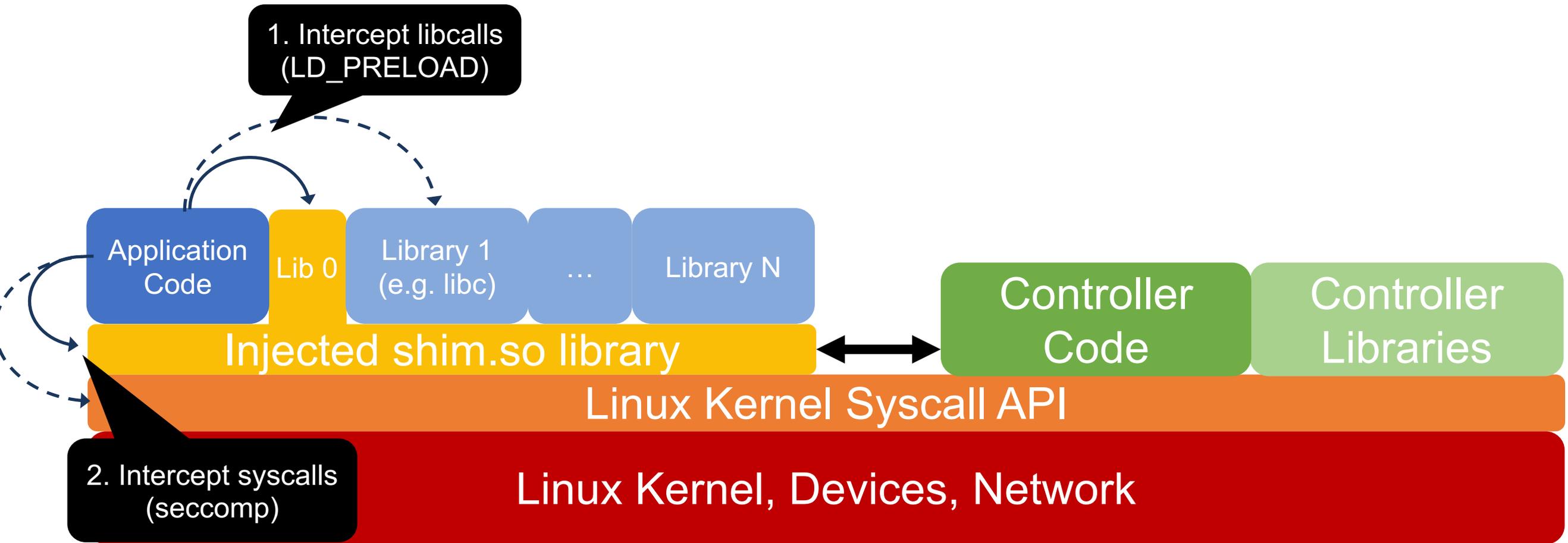5. IPC
6. process control
7. CPU affinity

# Syscall Interposition



1. Intercept libcalls (LD_PRELOAD)

2. Intercept syscalls (seccomp)

Application Code | Lib 0 | Library (e.g. lib...)

Injected shim

**blocking nanosleep**

- only seccomp: 15.34
- preload+seccomp: 13.37
- uni-process: 9.51

**nonblocking nanosleep**

- only seccomp: 9.64
- preload+seccomp: 6.8
- uni-process: 0

**1k write+read**

- only seccomp: 15.26
- preload+seccomp: 11.13
- uni-process: 7.78

Benchmark Time ($\mu s$)

1. parallel workers    5. IPC
2. direct execution    6. process control
3. syscall interposition    7. CPU affinity
4. **syscall emulation**

# Syscall Emulation



- File descriptors (files, sockets, pipes)
- Event notification (poll, epoll, select)
- Networking (buffers, protocols, ifaces)
- DNS and routing (latency, bandwidth)

Application Code
Lib 0
Library 1 (e.g. libc)
…
Library N
Injected shim.so library

IPC

Controller Code
Controller Libraries
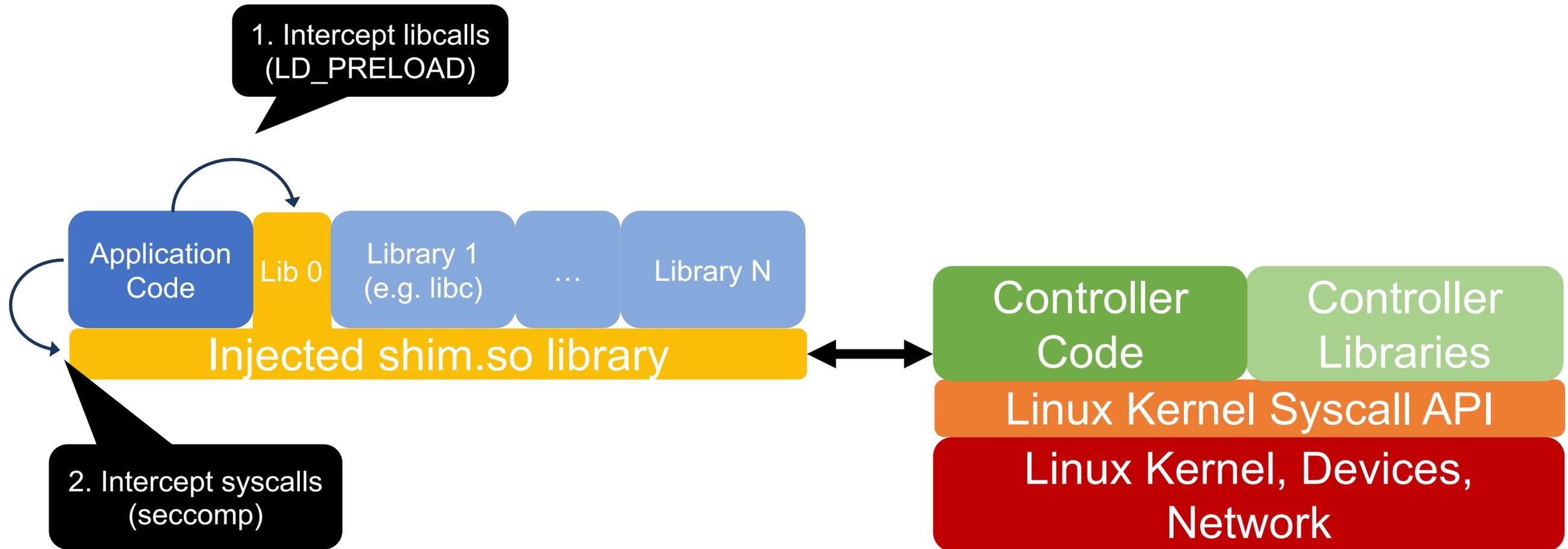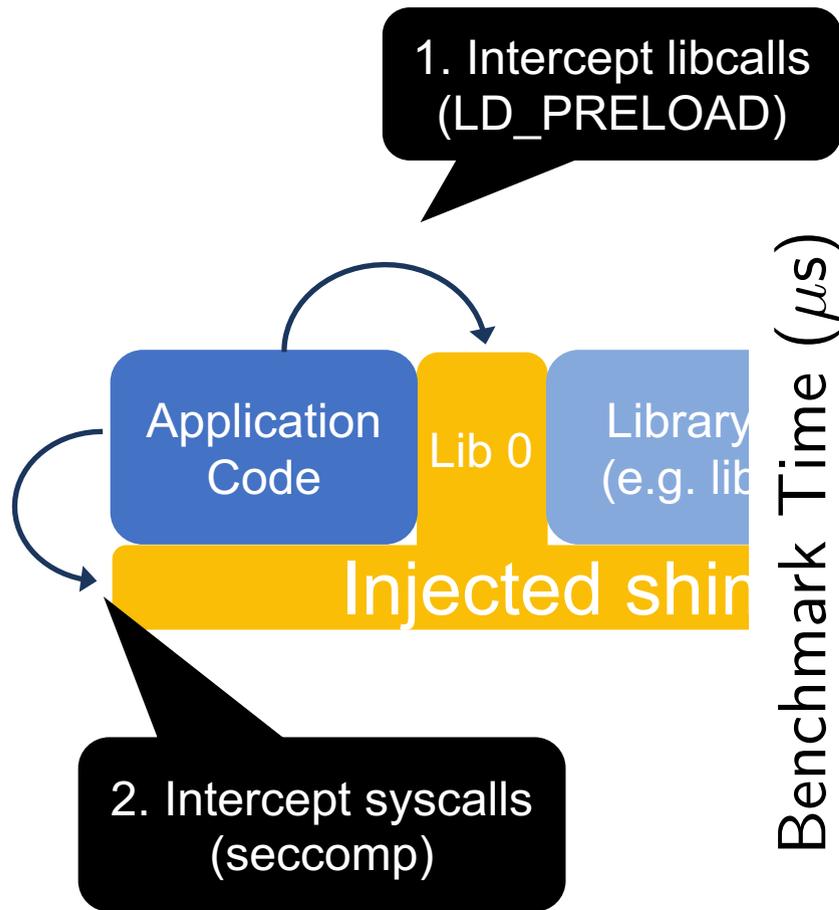Linux Kernel Syscall API
Linux Kernel, Devices, Network

# Syscall Emulation

1. parallel workers
2. direct execution
3. syscall interposition
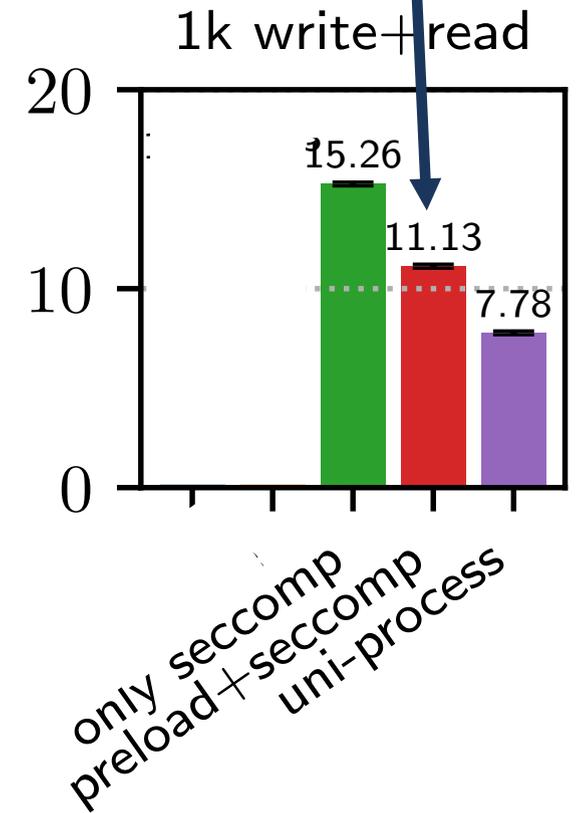4. **syscall emulation**
5. IPC
6. process control
7. CPU affinity

**Emulate in Shim**
- faster
- for hotpath syscalls (time)

Application Code

Lib 0

Library 1 (e.g. libc)

...

Library N

Injected shim.so library

IPC

Controller Code

Controller Libraries

Linux Kernel Syscall API

Linux Kernel, Devices, Network

- File descriptors (files, sockets, pipes)
- Event notification (poll, epoll, select)
- Networking (buffers, protocols, ifaces)
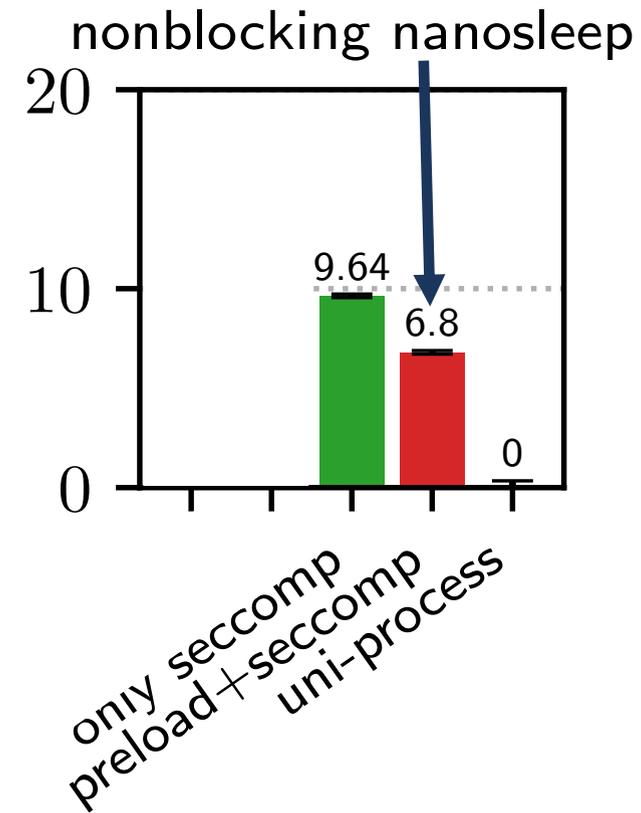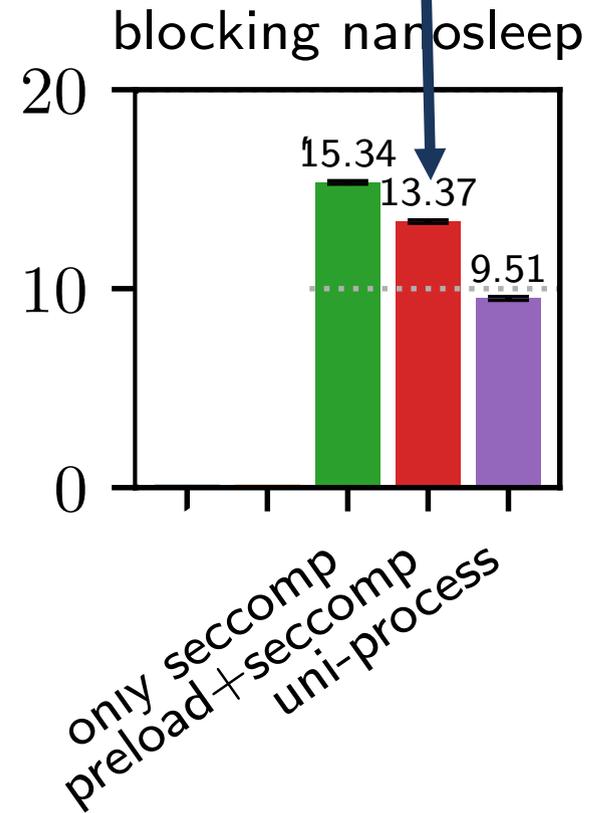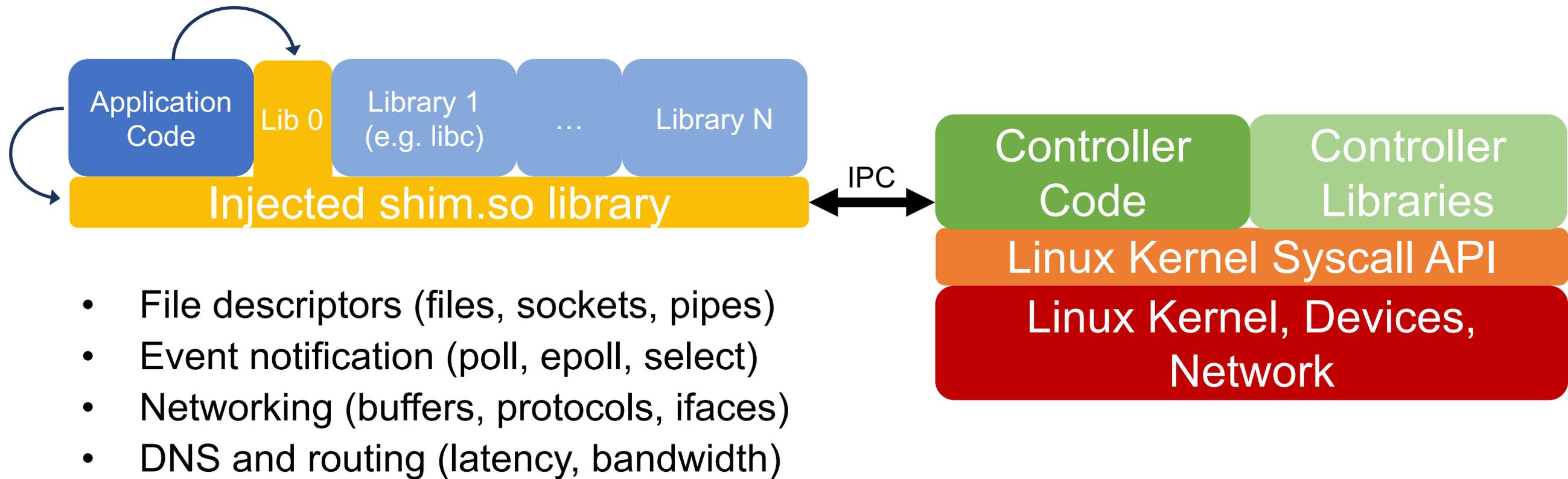- DNS and routing (latency, bandwidth)

# Syscall Emulation
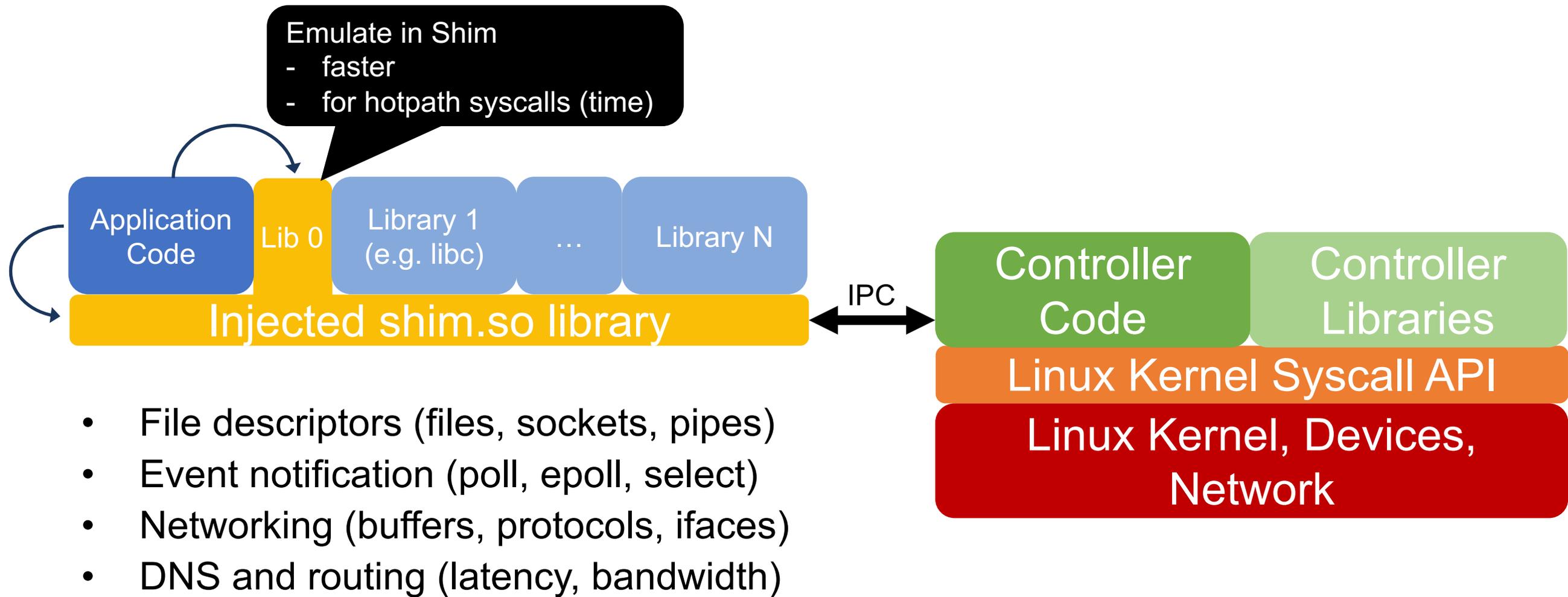
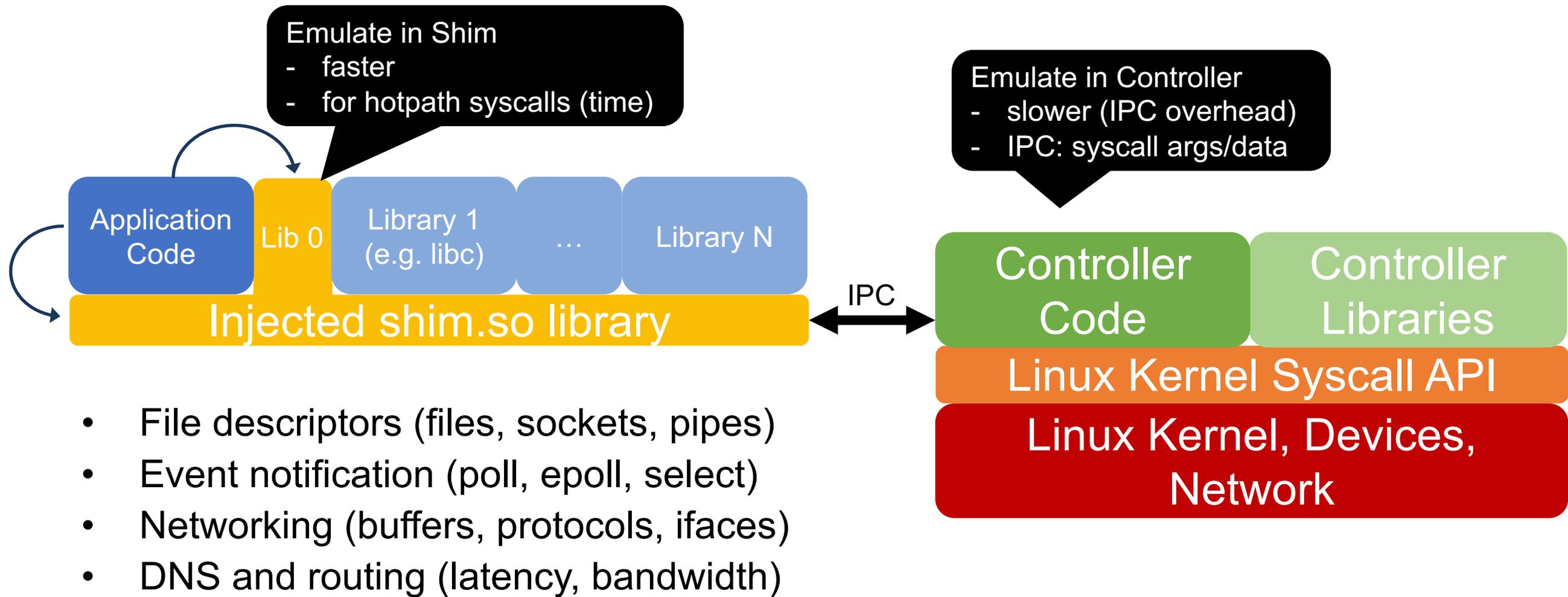1. parallel workers
2. direct execution
3. syscall interposition
4. **syscall emulation**
5. IPC
6. process control
7. CPU affinity

Emulate in Shim
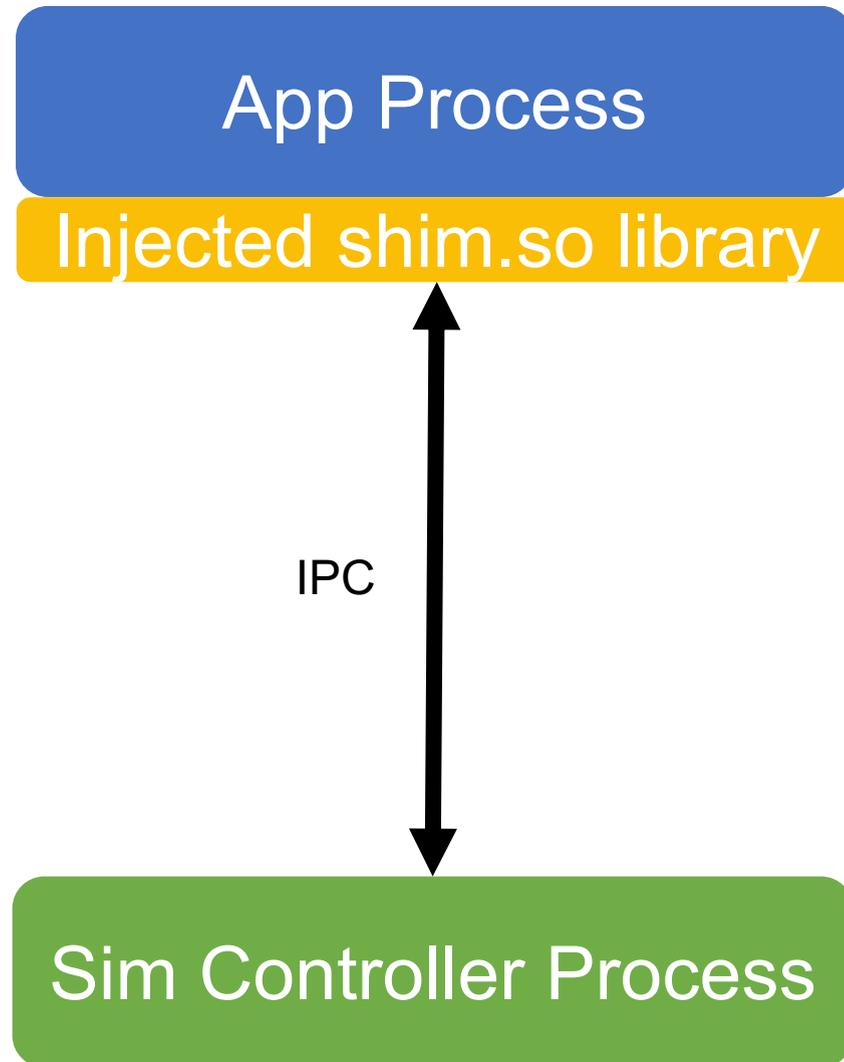- faster
- for hotpath syscalls (time)

Emulate in Controller
- slower (IPC overhead)
- IPC: syscall args/data

Application Code

Lib 0

Library 1 (e.g. libc)

…

Library N

Injected shim.so library

IPC

Controller Code

Controller Libraries

Linux Kernel Syscall API

Linux Kernel, Devices, Network

- File descriptors (files, sockets, pipes)
- Event notification (poll, epoll, select)
- Networking (buffers, protocols, ifaces)
- DNS and routing (latency, bandwidth)

1. parallel workers          **5. IPC**
2. direct execution          6. process control
3. syscall interposition     7. CPU affinity
4. syscall emulation

# Inter-Process Communication

App Process

Injected shim.so library

IPC

Sim Controller Process

1. parallel workers
2. direct execution
3. syscall interposition
4. syscall emulation

5. **IPC**
6. process control
7. CPU affinity

# Inter-Process Communication

App Process

Injected shim.so library

shared memory

Sim Controller Process

# Inter-Process Communication

1. parallel workers
2. direct execution
3. syscall interposition
4. syscall emulation

**5. IPC**
6. process control
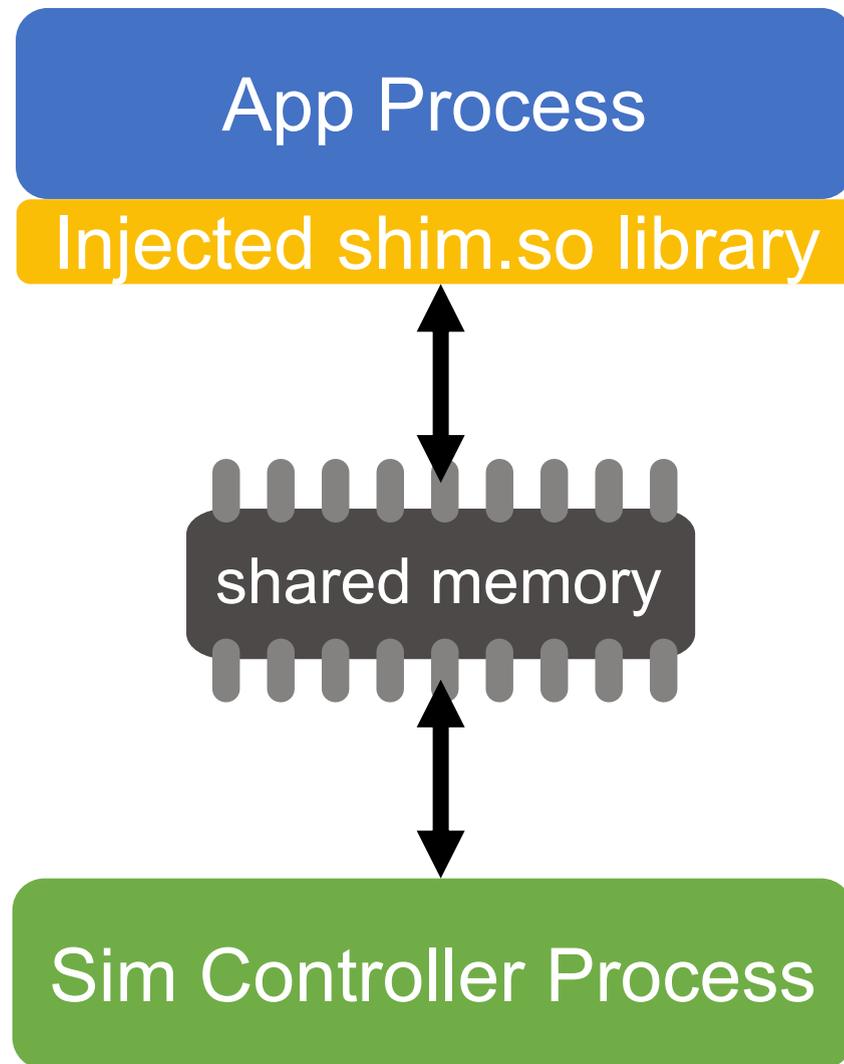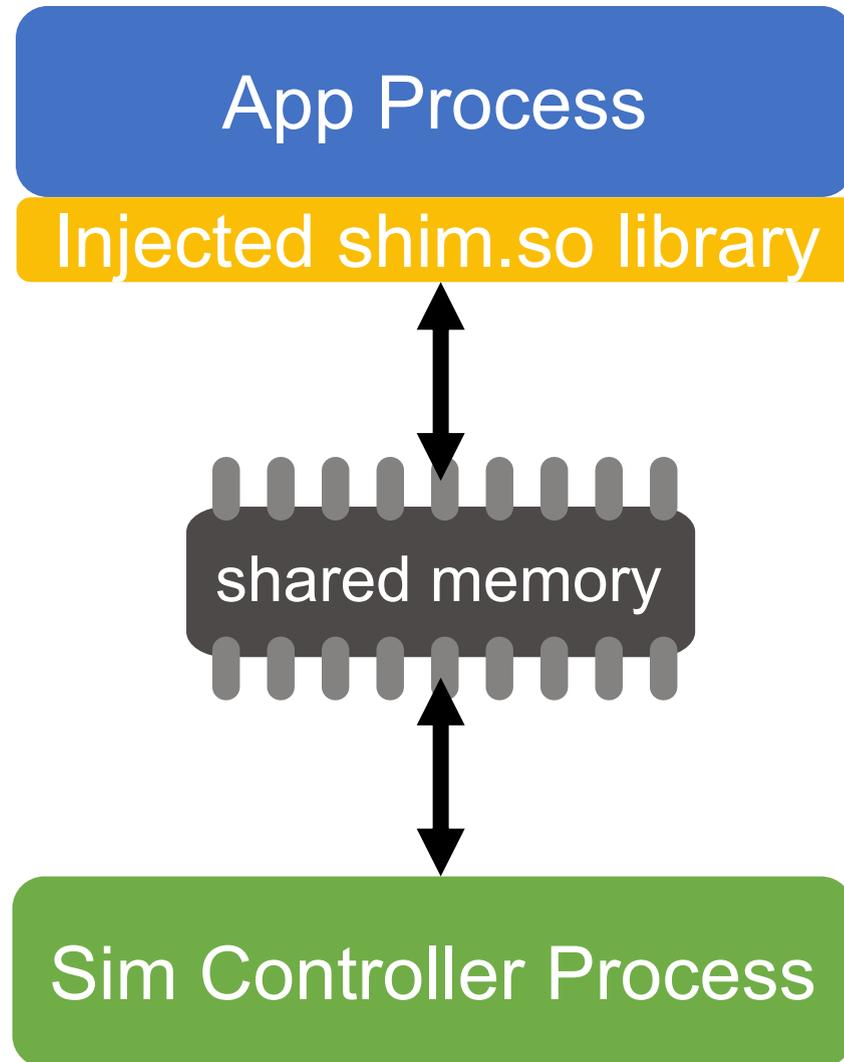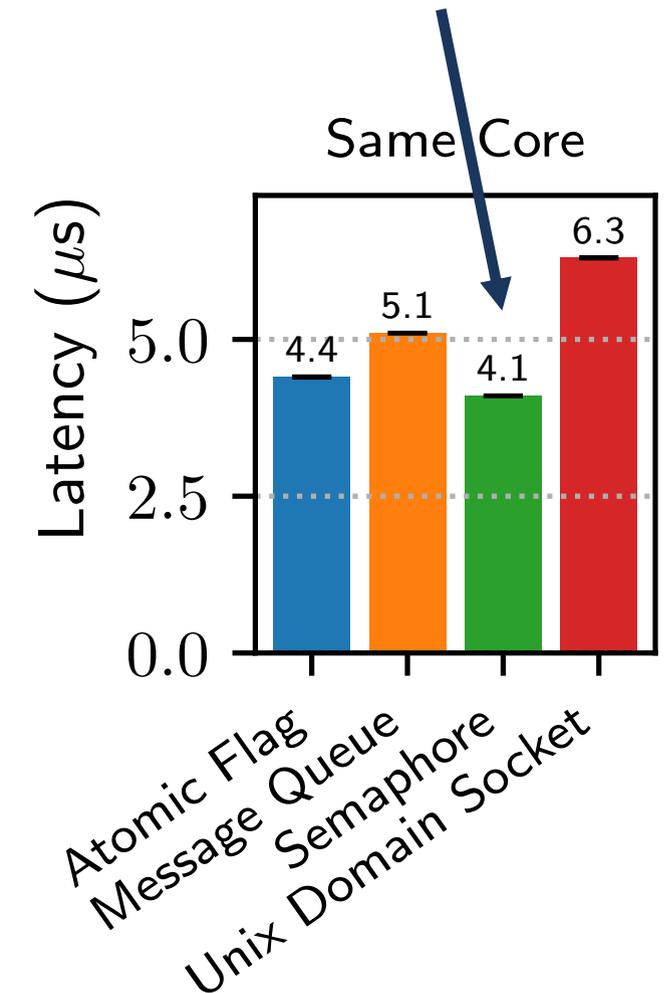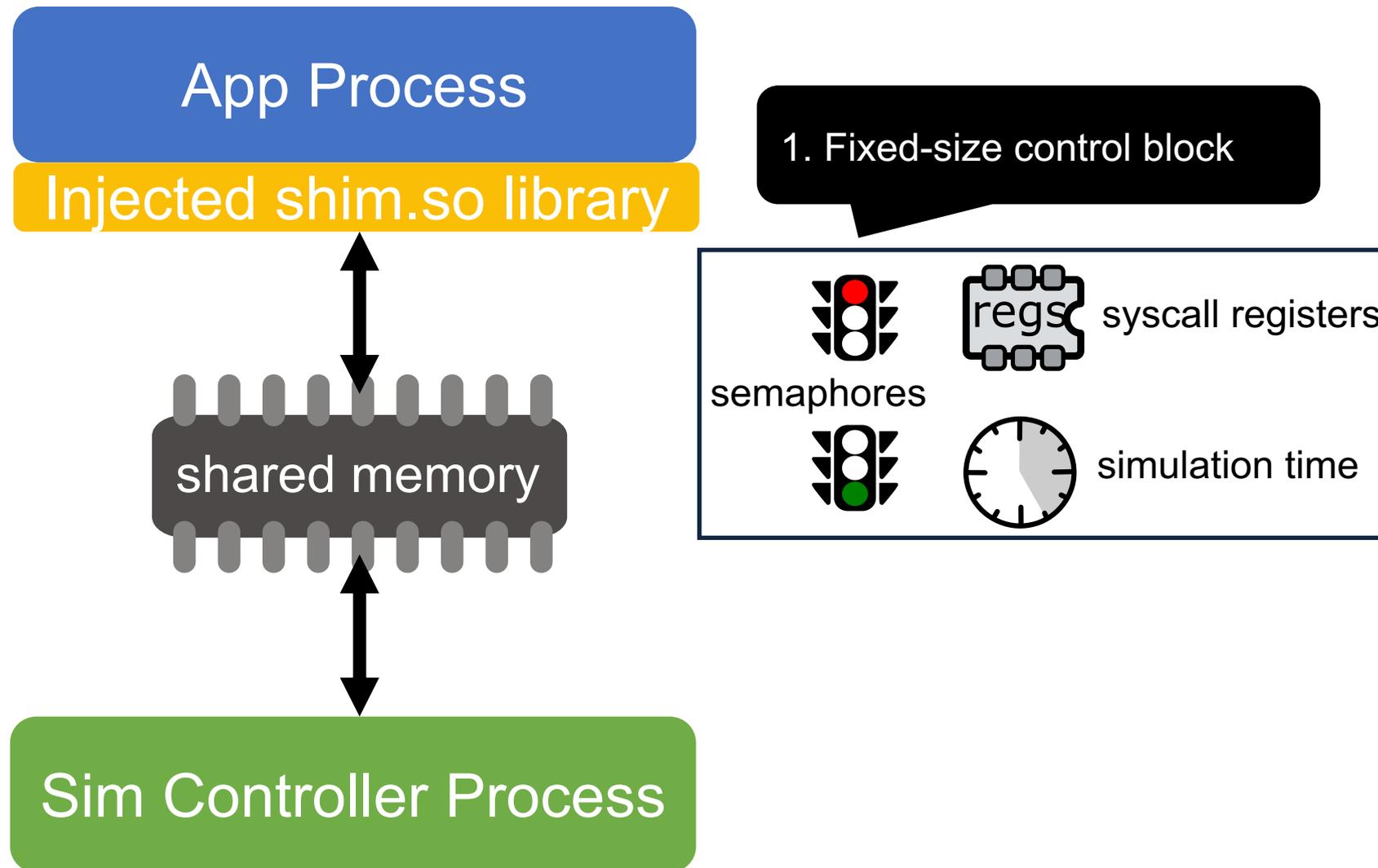7. CPU affinity

App Process

Injected shim.so library

shared memory

Sim Controller Process

Shared memory + semaphores is the fastest IPC method for two processes running on the same core

Same Core

6.3

5.1

4.4

4.1

Latency (μs)

5.0

2.5

0.0

Atomic Flag
Message Queue
Semaphore
Unix Domain Socket

1. parallel workers
2. direct execution
3. syscall interposition
4. syscall emulation
5. **IPC**
6. process control
7. CPU affinity

# Inter-Process Communication

**App Process**

**Injected shim.so library**

**shared memory**

**Sim Controller Process**

1. Fixed-size control block

semaphores · regs syscall registers

simulation time

1. parallel workers
2. direct execution
3. syscall interposition
4. syscall emulation
5. IPC
6. **process control**
7. CPU affinity

# Process Control

App Process

Injected shim.so library

shared memory

Sim Controller Process

semaphores

regs → syscall registers

simulation time

app stack

app heap

1. parallel workers    5. IPC
2. direct execution    **6. process control**
3. syscall interposition    7. CPU affinity
4. syscall emulation

# Process Control

App Process

Injected shim.so library

1. Write syscall registers
2. Signal controller semaphore
3. Wait on app semaphore

shared memory

Sim Controller Process

semaphores

syscall registers

simulation time

app stack    app heap

1. parallel workers
2. direct execution
3. syscall interposition
4. syscall emulation
5. IPC
6. **process control**
7. CPU affinity
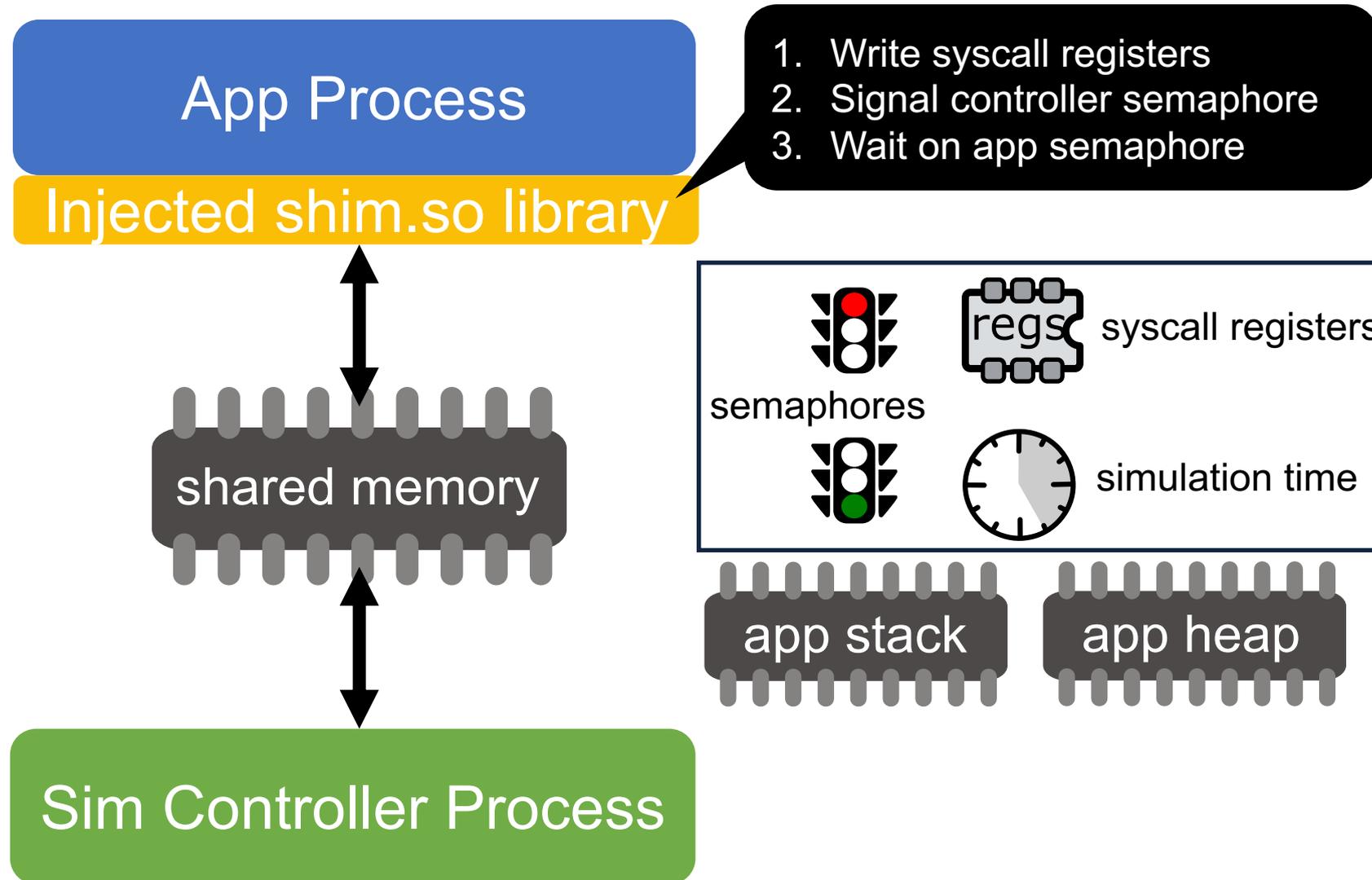
# Process Control

**U.S. NAVAL RESEARCH LABORATORY**



App Process

Injected shim.so library

1. Write syscall registers
2. Signal controller semaphore
3. Wait on app semaphore

shared memory

semaphores

regs  syscall registers

simulation time

app stack    app heap

Sim Controller Process

3. Read syscall registers
4. Return if nonblocking else leave app idle and return when ready
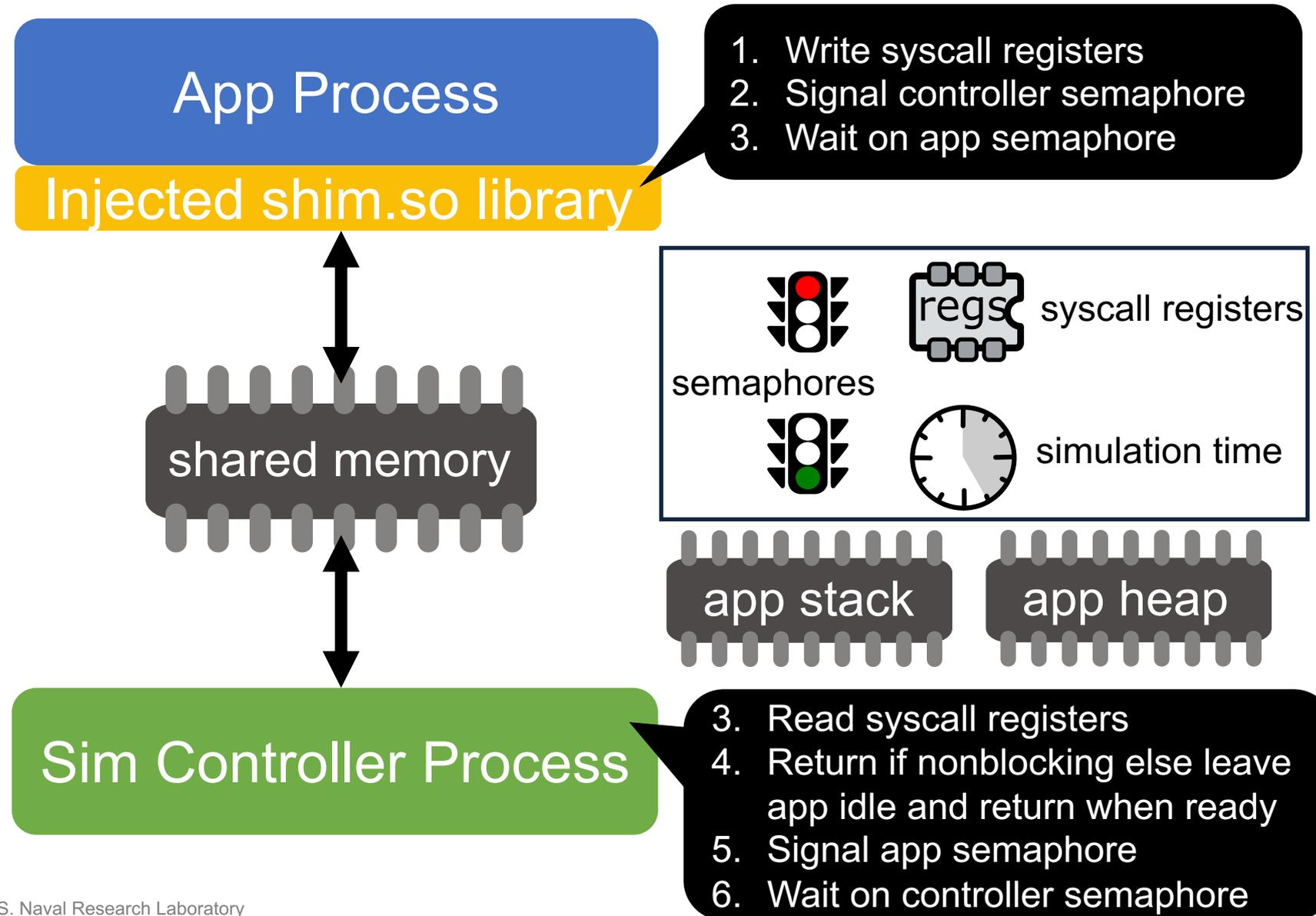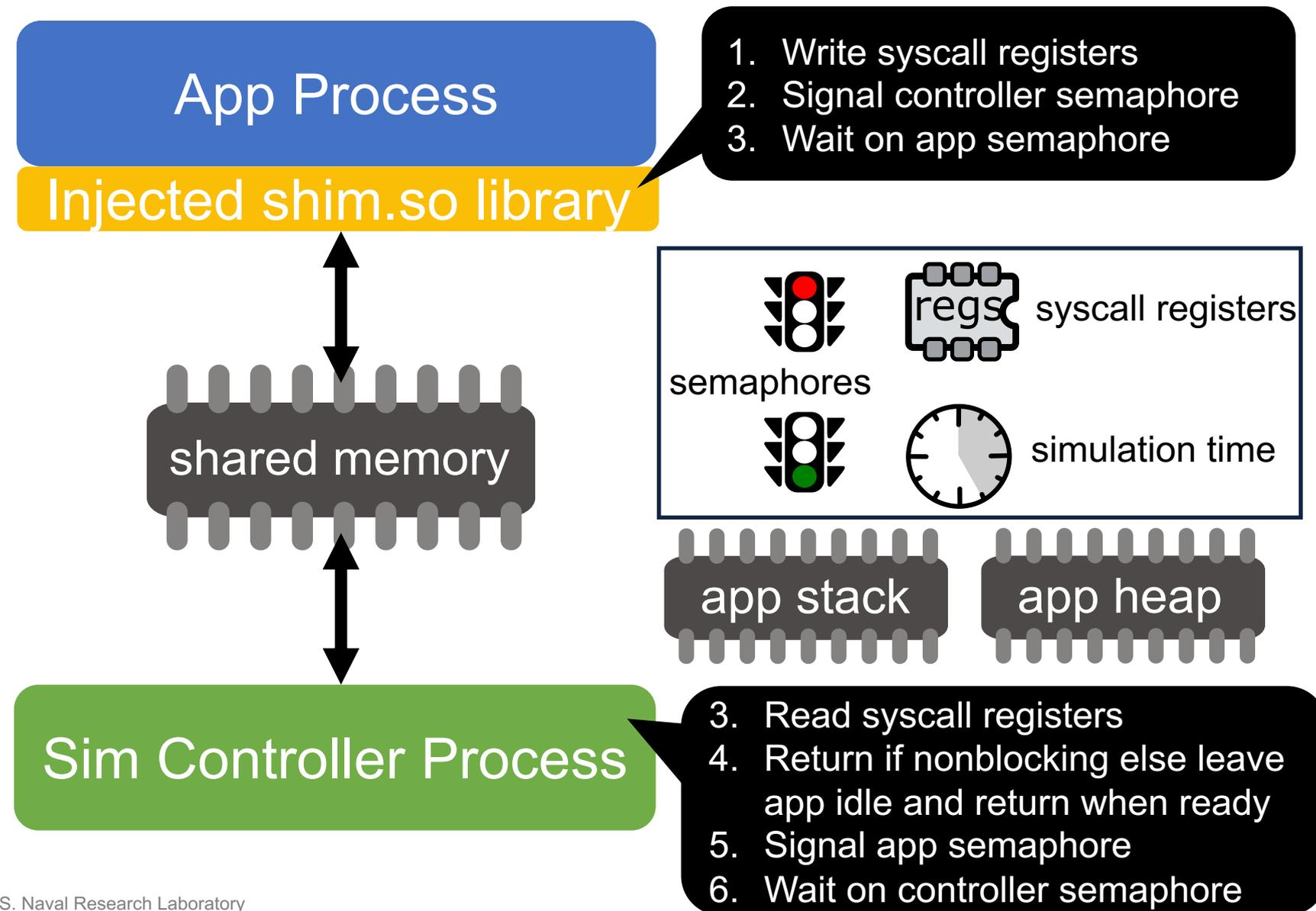5. Signal app semaphore
6. Wait on controller semaphore

1. parallel workers
2. direct execution
3. syscall interposition
4. syscall emulation
5. IPC
6. **process control**
7. CPU affinity

# Process Control

**App Process**

**Injected shim.so library**

1. Write syscall registers
2. Signal controller semaphore
3. Wait on app semaphore

shared memory

semaphores

regs — syscall registers

simulation time

app stack          app heap

**Sim Controller Process**

3. Read syscall registers
4. Return if nonblocking else leave app idle and return when ready
5. Signal app semaphore
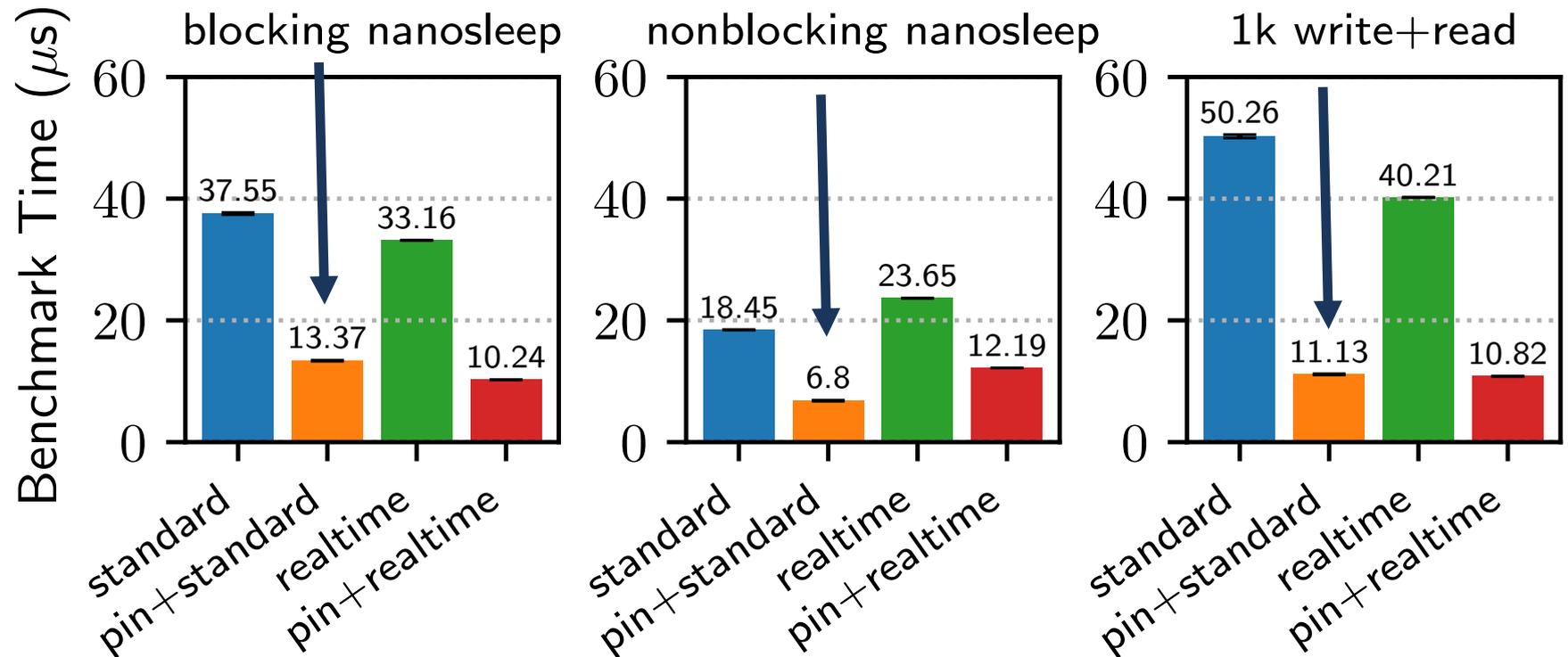6. Wait on controller semaphore

## Properties:

- Controller worker thread and its app process run synchronously

- Ensures nonconcurrent access to app stack and heap memory

1. parallel workers
2. direct execution
3. syscall interposition
4. syscall emulation
5. IPC
6. process control
**7. CPU affinity**

# Process Control

App Process

Injected shim.so library

shared memory

Sim Controller Process

Use CPU pinning to pin each worker thread and all of its managed processes to the same core



blocking nanosleep

nonblocking nanosleep

1k write+read

Benchmark Time ($\mu s$)

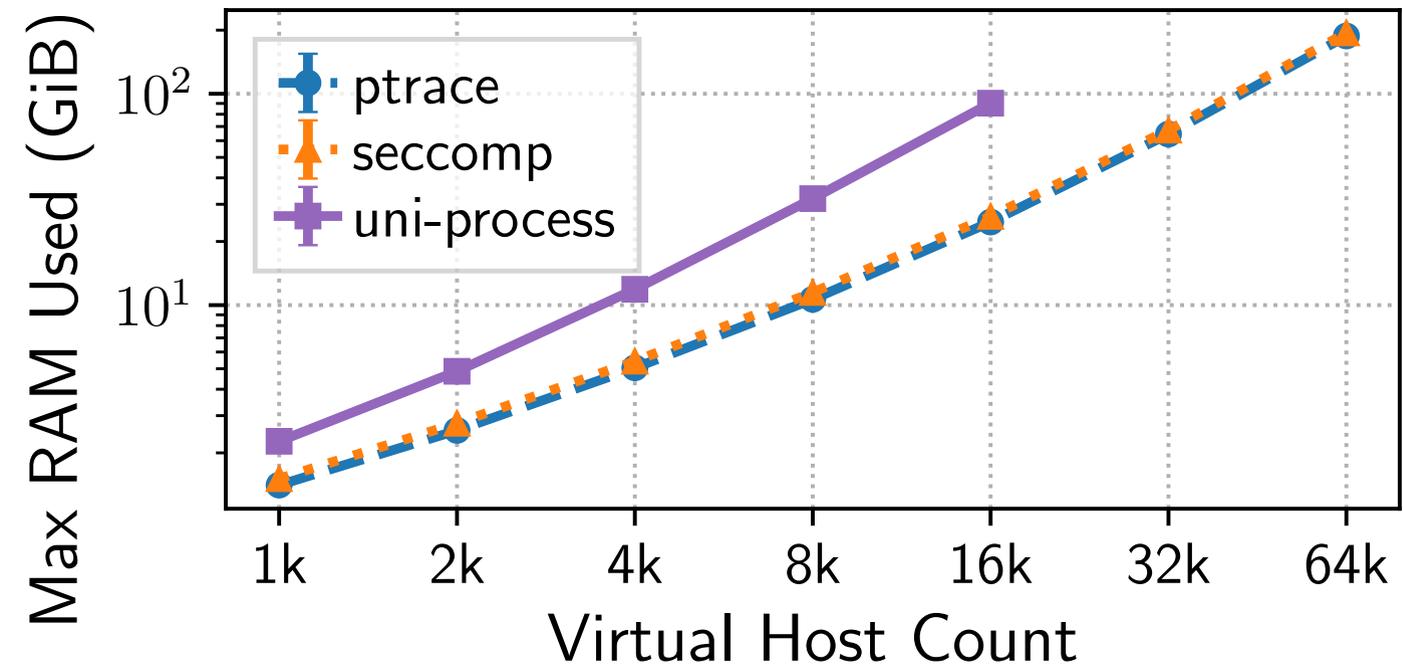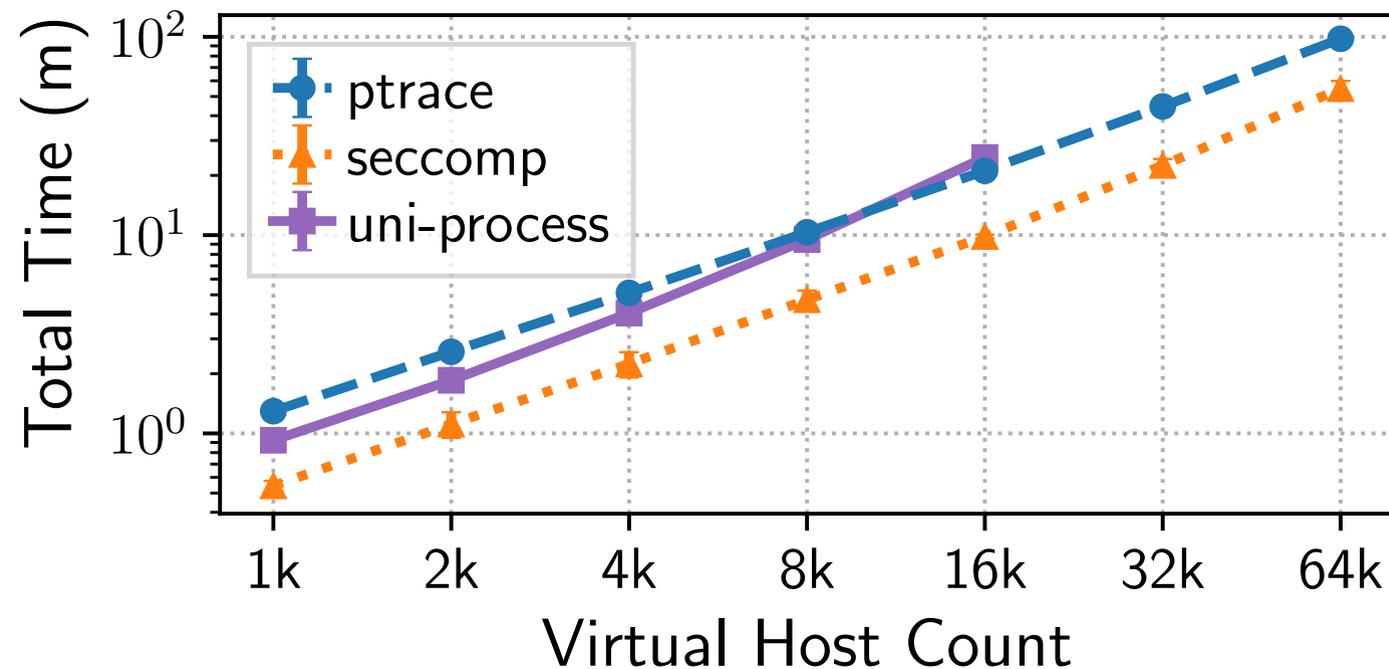| | blocking nanosleep | nonblocking nanosleep | 1k write+read |
|---|---|---|---|
| standard | 37.55 | 18.45 | 50.26 |
| pin+standard | 13.37 | 6.8 | 11.13 |
| realtime | 33.16 | 23.65 | 40.21 |
| pin+realtime | 10.24 | 12.19 | 10.82 |

# Outline

motivation
design
**evaluation**

# Evaluation: Large P2P Benchmarks

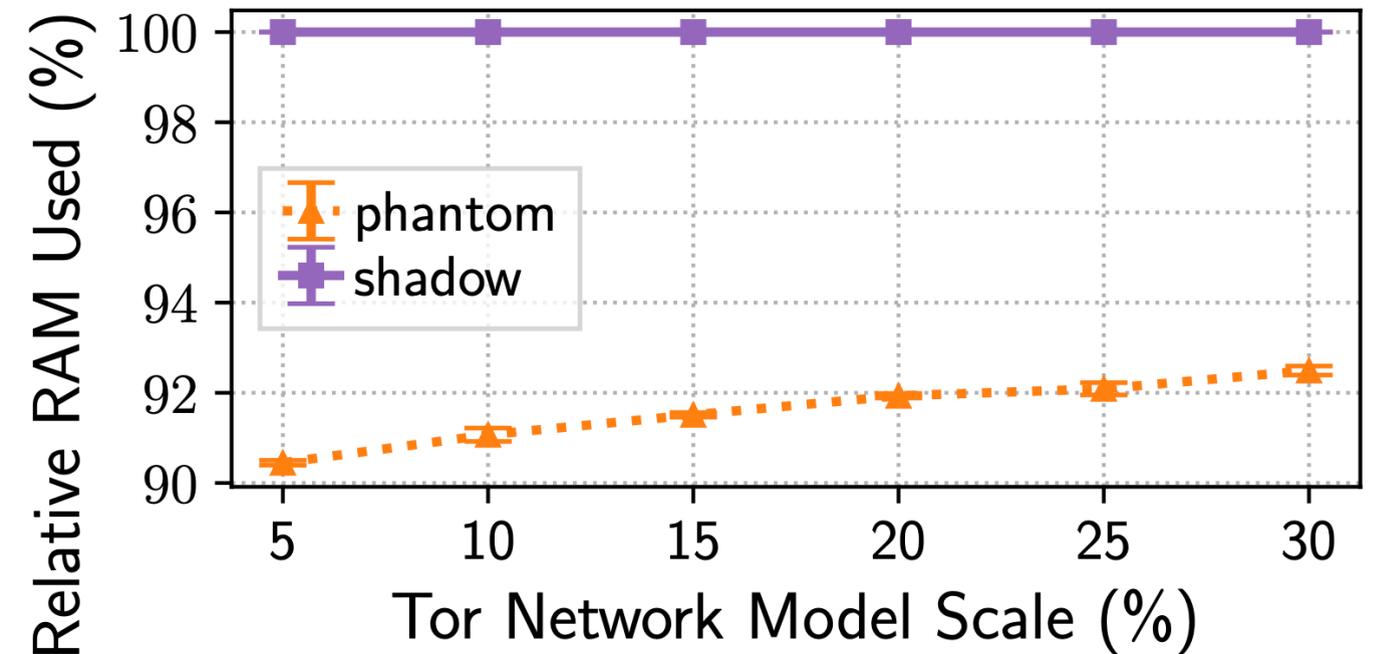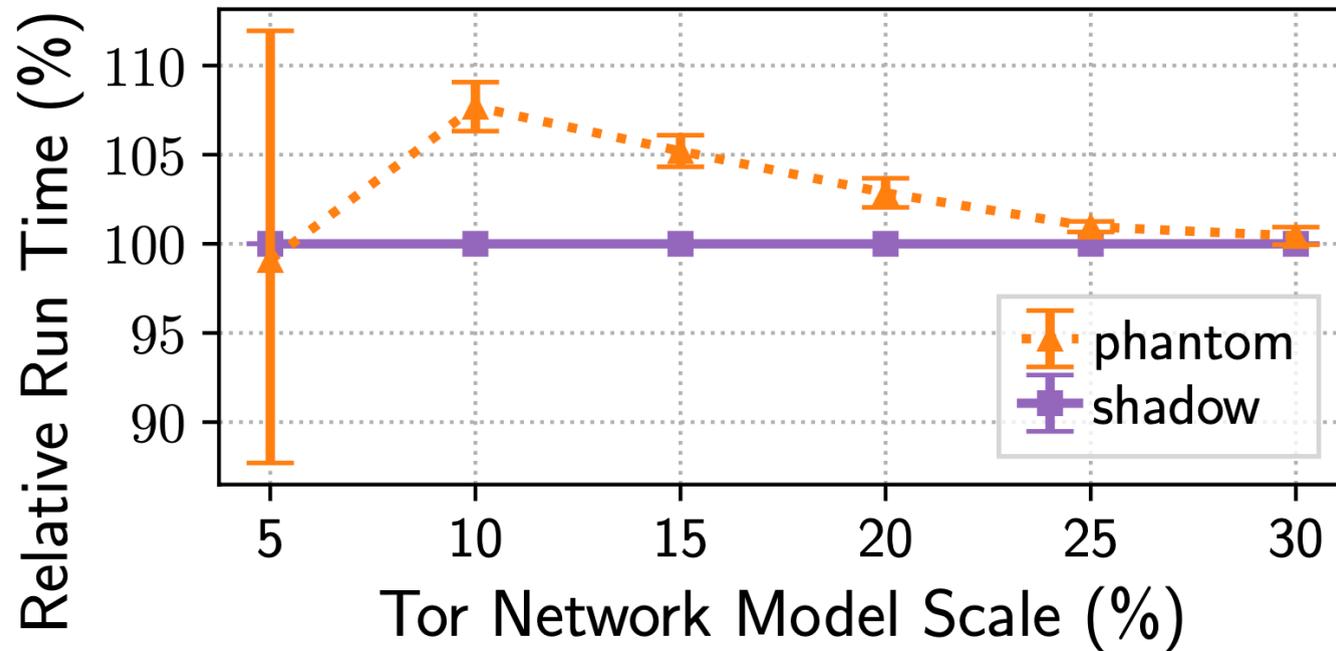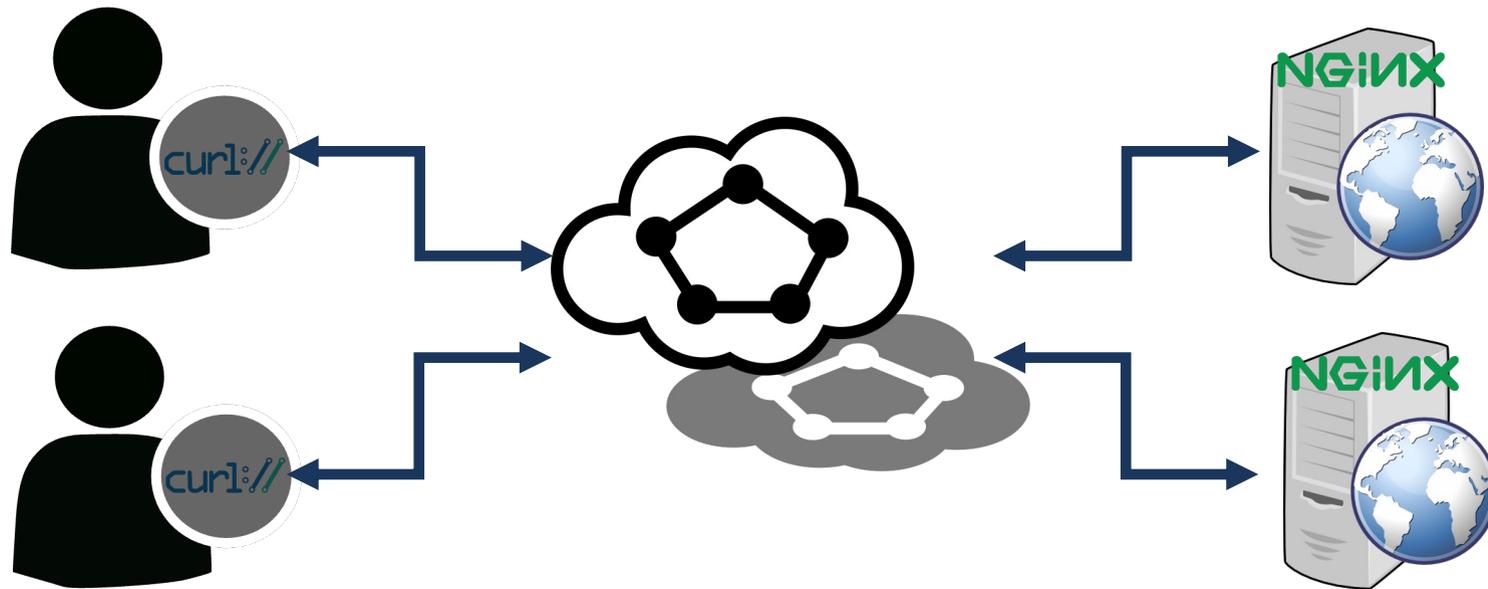Faster and more scalable than the uni-process, plugin architecture!

Uses significantly less memory than the uni-process, plugin architecture!

# Designed a new, hybrid network simulator/emulator

- co-opts Linux processes into a discrete-event network simulation that emulates kernel functionality
- enables large-scale, distributed system experiments

- 2.3x faster than Shadow v1
- 3.4x faster than NS-3
- 43x faster than gRaIL [ToN'19]



$10^3$ Hosts

- Merged into the open-source Shadow project and synonymous with Shadow v2
- Artifacts: https://netsim-atc2022.github.io