

Enabling Hybrid Parallel Runtimes Through Kernel and Virtualization Support

Kyle C. Hale and Peter Dinda



**Northwestern
University**

HOBBS
xstack.sandia.gov/hobbes

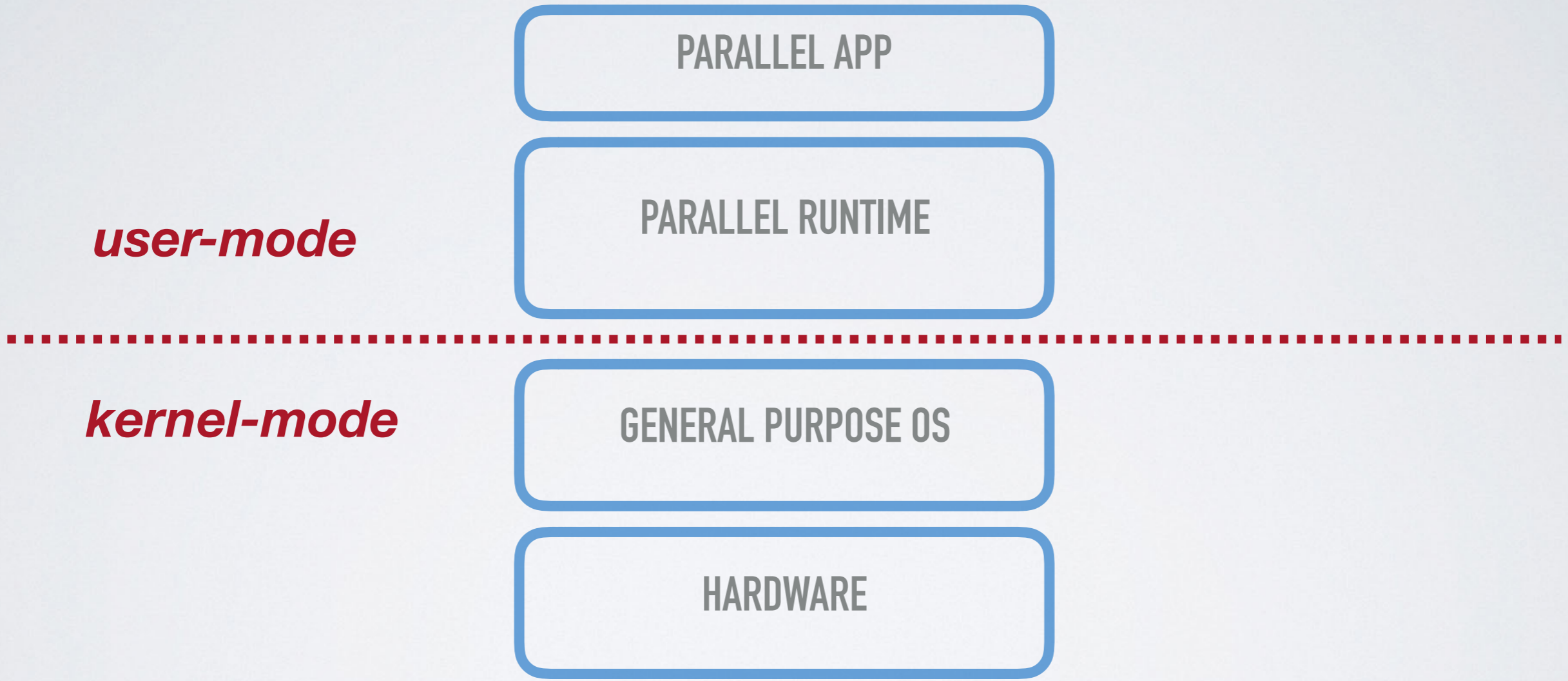
Hybrid Runtimes

the runtime IS the kernel

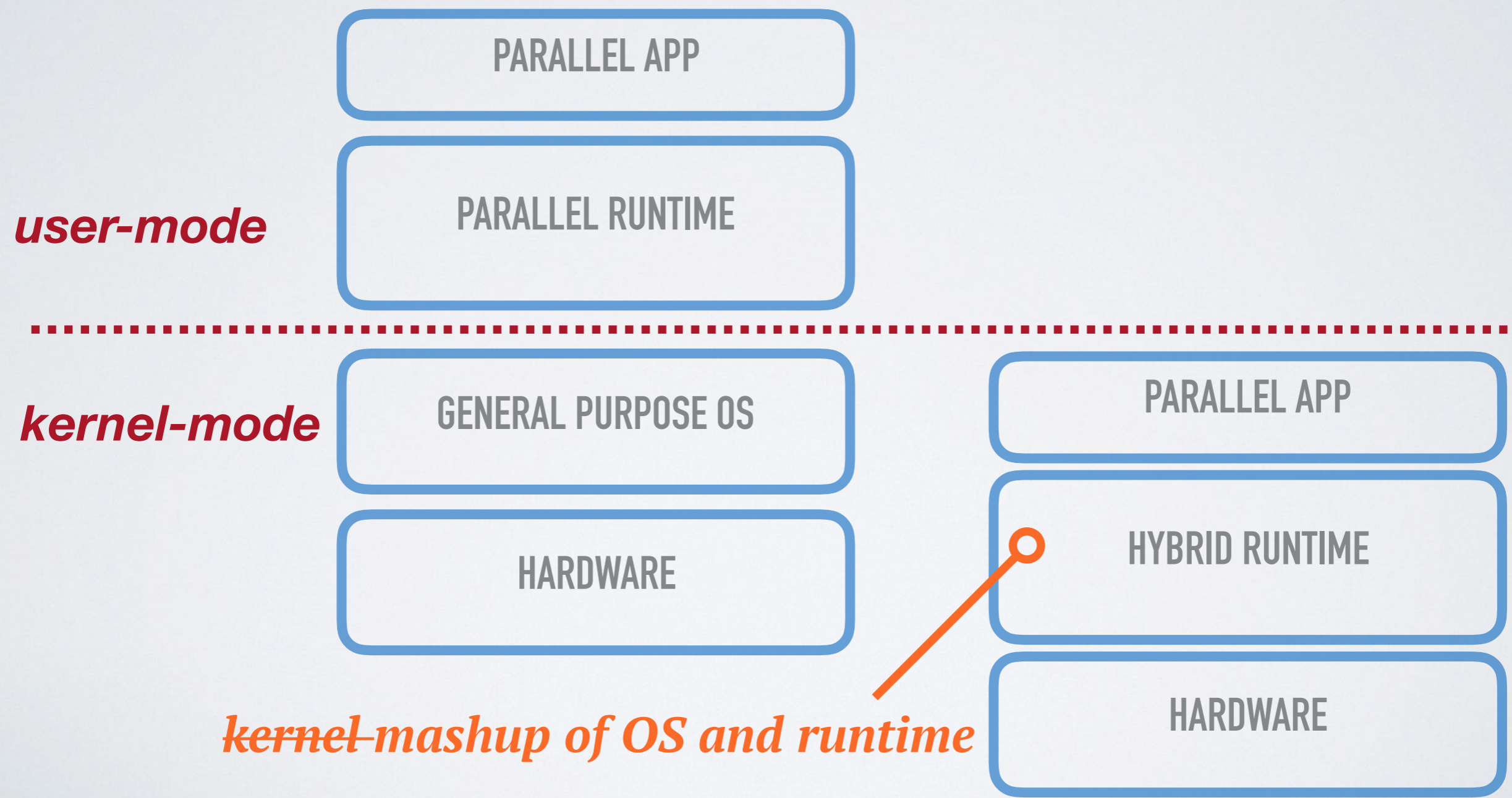
**runtime not limited to abstractions exposed by
syscall interface**

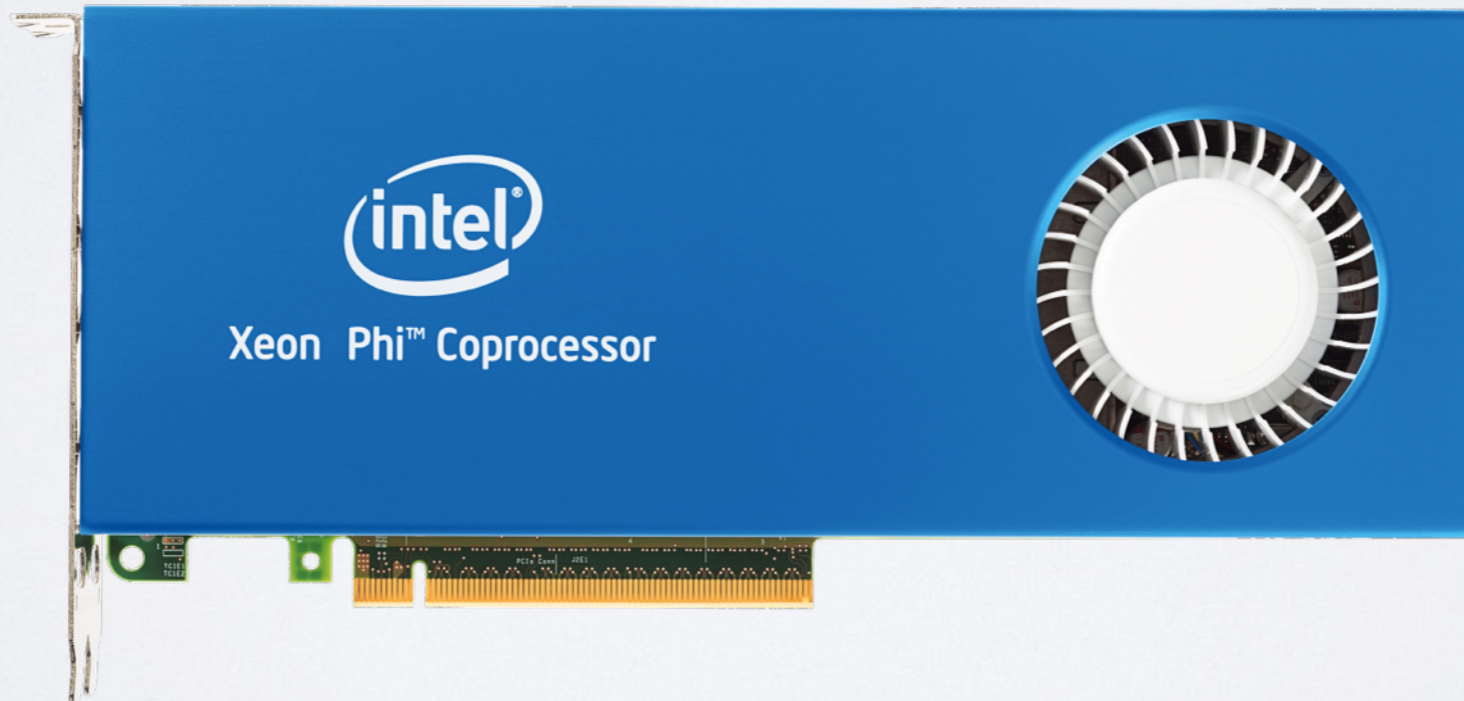
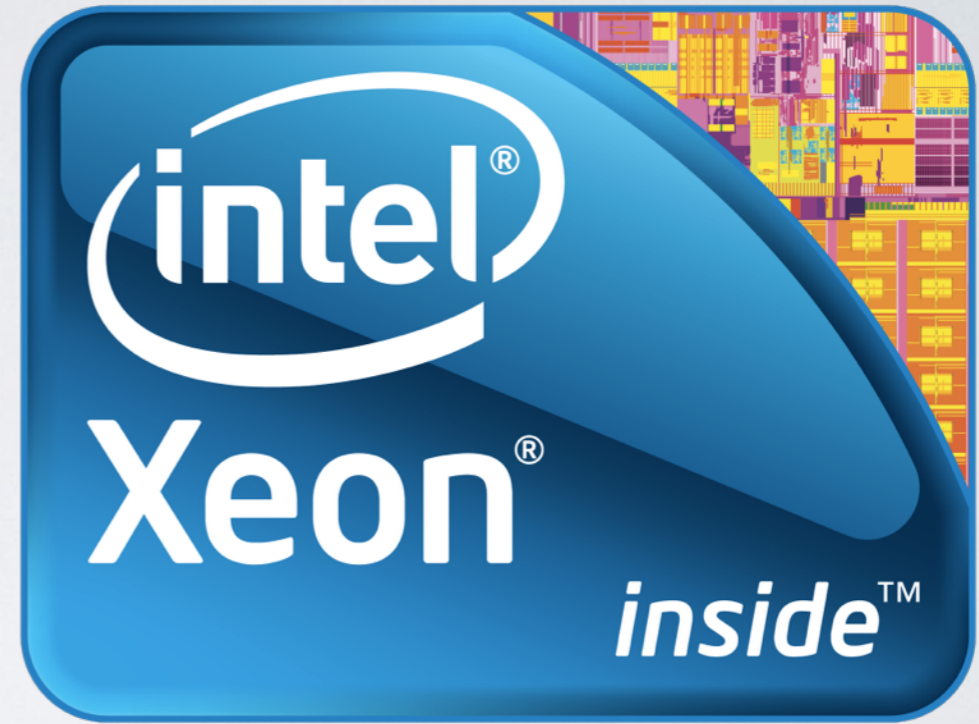
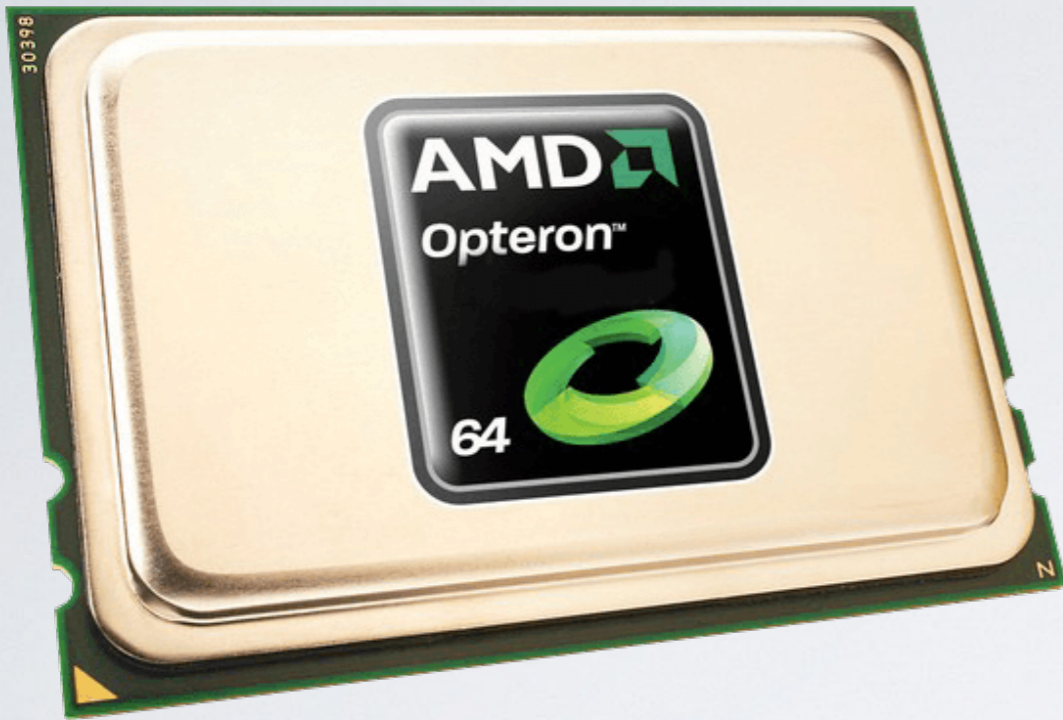
**opportunity for leveraging privileged
HW features**

CURRENT OS/R MODEL



THE HYBRID RUNTIME







Click on any of the above.

NESL: A Parallel Programming Language

NESL

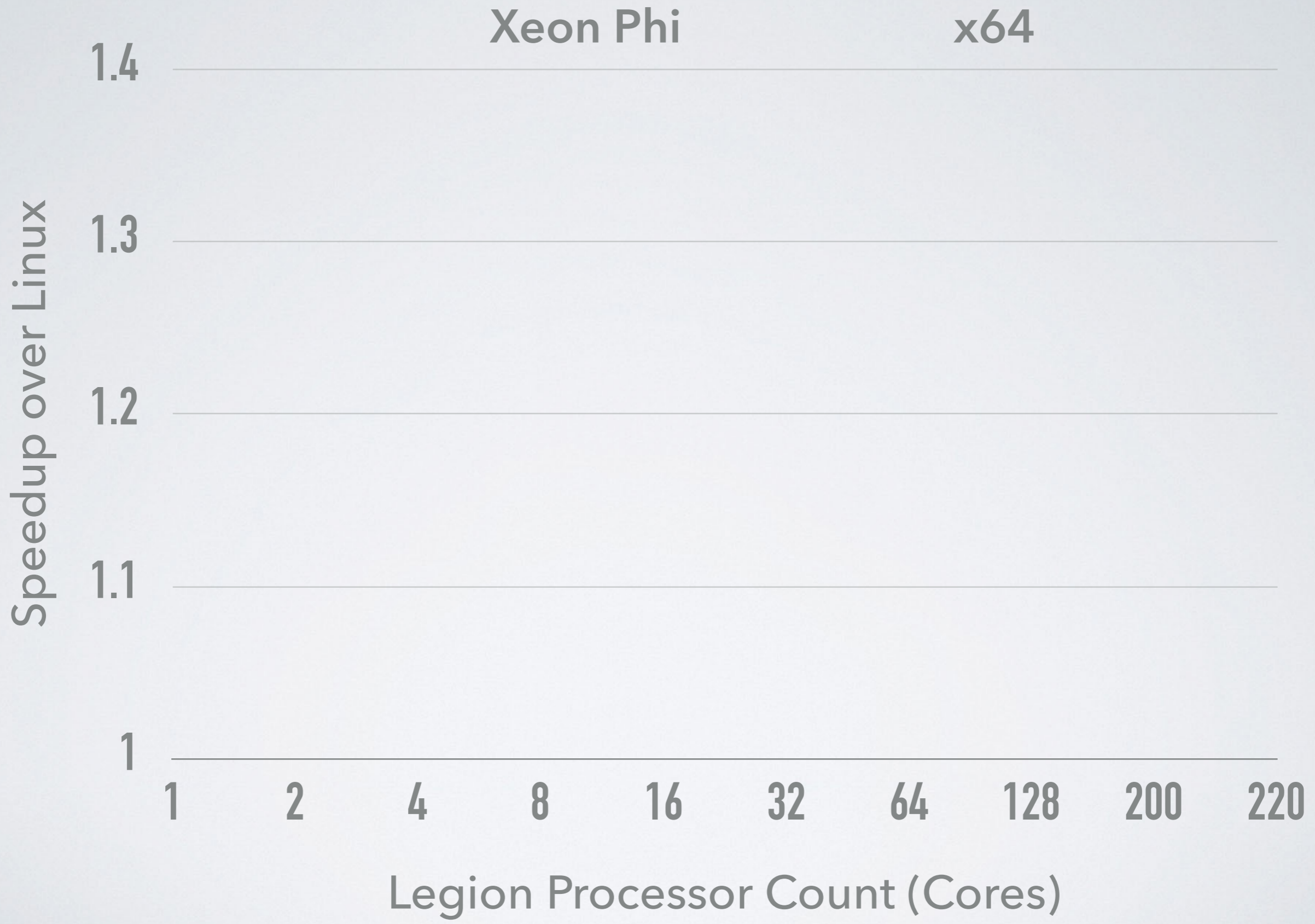
NDPC

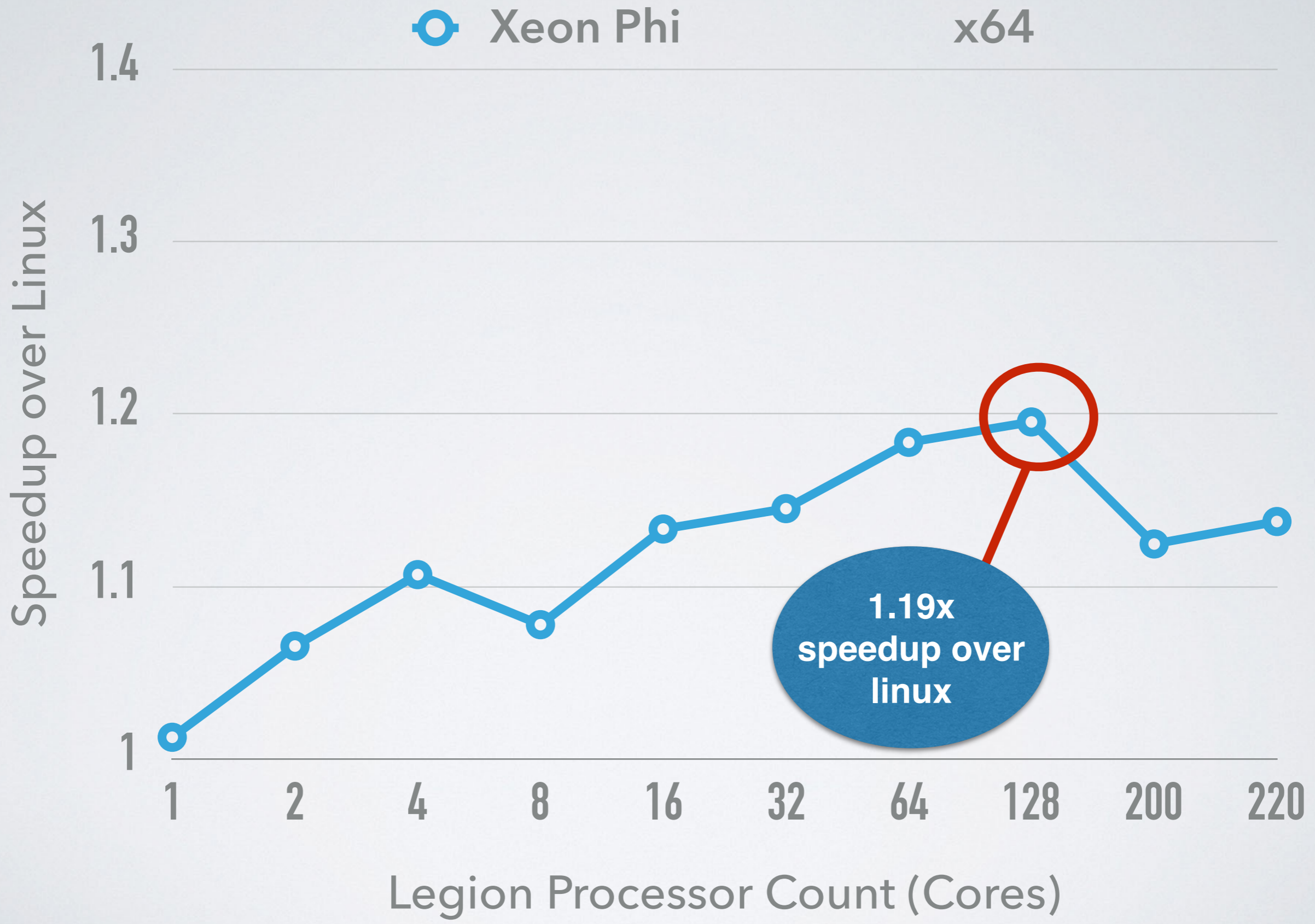


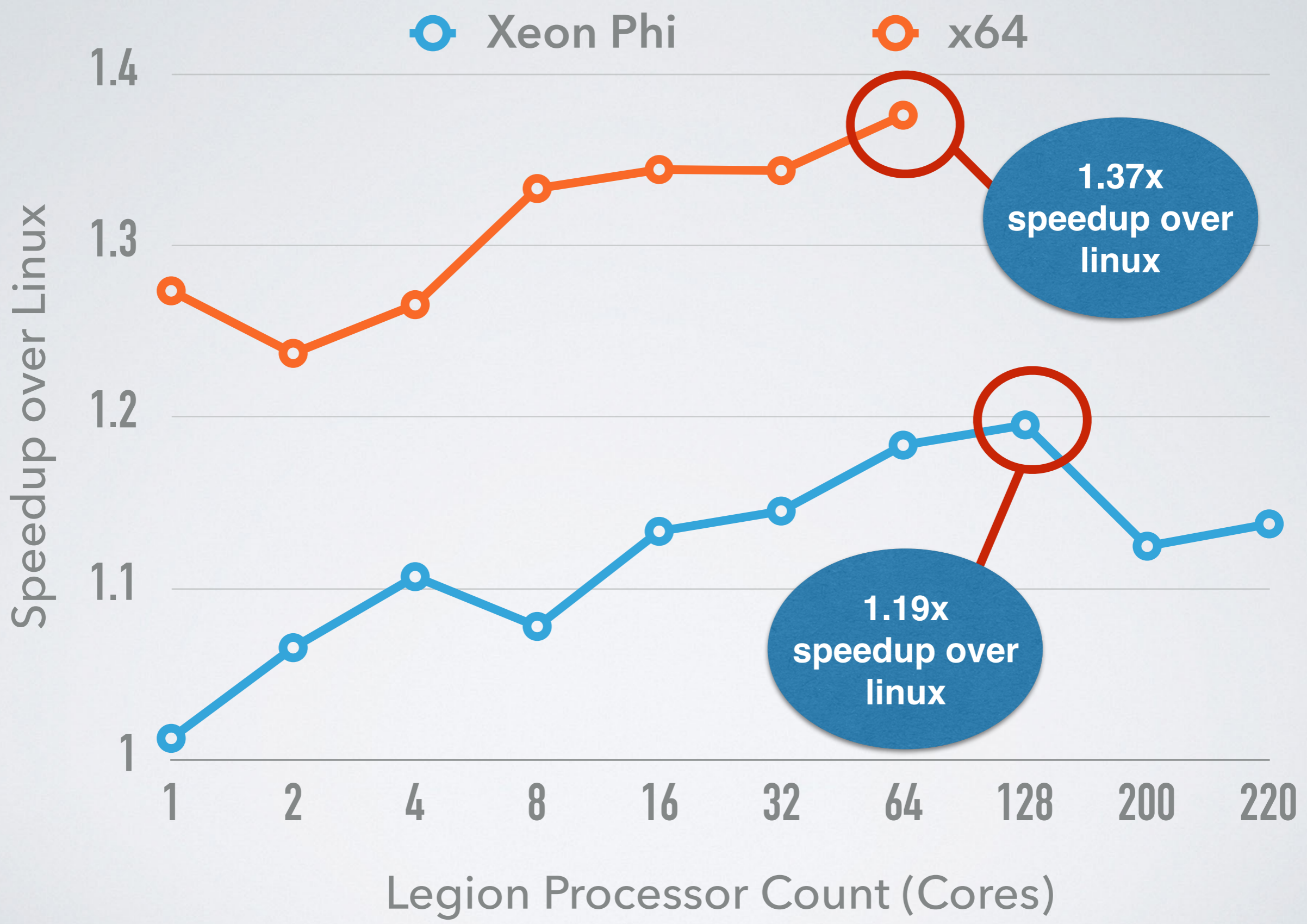
Racket

Legion





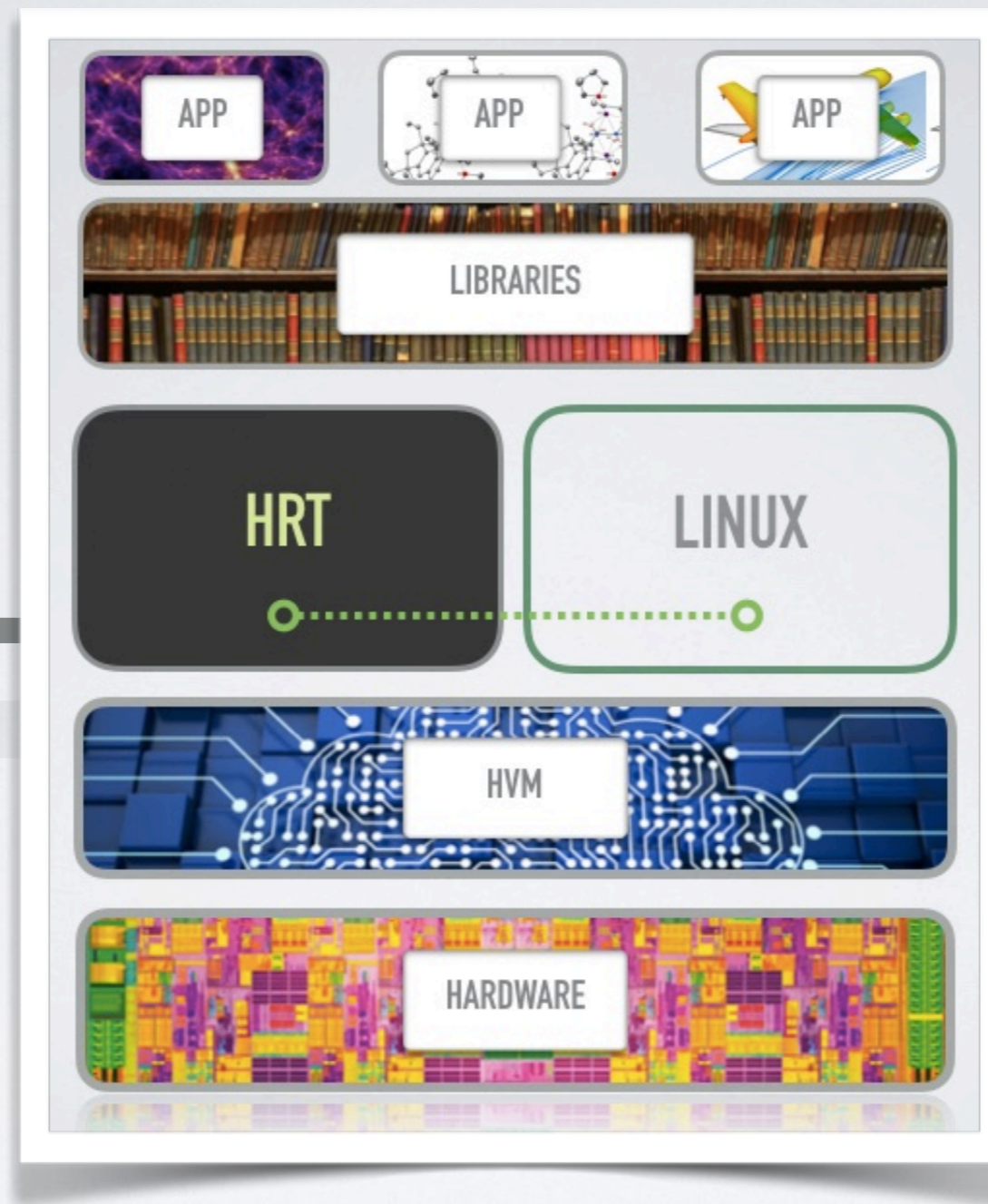




TWO ENABLING TOOLS



kernel framework



bridge HRT with legacy OS

OUTLINE

Background/Overview

○ Nautilus

Deployment Models

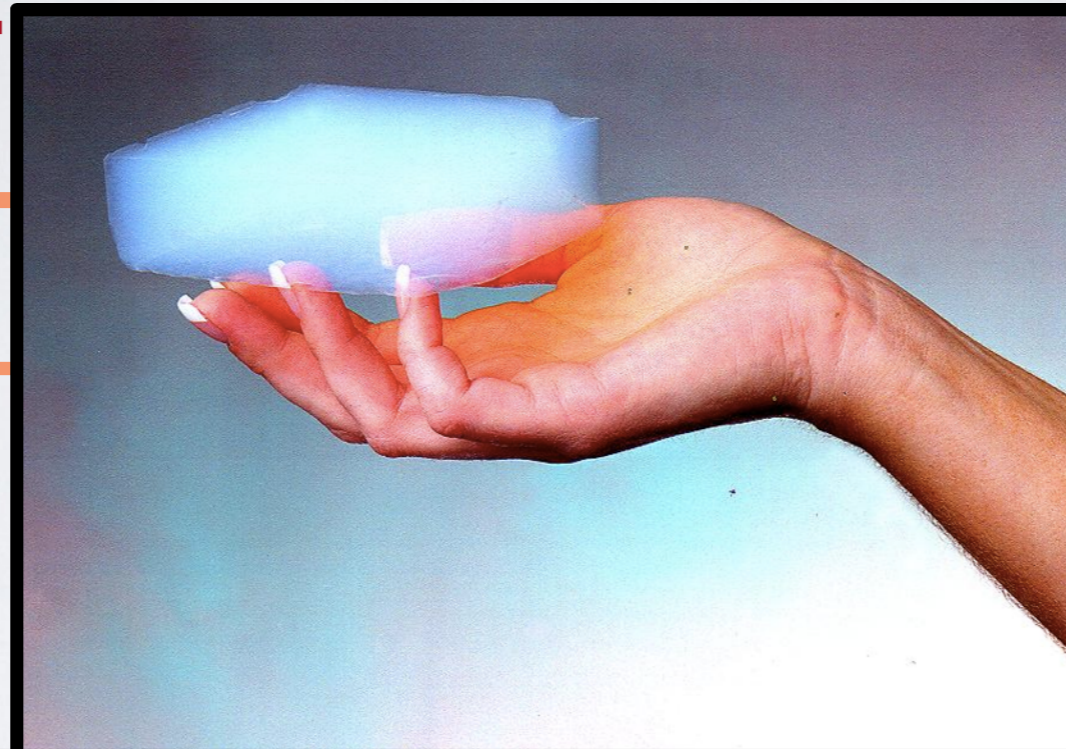
Hybrid Virtual Machine

Multiverse & Future Work

NAUTILUS

user-mode

kernel-mode



aerokernel

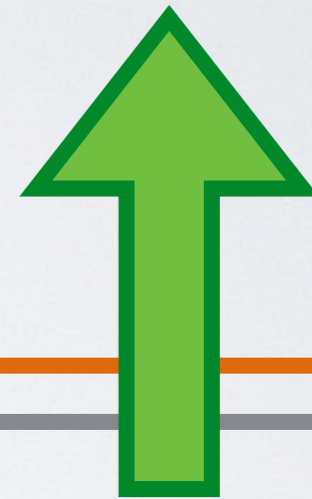
HARDWARE

Nautilus primitives & utilities (HRT can use or not use any of them)

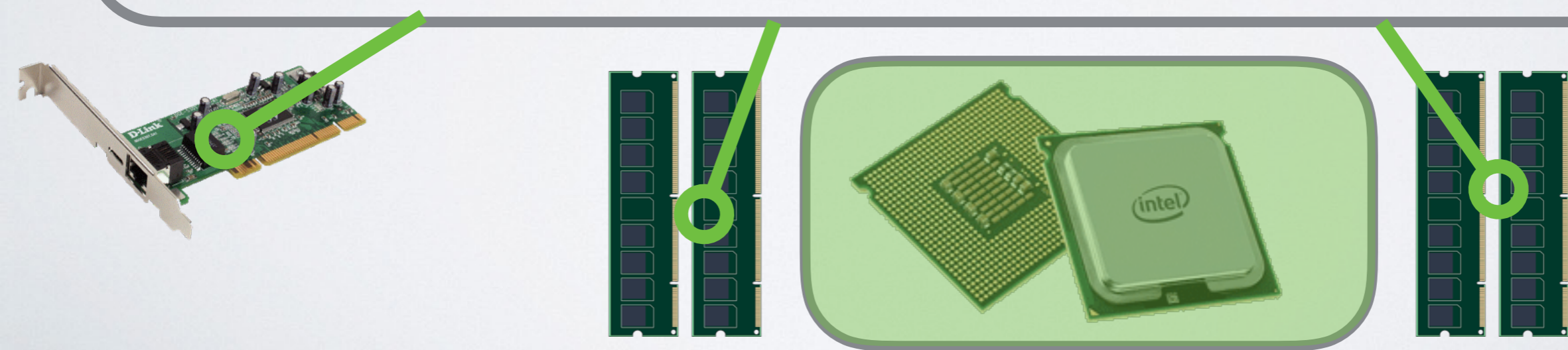
Nautilus

under the hood

Parallel Runtime System



Nautilus



**kernel primitives should be
SIMPLE and FAST**

**runtime developer can easily
reason about them!**

start with familiar interfaces

threads

condition variables

mutexes/locks

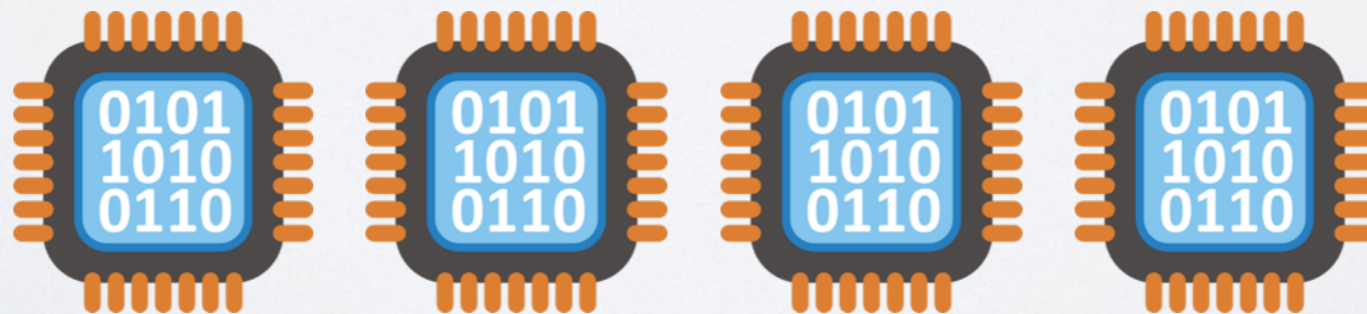
memory management

fork/join parallelism

**map runtime's logical
view of machine
Parallel Runtime System
onto physical HW**



logical CPUs



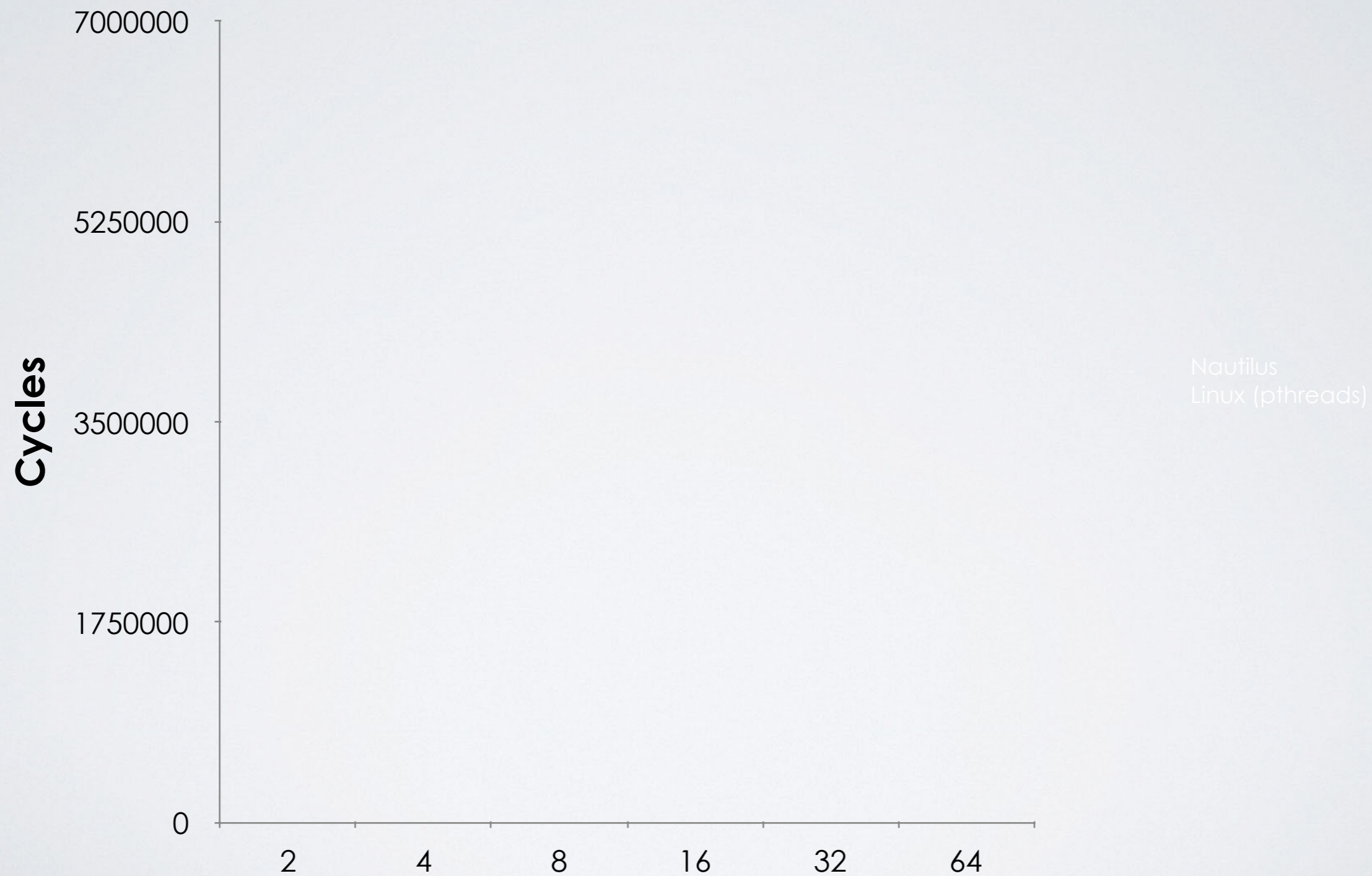
threads

unified, shared address space

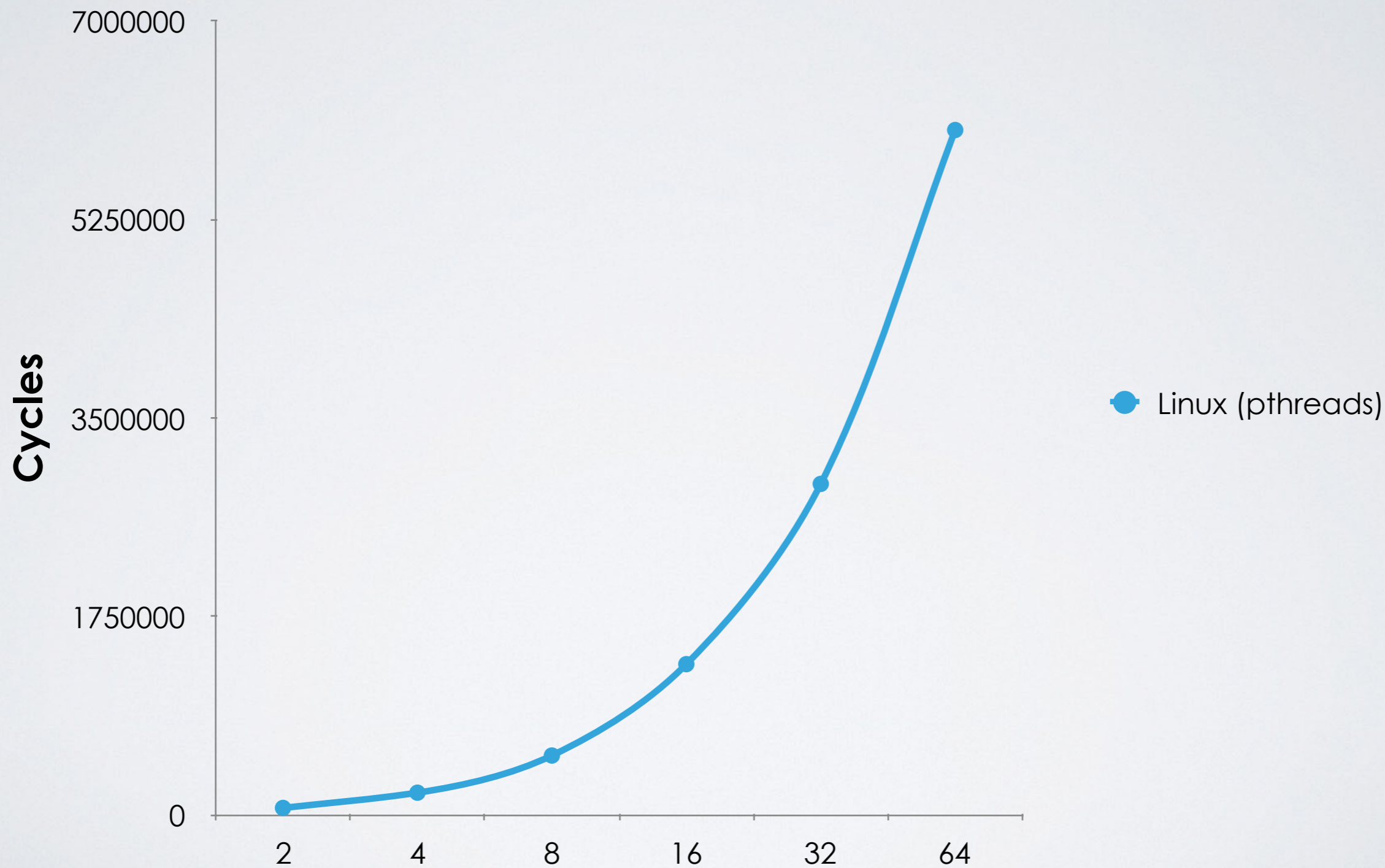
**threads can operate preemptively or
cooperatively***

***for runtimes that require more determinism**

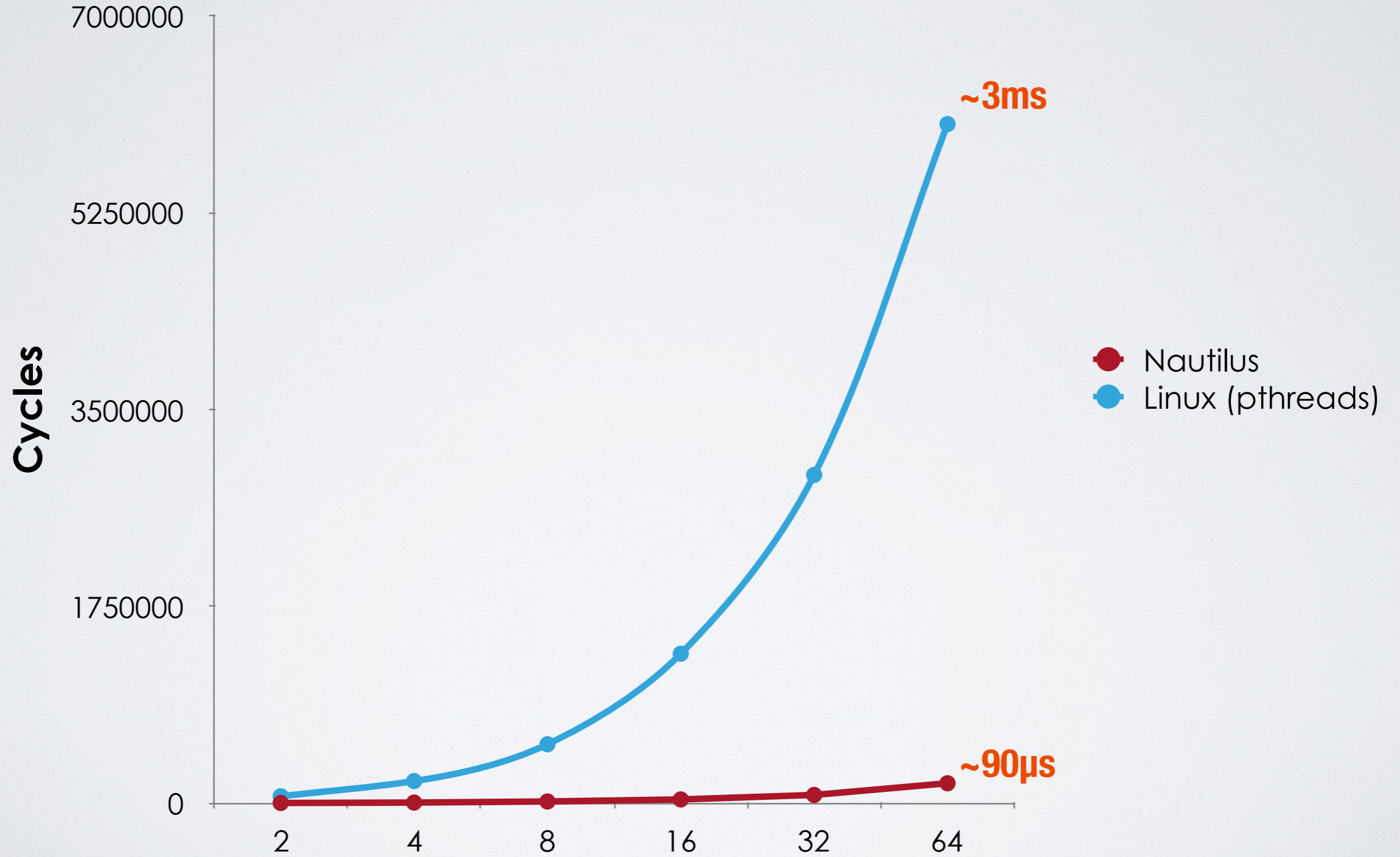
thread creation is FAST



thread creation is FAST



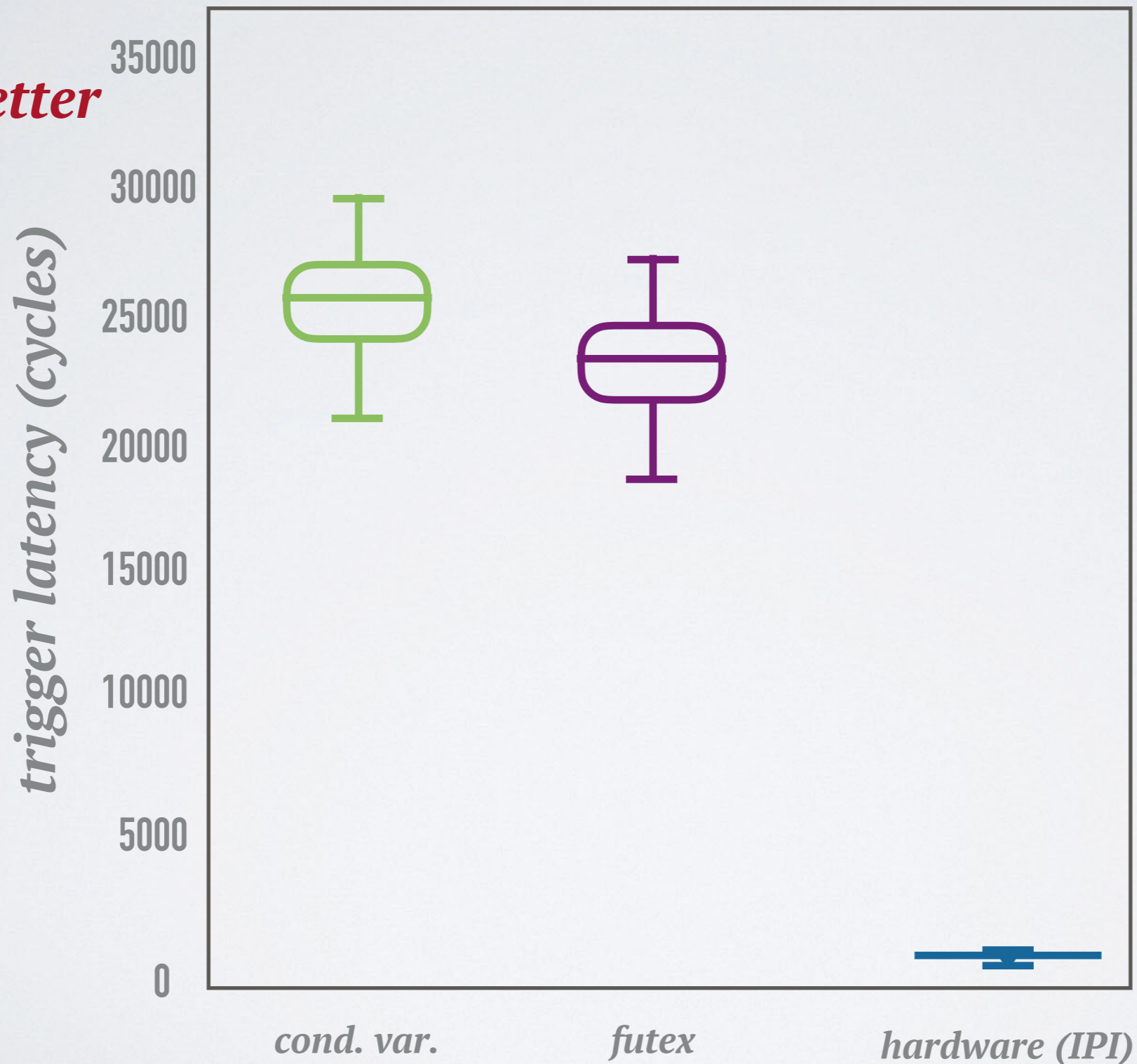
thread creation is FAST



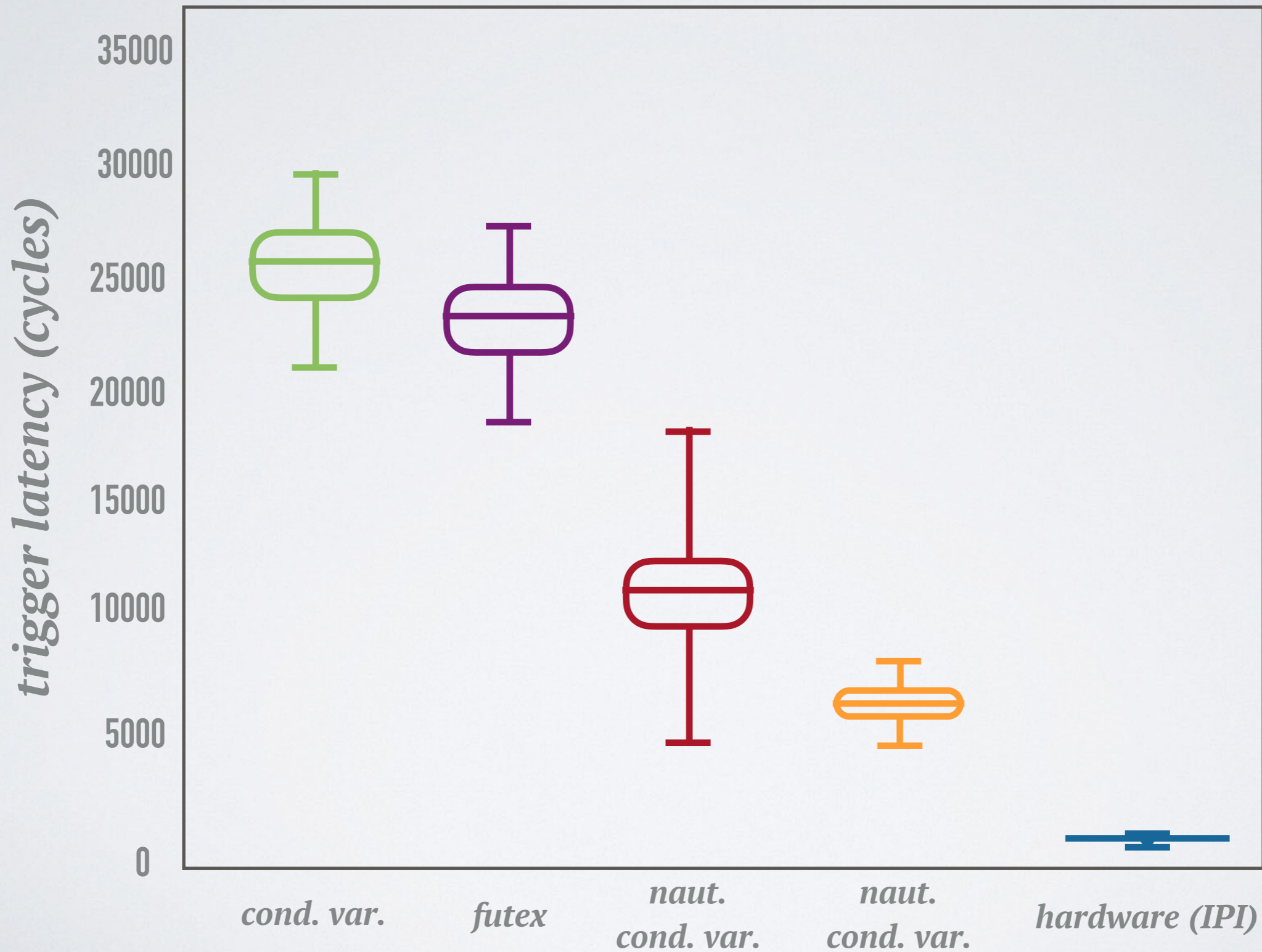
software events are **SLOW**

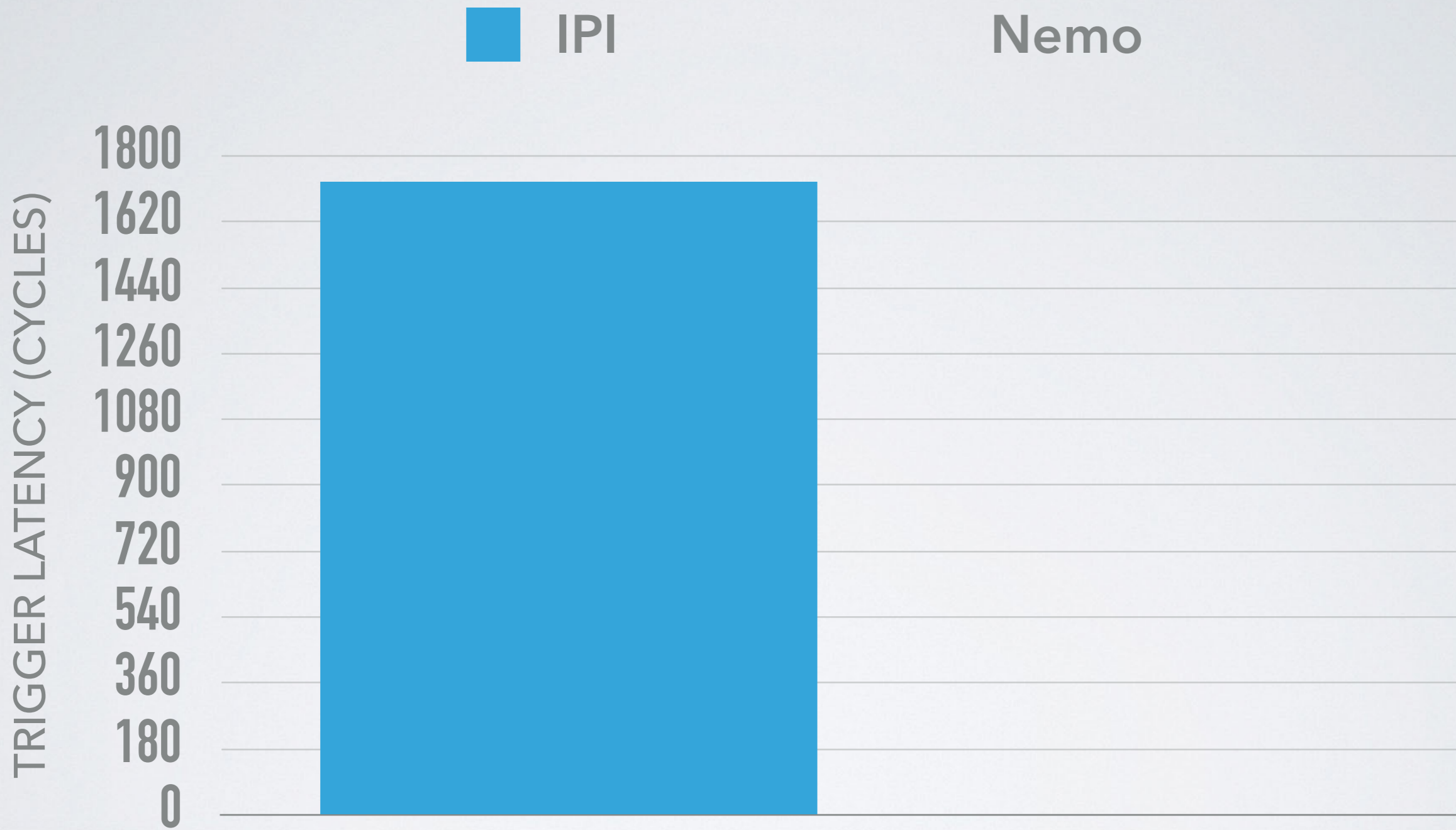
24

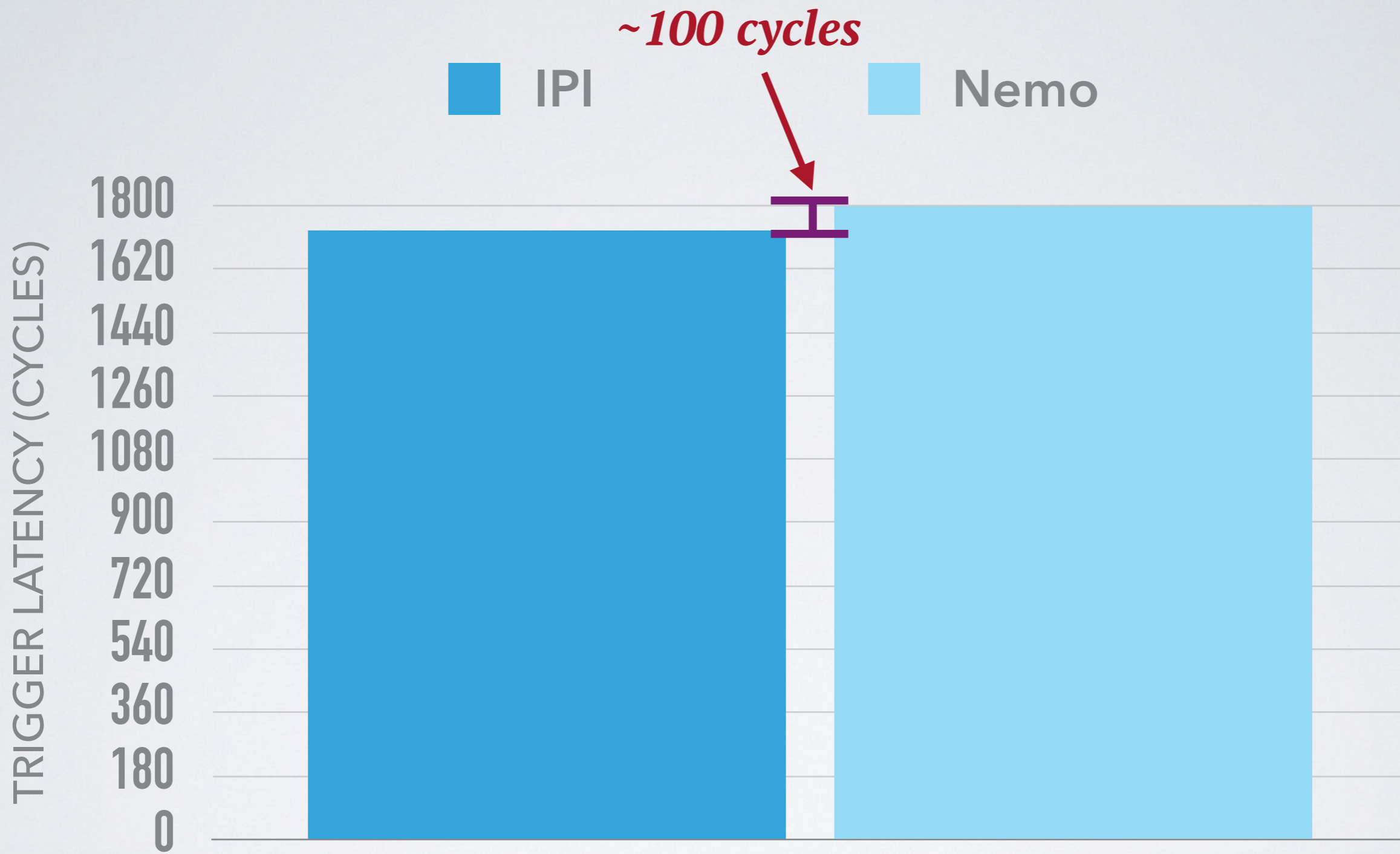
lower is better



*nautilus events triggers are **FASTER***

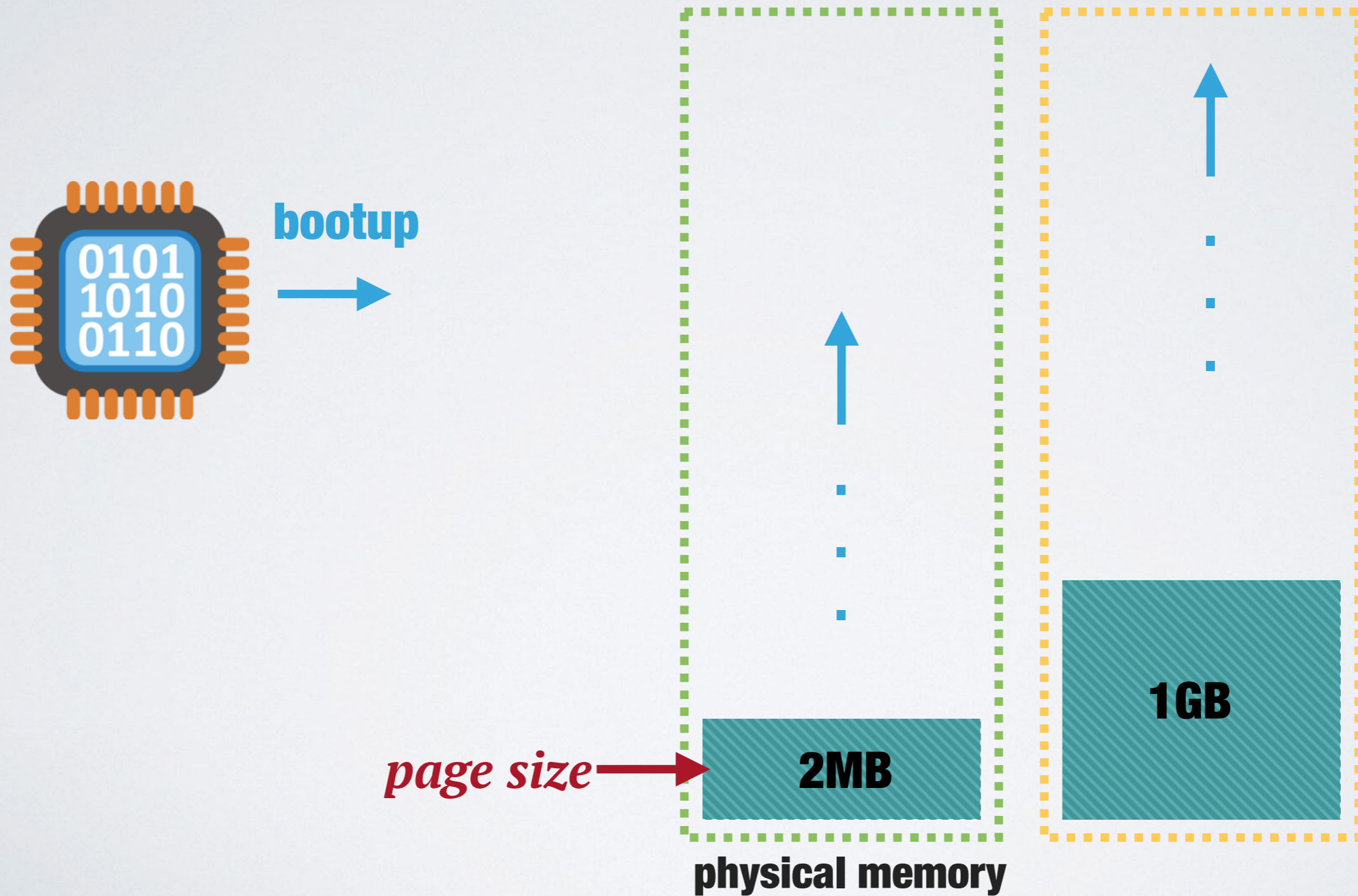






memory

static identity map



eliminates expensive page faults

reduces TLB misses + shutdowns

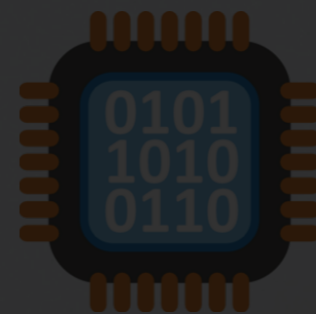
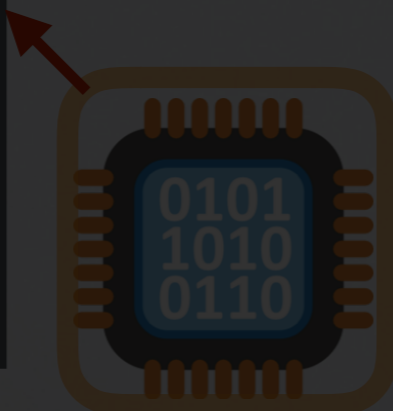
*increases performance under
virtualization*

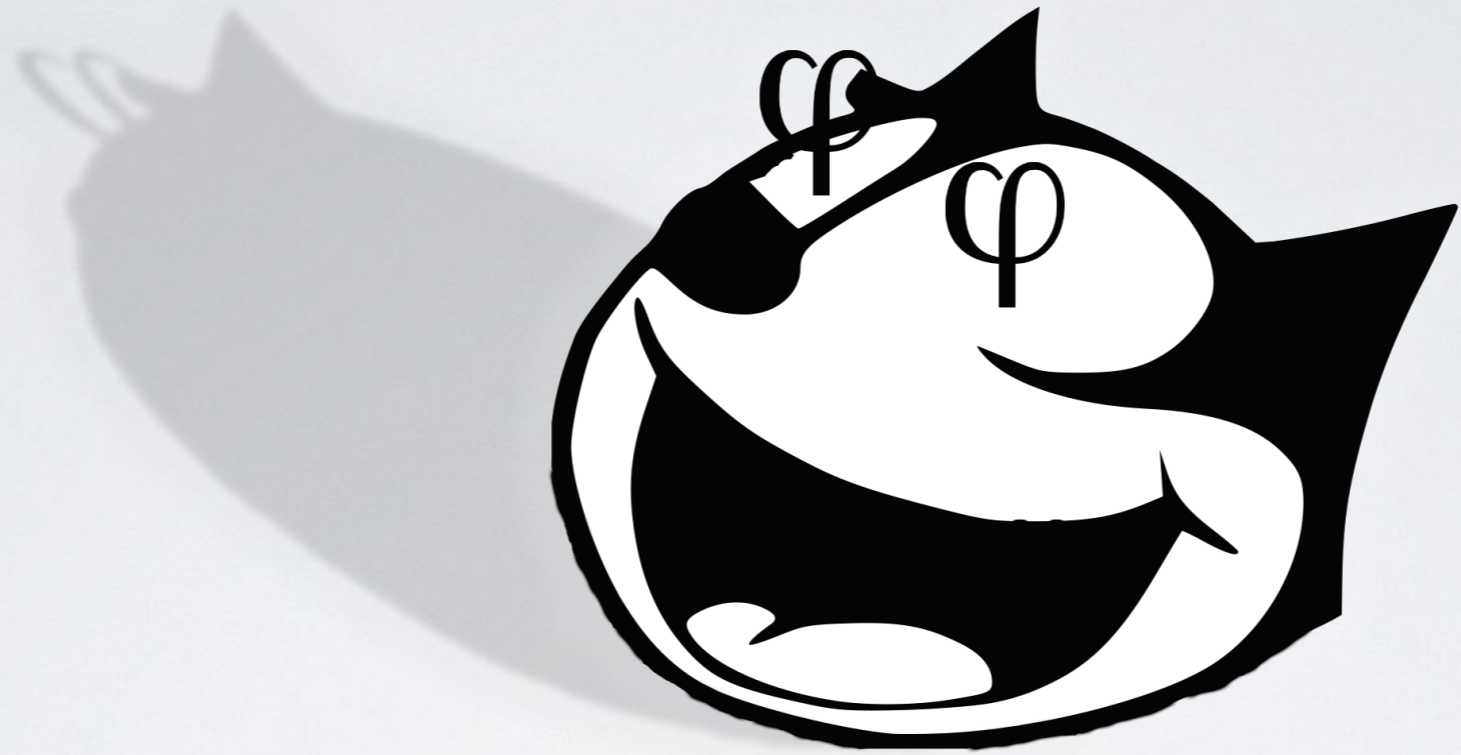
NUMA

~~first-touch~~

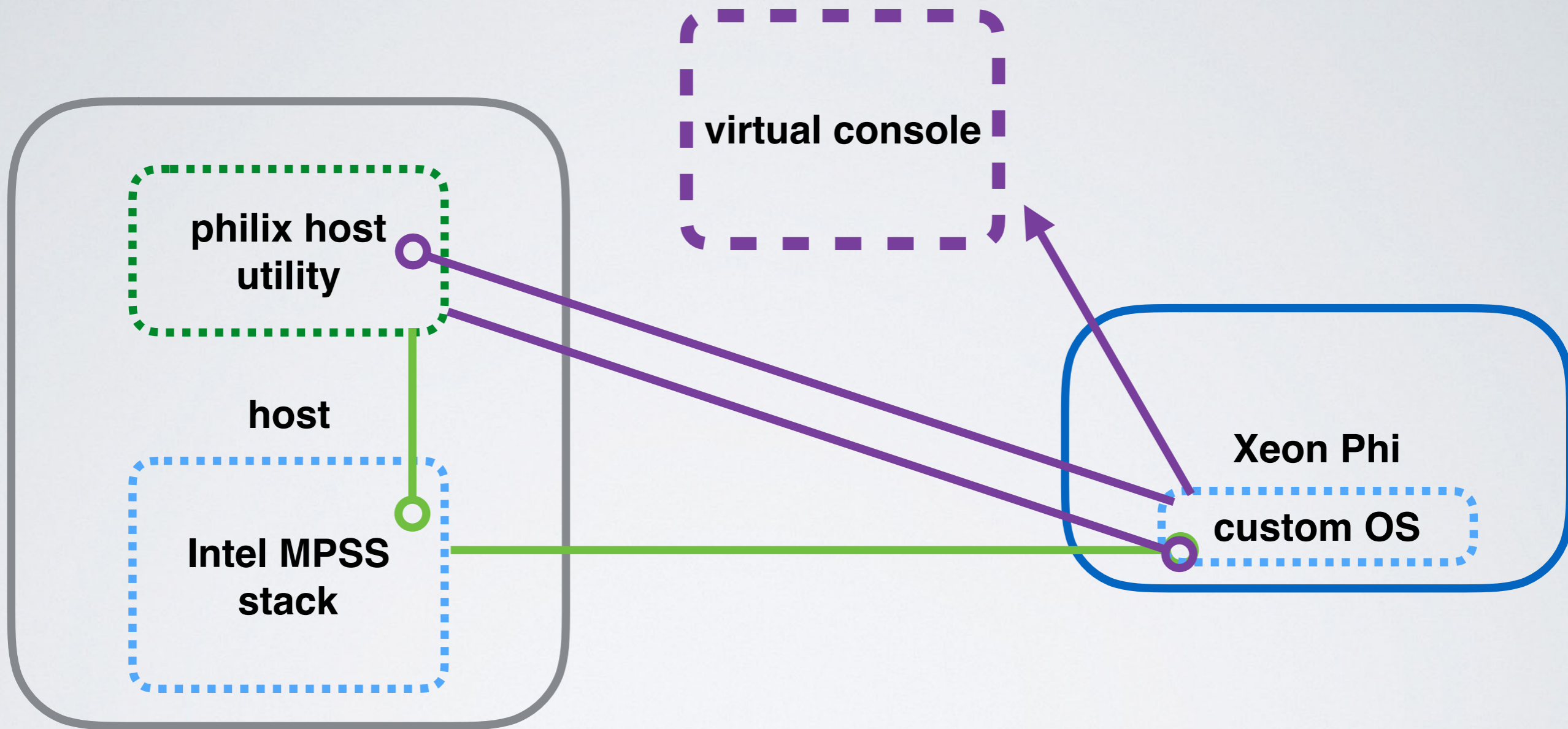
NUMA domain NUMA domain

runtime has FULL control
over thread placement and
memory layout





philix



Nautilus is fairly small

Nautilus

25,000 lines of mostly C

Legion

43,000 lines of C++

Additions for Legion

800 lines of C/C++

Additions for Xeon Phi

1350 lines of C

What do we lose?

```
[root@localhost ~]# ping -q fa.wikipedia.org
PING text.pmtpa.wikimedia.org (208.80.152.2) 56(84) bytes of data.
^C
--- text.pmtpa.wikimedia.org ping statistics ---
1 packets transmitted, 1 received, 0% packet loss, time 0ms
rtt min/avg/max/mdev = 540.528/540.528/540.528/0.000 ms
```

```
[root@localhost ~]# pwd
/root
[root@localhost ~]# cd /var
[root@localhost var]# ls -la
```

```
total 72
drwxr-xr-x. 18 root root 4096 Jul 30 22:43 .
drwxr-xr-x. 23 root root 4096 Sep 14 20:42 ..
drwxr-xr-x.  2 root root 4096 May 14 00:15 account
drwxr-xr-x. 11 root root 4096 Jul 31 22:26 cache
drwxr-xr-x.  3 root root 4096 May 18 16:03 db
drwxr-xr-x.  3 root root 4096 May 18 16:03 empty
drwxr-xr-x.  2 root root 4096 May 18 16:03 games
drwxr-xr-x.  2 root gdm 4096 Jul  2 18:39 gdm
drwxr-xr-x. 28 root root 4096 May 18 16:03 lib
drwxr-xr-x.  2 root root 4096 May 18 16:03 local
lrwxrwxrwx.  1 root root    1 May 14 00:12 lock -> ../run/lock
drwxr-xr-x. 14 root root 4096 Sep 14 20:42 log
lrwxrwxrwx.  1 root root    10 Jul 30 22:43 mail -> spool/mail
drwxr-xr-x.  2 root root 4096 May 18 16:03 nis
drwxr-xr-x.  2 root root 4096 May 18 16:03 opt
drwxr-xr-x.  2 root root 4096 May 18 16:03 preserve
drwxr-xr-x.  2 root root 4096 Jul  1 22:11 report
lrwxrwxrwx.  1 root root    6 May 14 00:12 run -> ../run
drwxr-xr-x. 14 root root 4096 May 18 16:03 spool
drwxrwxrwt.  4 root root 4096 Sep 12 23:50 tmp
drwxr-xr-x.  2 root root 4096 May 18 16:03 yp
```

```
[root@localhost var]# yum search wiki
Loaded plugins: langpacks, presto, refresh-packagekit, remove-with-leaves
```

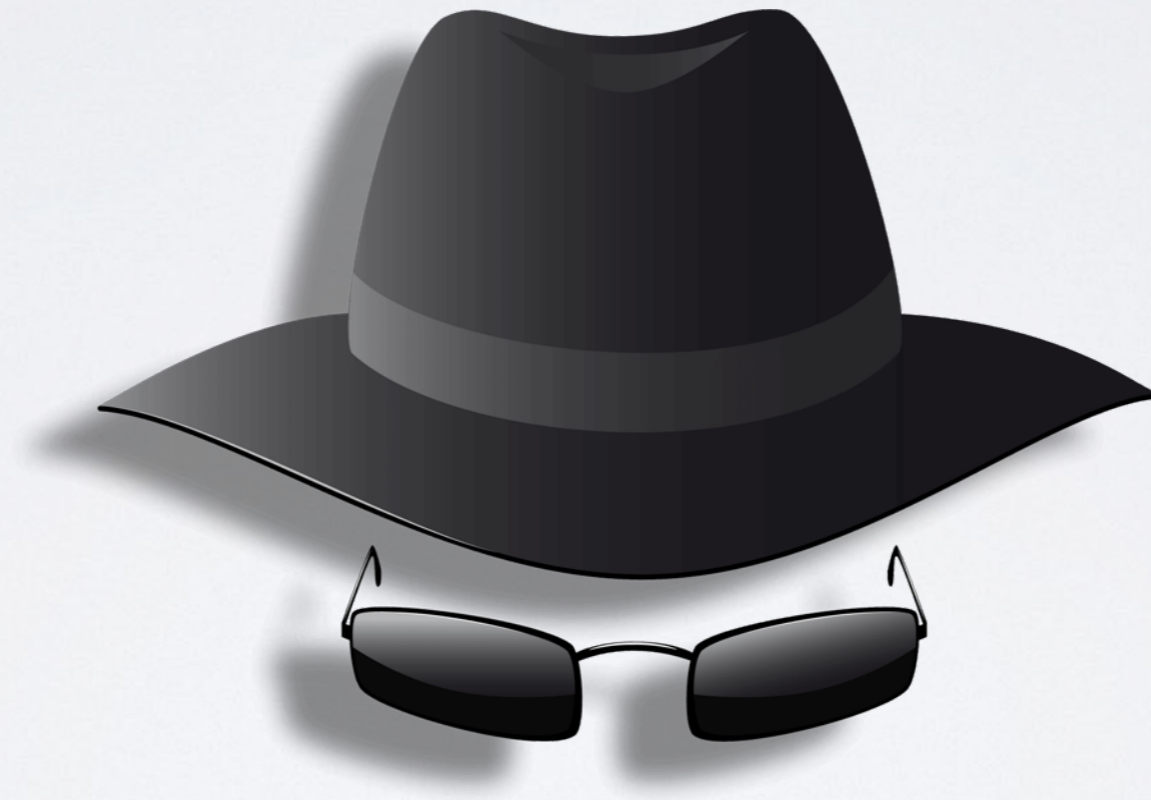
```
rpmfusion-free-updates | 2.7 kB | 00:00
rpmfusion-free-updates/primary_db | 206 kB | 00:04
rpmfusion-nonfree-updates | 2.7 kB | 00:00
updates/metalink | 5.9 kB | 00:00
updates | 4.7 kB | 00:00
updates/primary_db | 2.6 MB | 00:15 ETA
73% [=====] ] 62 kB/s
```

familiar environment

driver ecosystem



protection/isolation



OUTLINE

Background/Overview

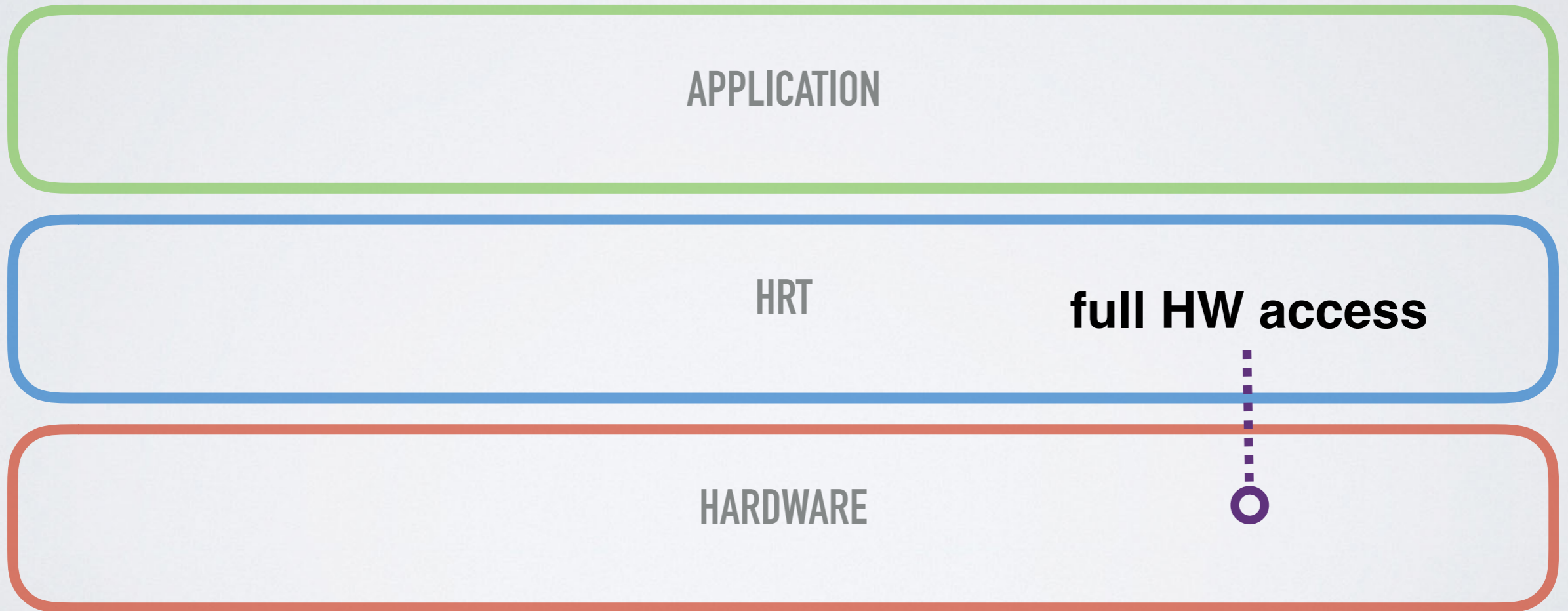
Nautilus

Deployment Models

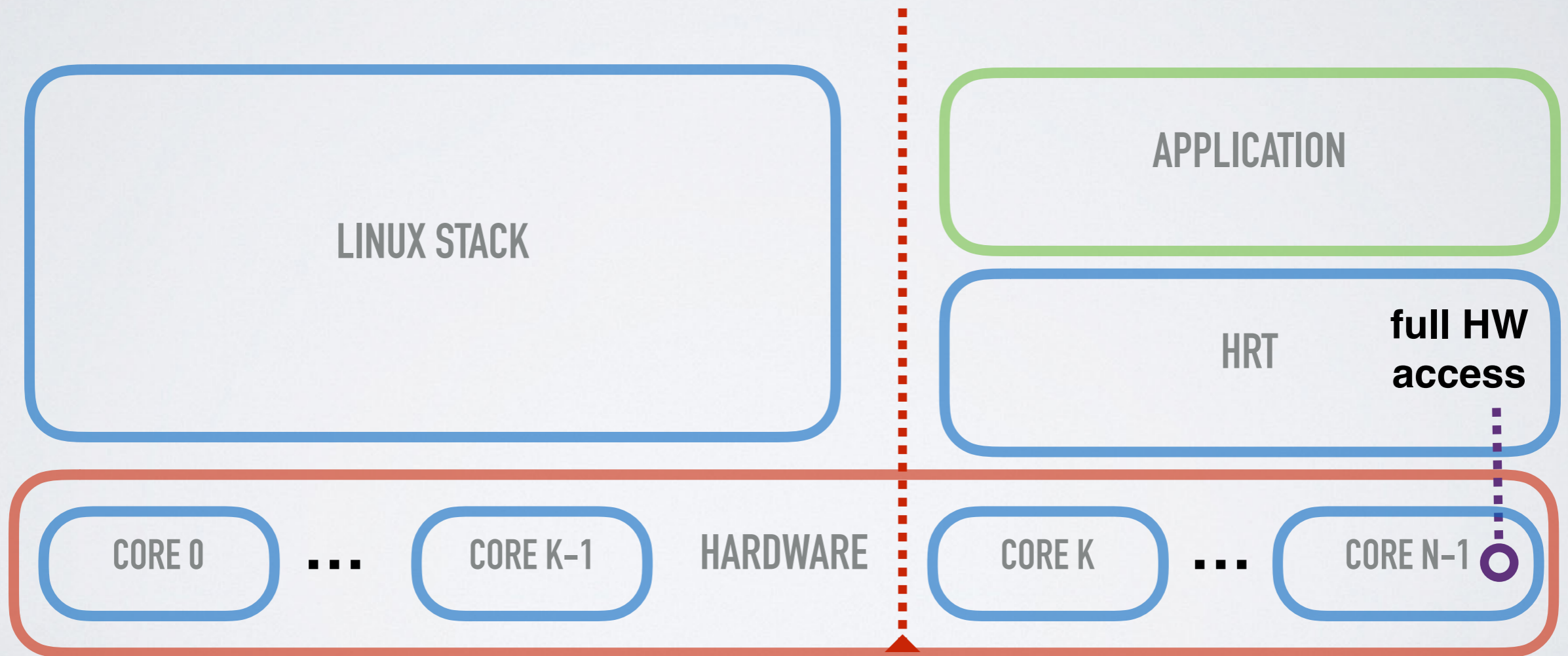
Hybrid Virtual Machine

Multiverse & Future Work

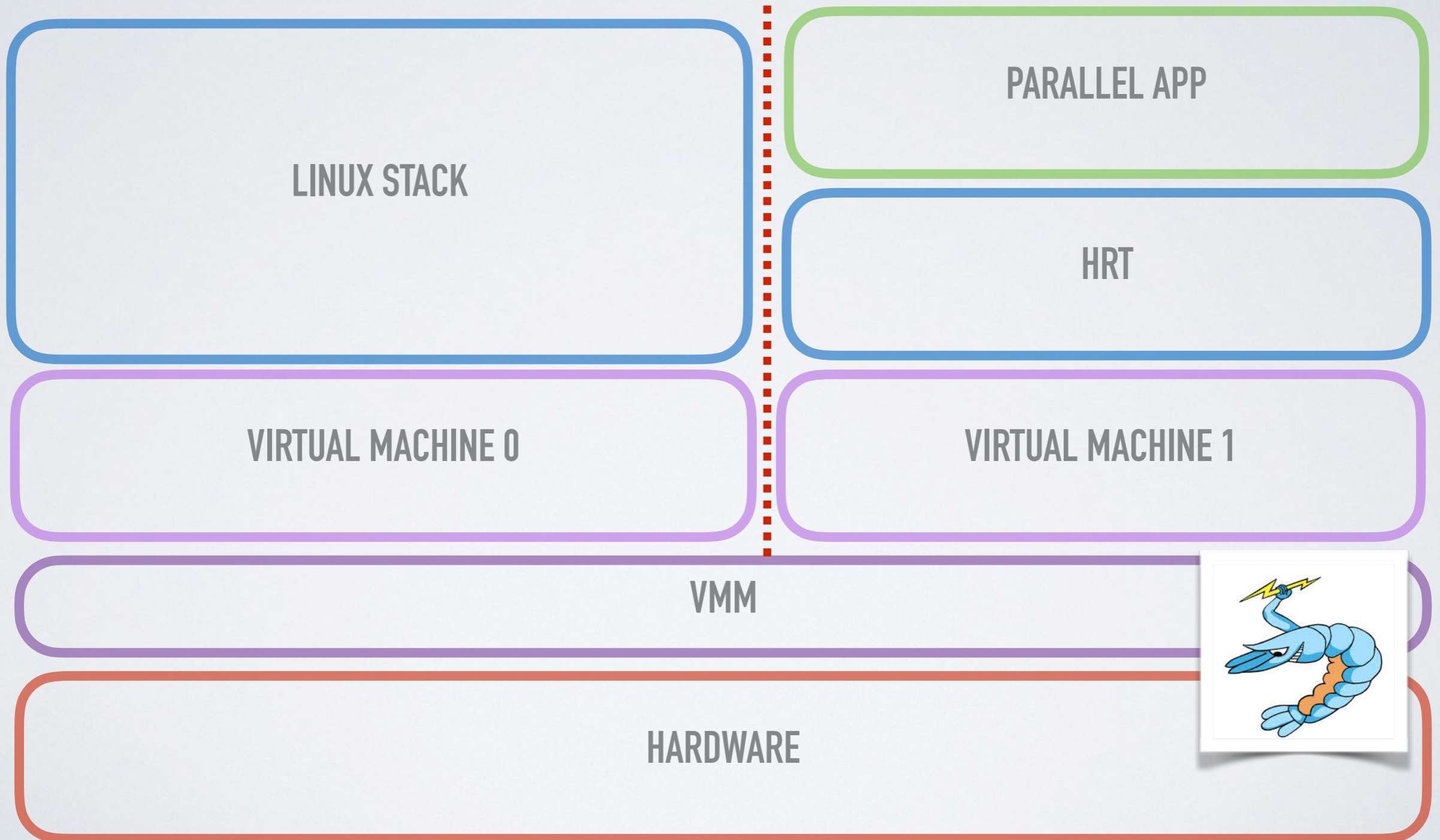
DEDICATED



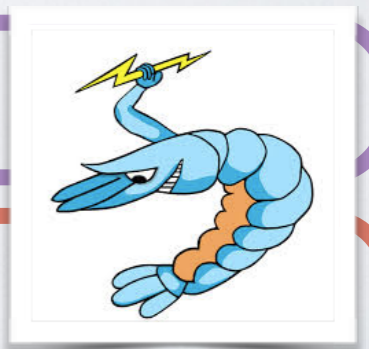
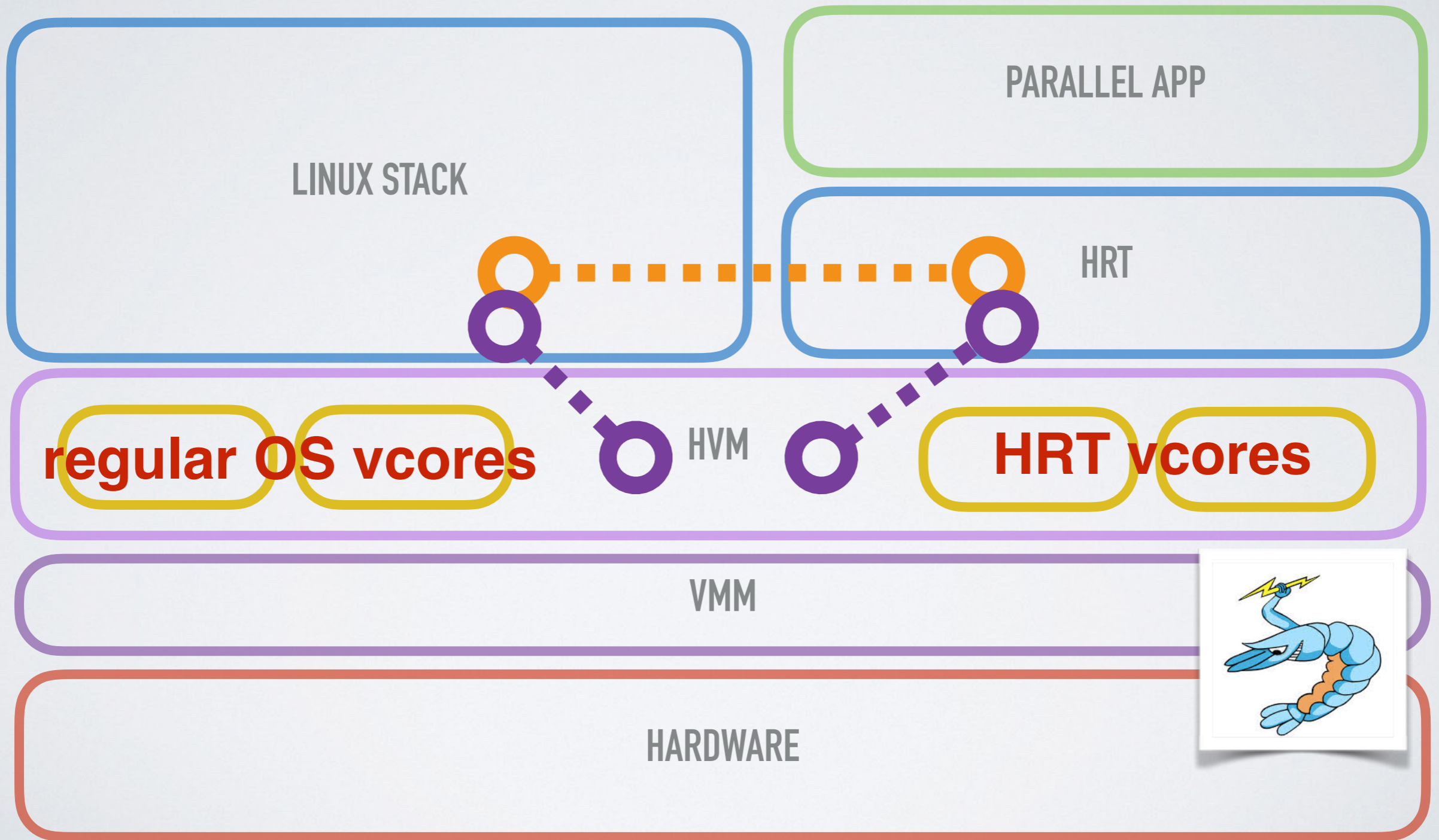
PARTITIONED



VM



Hybrid Virtual Machine



OUTLINE

Background/Overview

Nautilus

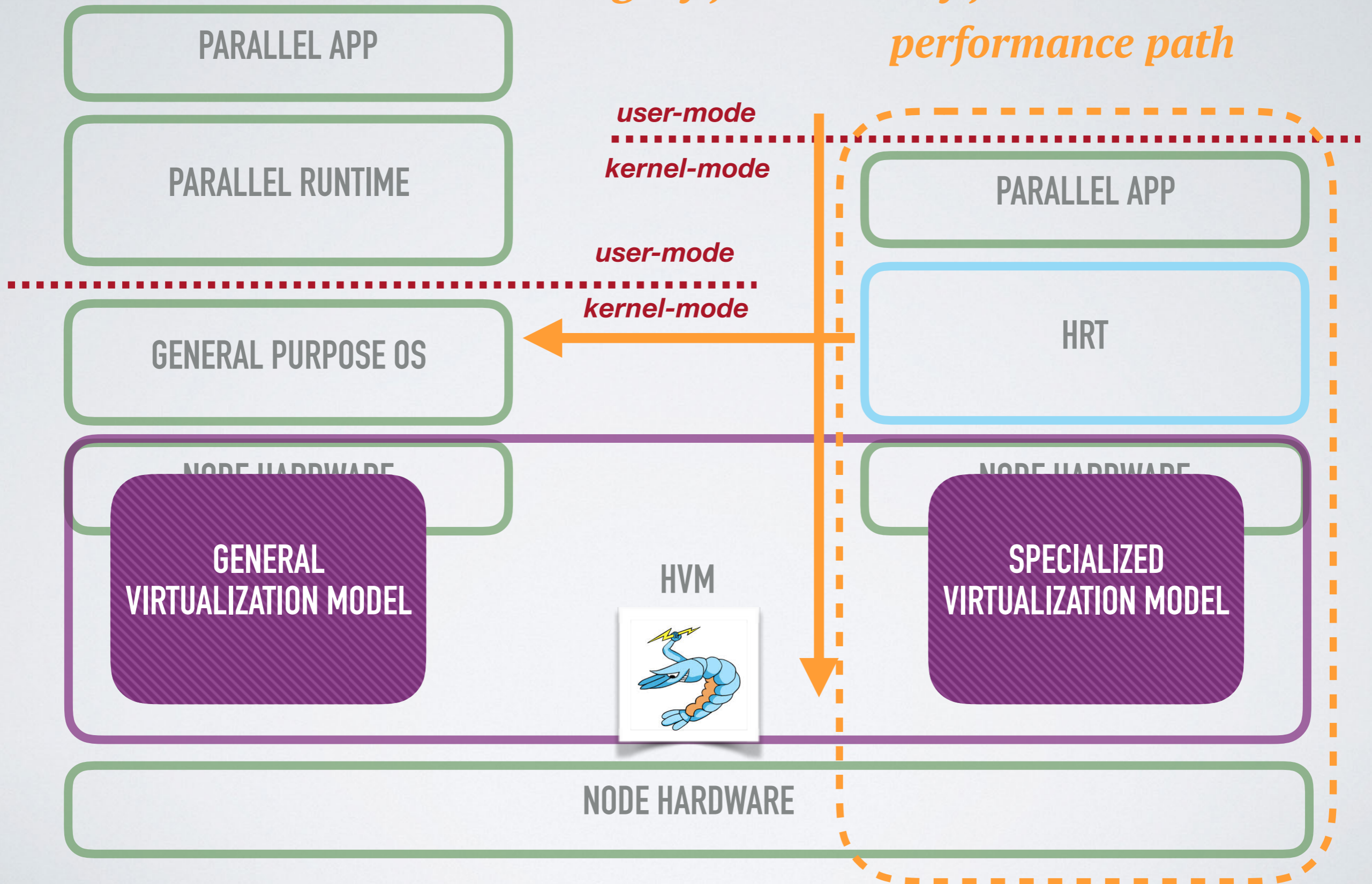
Deployment Models

○ Hybrid Virtual Machine

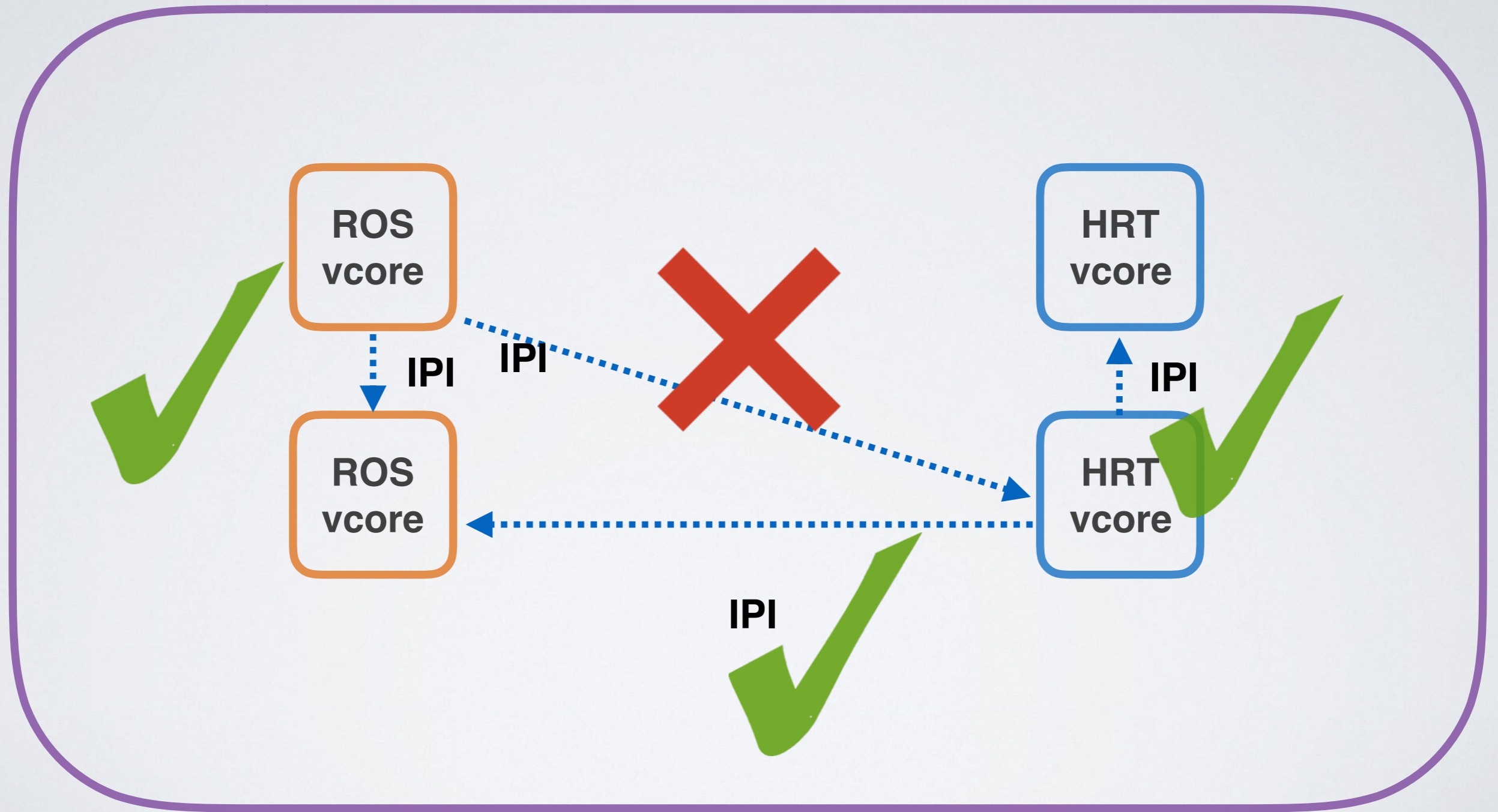
Multiverse & Future Work

regular OS (ROS)

*legacy functionality from ROS via HVM
performance path*

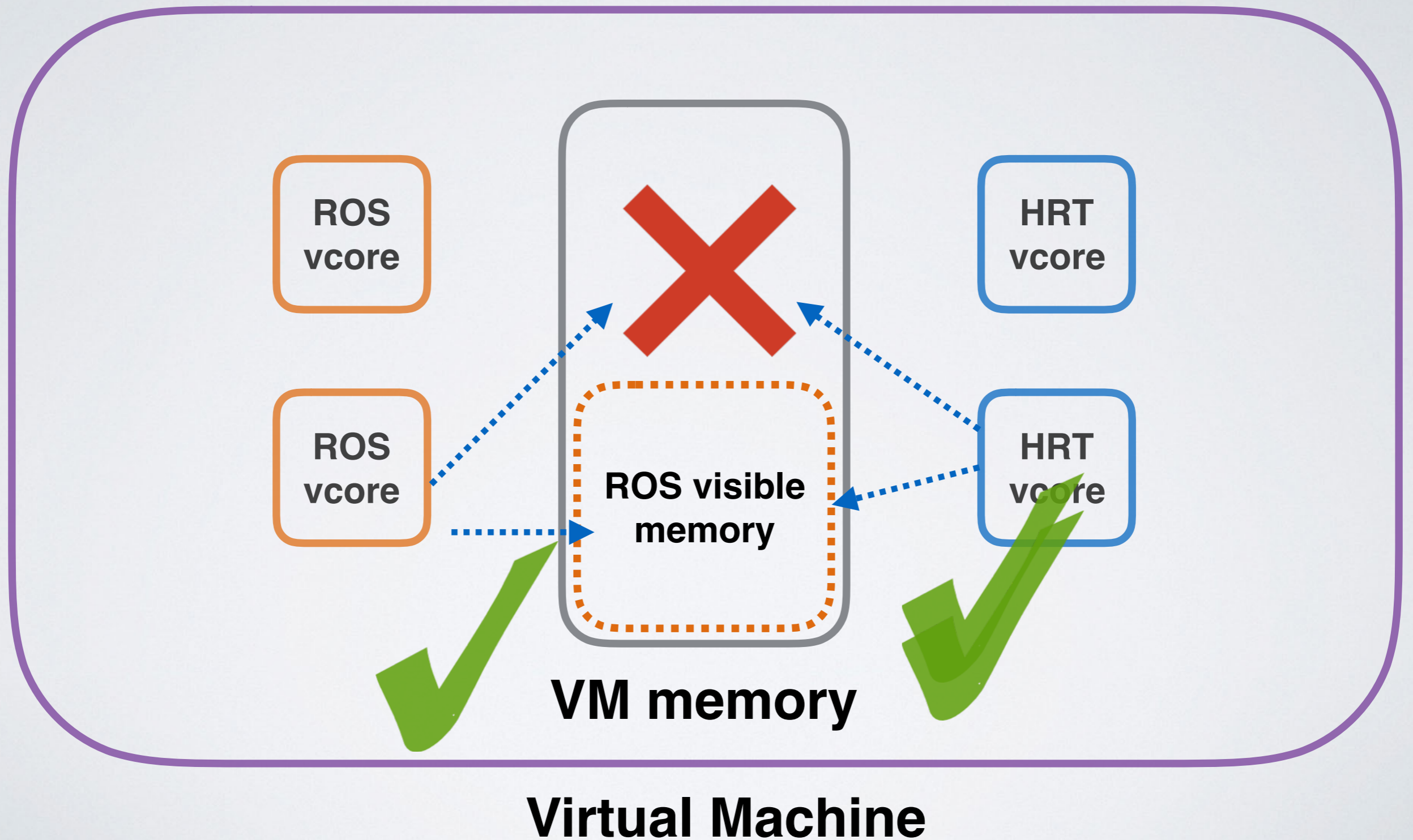


ROS/HRT setup

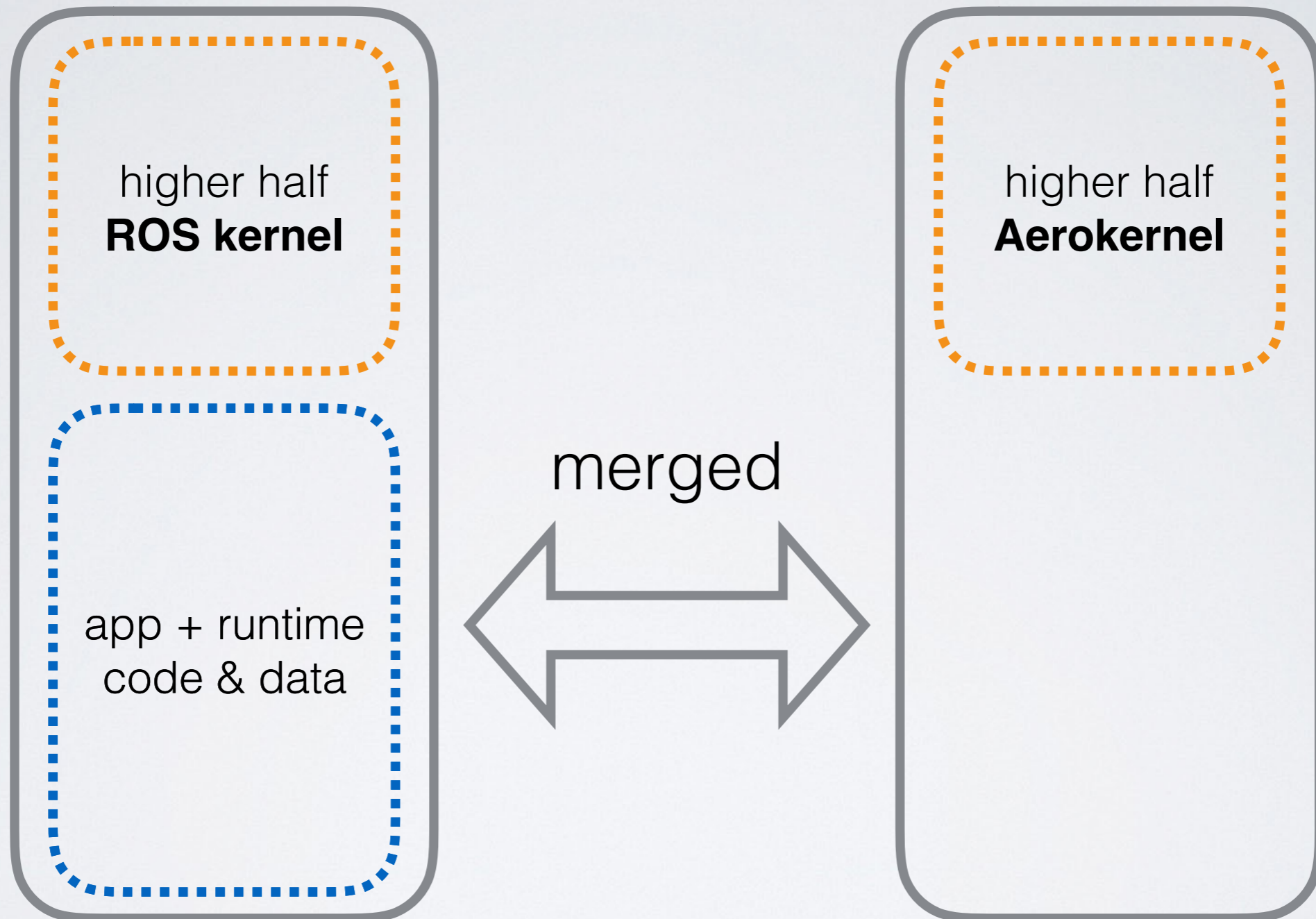


Virtual Machine

ROS/HRT setup



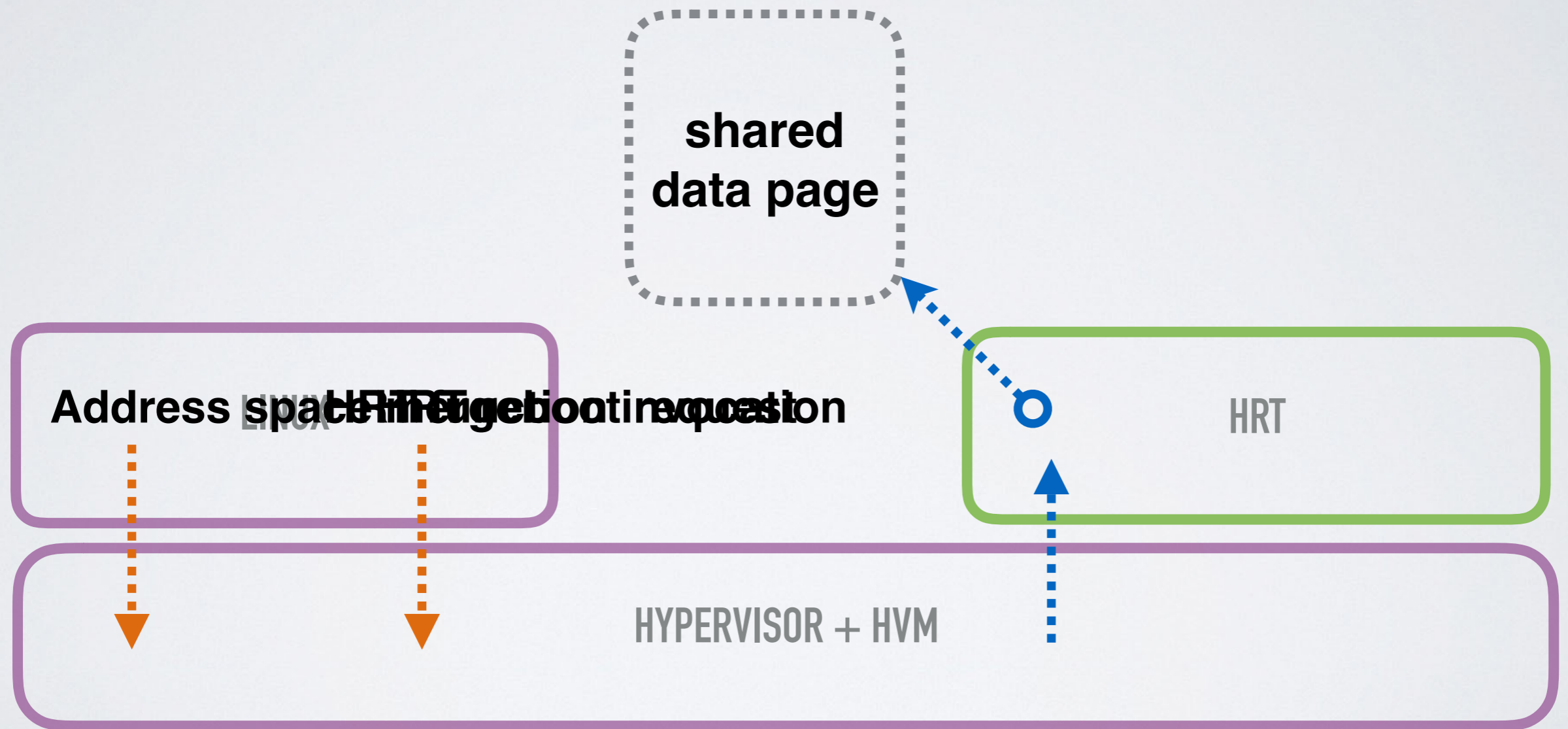
merged address space



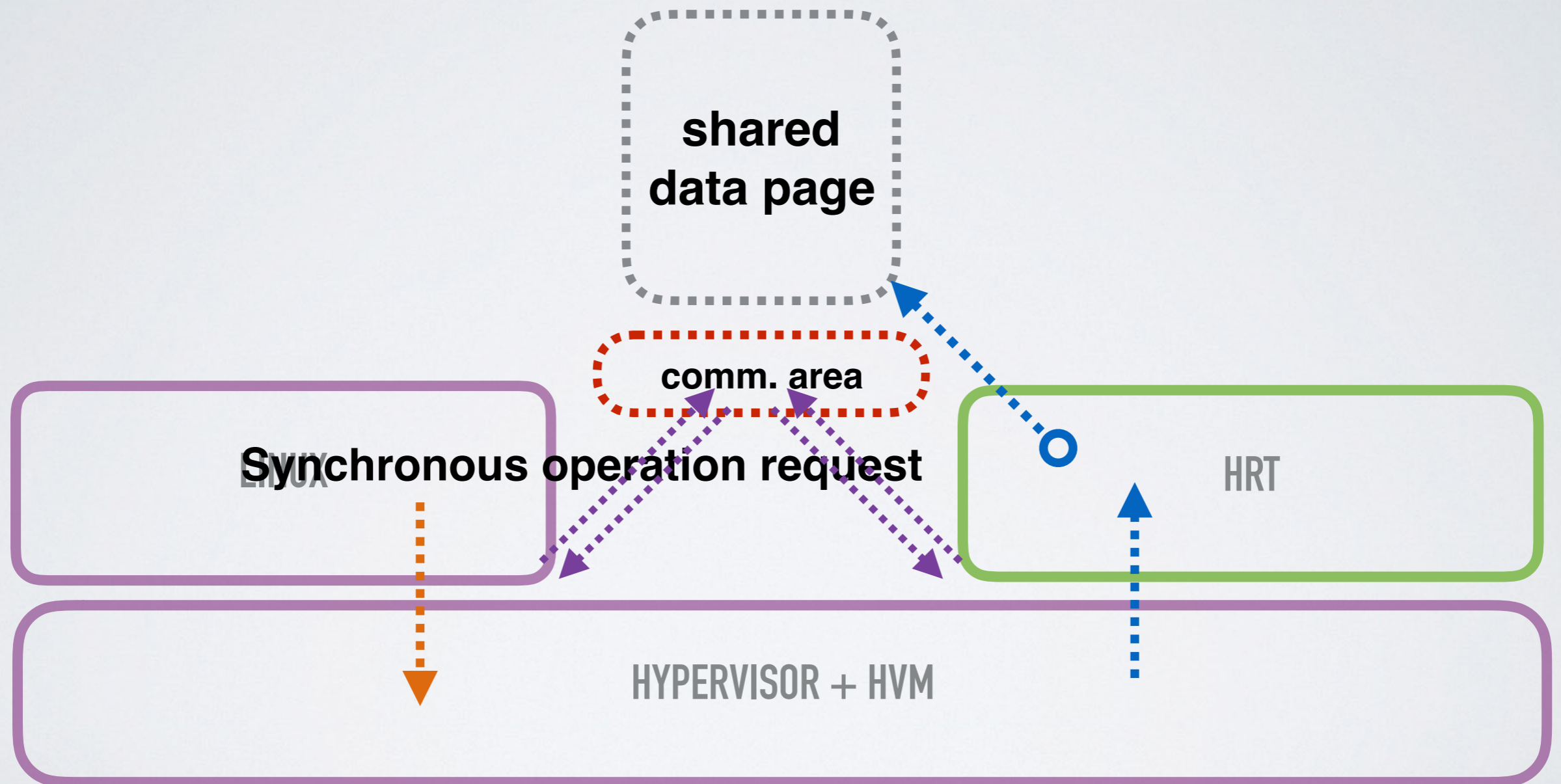
ROS vaddr space

HRT vaddr space

ROS/HRT communication ⁴⁹



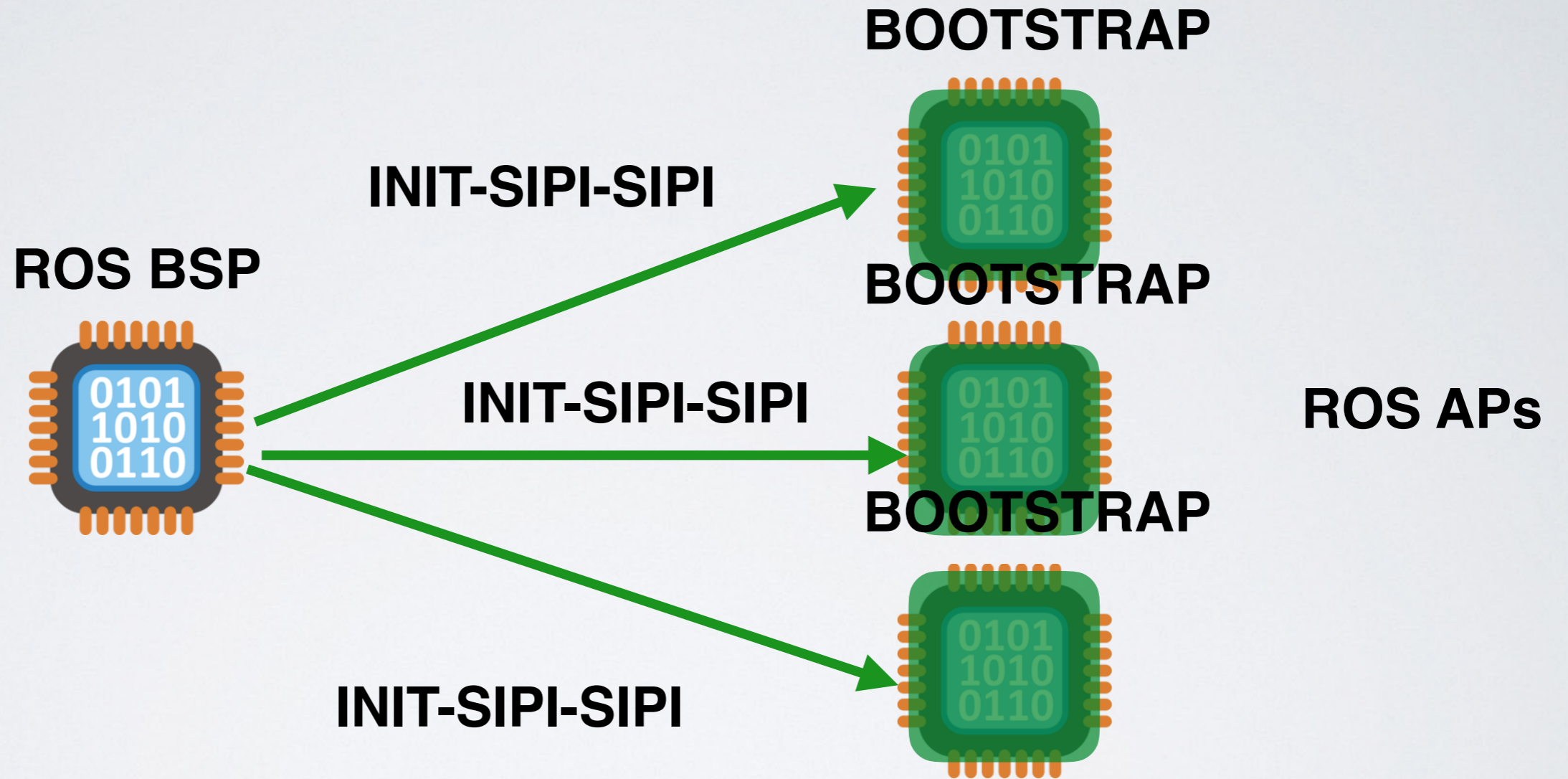
ROS/HRT communication ⁵⁰



ROS/HRT communication⁵¹

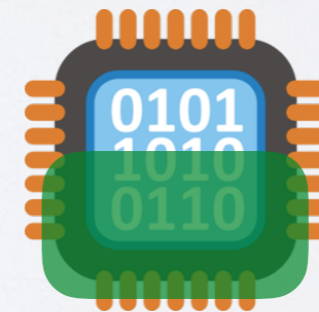
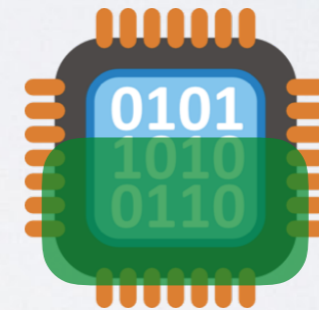
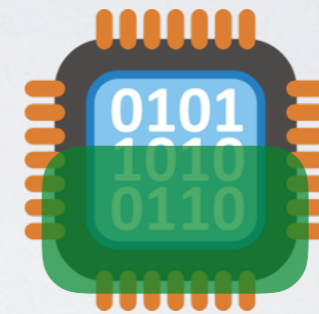
	Cycles	Time
Addr space merge	~33K	15 μ s
Asynch function invocation	~25K	11 μ s
Synch function invocation (remote socket)	~1060	482ns
Synch function invocation (same socket)	~790	359ns

ROS boot



HRT boot

**HVM set up
registers, in
systematically**



GDT

IDT

TSS

REGISTERS

INIT
PAGE
TABLES

***no need for INIT-SIPI-SIPI sequence**

LINUX FORK + EXEC ~ 714 μ s

HVM + HRT CORE BOOT ~ 61 μ s

*how do we take a
legacy runtime to
the HRT + X model?*

PORT

***porting a runtime/app
to an a new OS
environment is...***

DIFFICULT

TIME-CONSUMING

ERROR-PRONE

development cycle:

do {

ADD FUNCTION

REBUILD

BOOT HRT

} while (HRT falls over)

*much of the
functionality is*

NOT ON THE CRITICAL PATH

***we want to make this
easier***

OUTLINE

Background/Overview

Nautilus

Deployment Models

Hybrid Virtual Machine

○ Multiverse & Future Work

give us your legacy
runtime

*we automatically
transform it to run*

as an HRT

in kernel mode

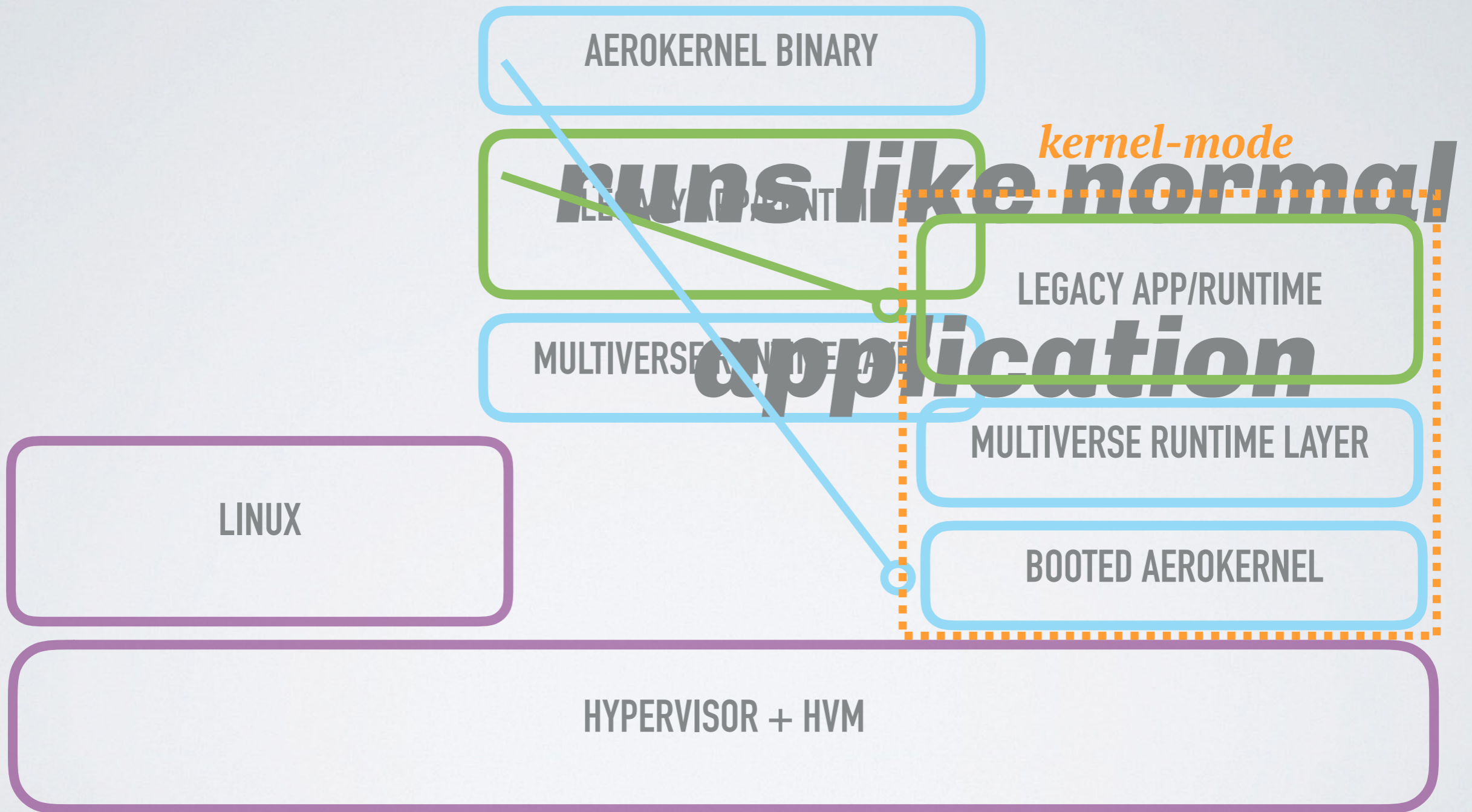
bridged with a legacy OS

(AUTOMATIC HYBRIDIZATION)

LEGACY APP/RUNTIME



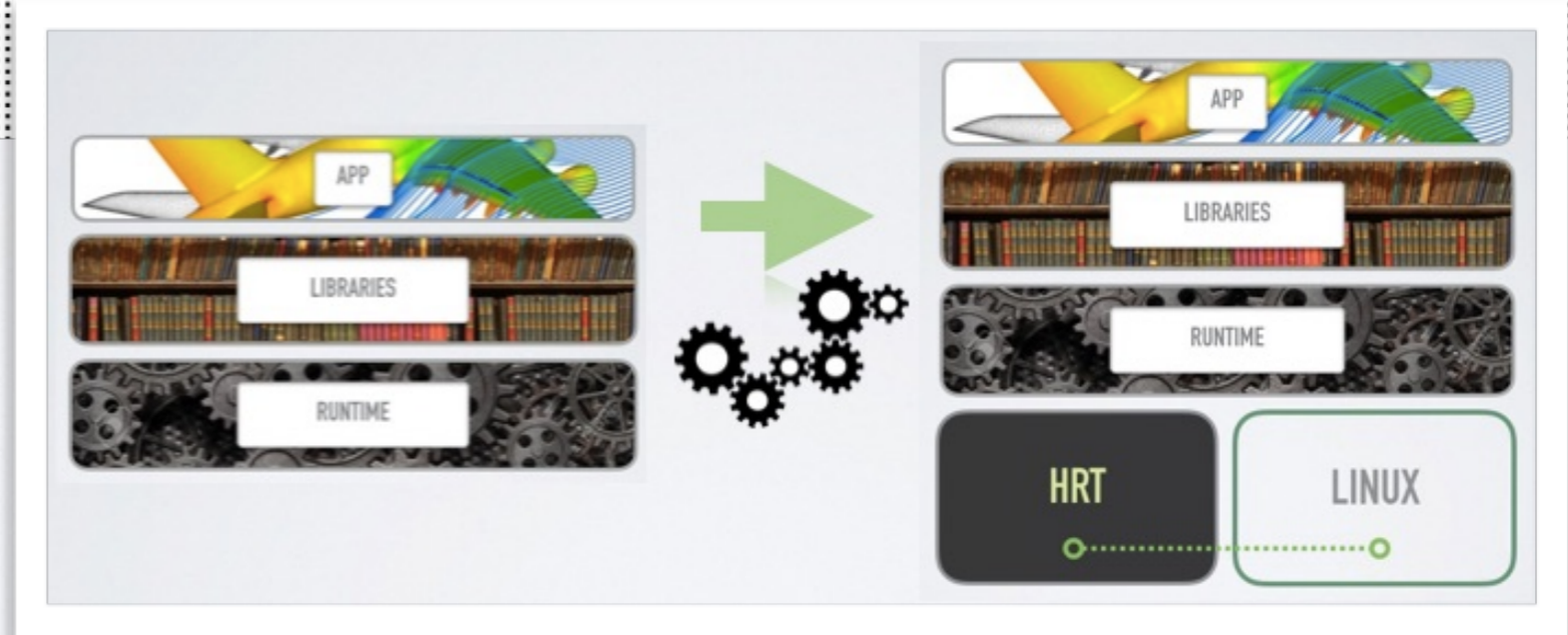
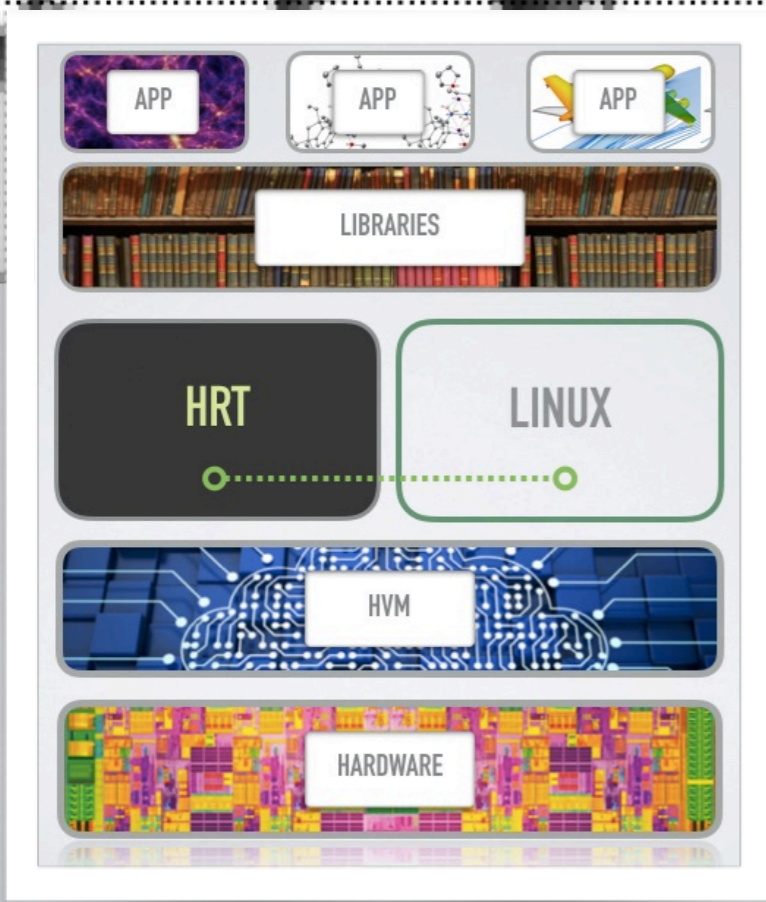
rebuild with our toolchain




```
#  
# ls  
bench-write.out      go      mracket-GOLD  
binary-tree-2.rkt   intsum-native  multiverse-racket  
bytes               ism     multiverse.log  
collects           isn     nbody.rkt  
doall.sh           lgn-hpcg    racket  
doruns.sh          lgo     results  
fannkuch-redux.rkt lost+found  spectral-norm.rkt  
fasta-3.rkt        lpm     test.out  
fasta.rkt          lpn     test.t  
g                 mandelbrot-2.rkt  
# █
```

4 parallel runtimes (Legion, Racket, NDPC, NESL)

nautilus



Multiverse

automatic transformation: legacy app+runtime → HRT

Hybrid Virtual Machine
bridge HRT with legacy OS

thank you



nautilus

my webpage: halek.co

lab: presciencelab.org

download nautilus: nautilus.halek.co

v3vee project: v3vee.org

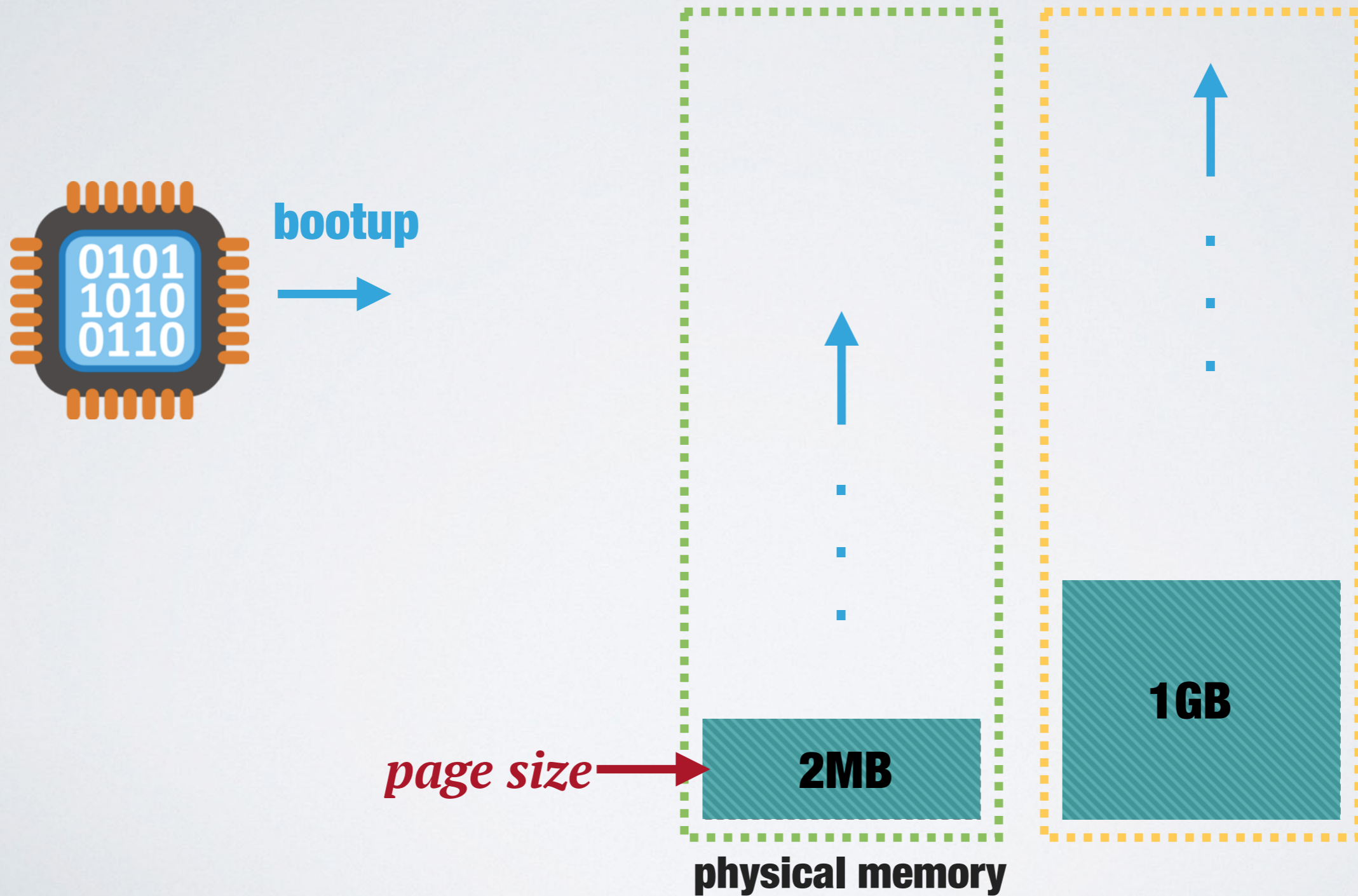
backups

thread fork

interrupt driven execution

memory and paging

static identity map



eliminates expensive page faults

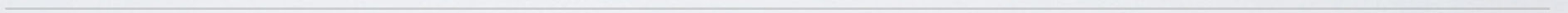
reduces TLB misses + shutdowns

*increases performance under
virtualization*

log scale

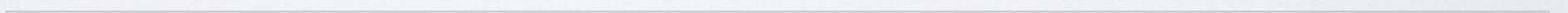


1000000



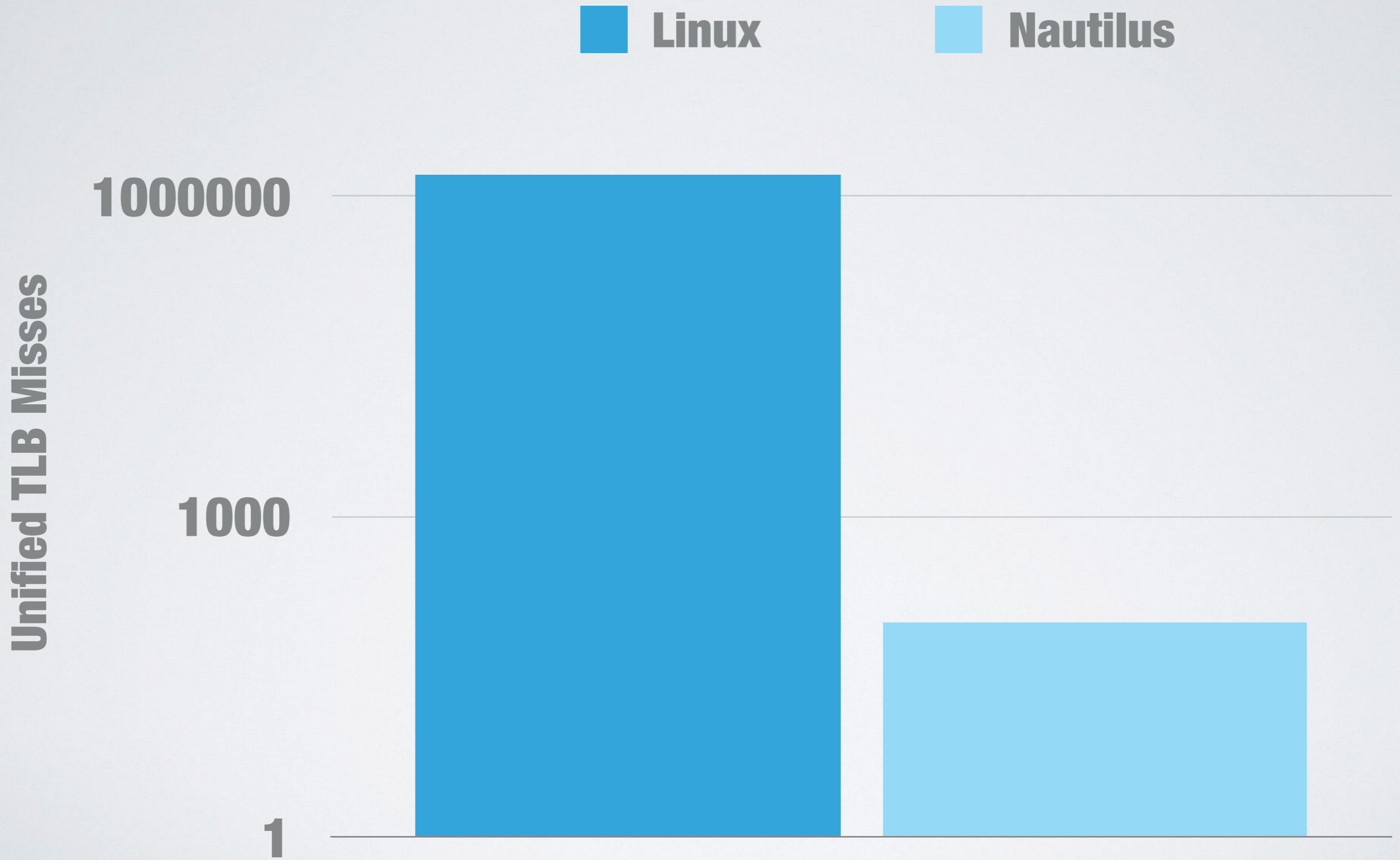
Unified TLB Misses

1000

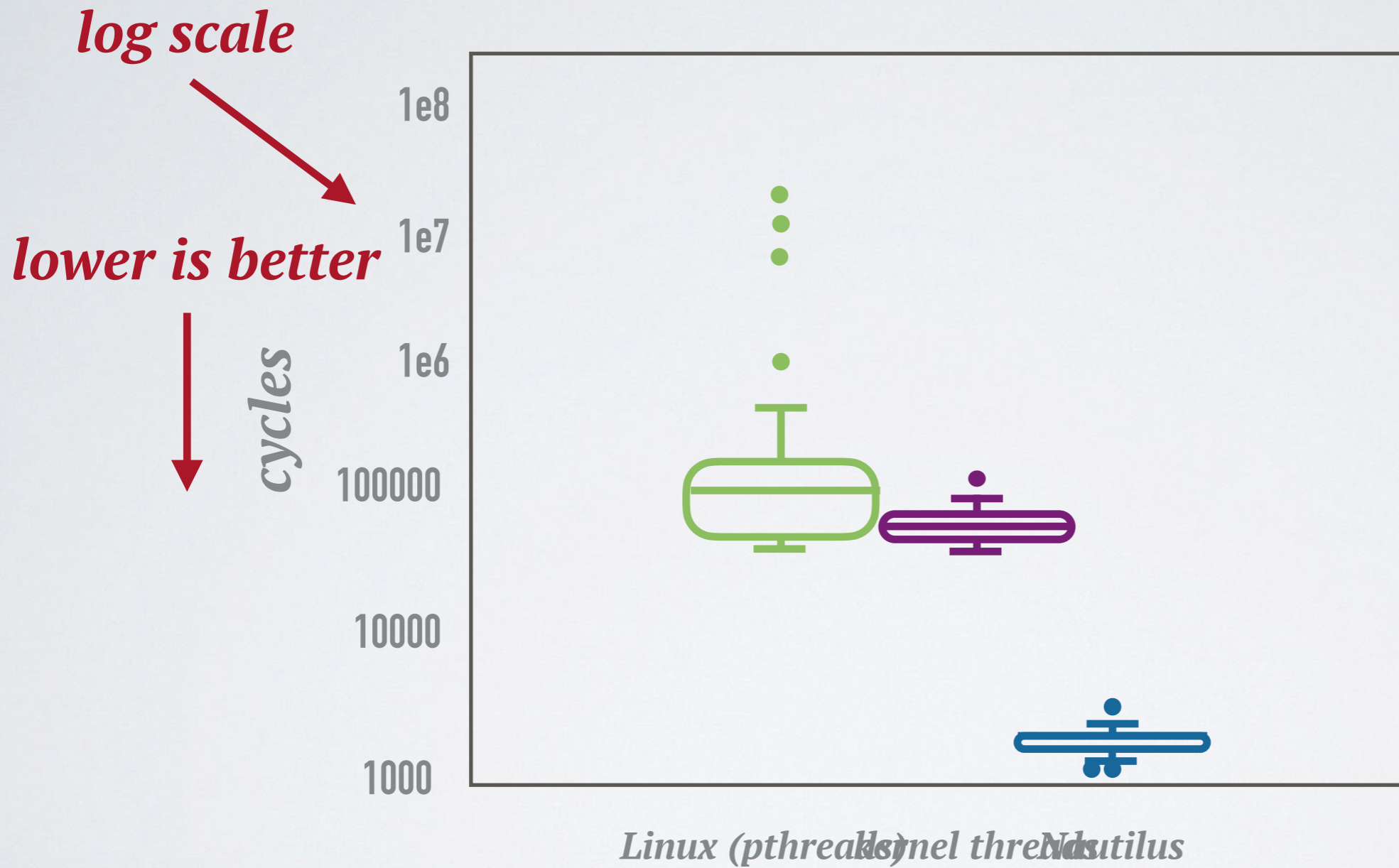


1

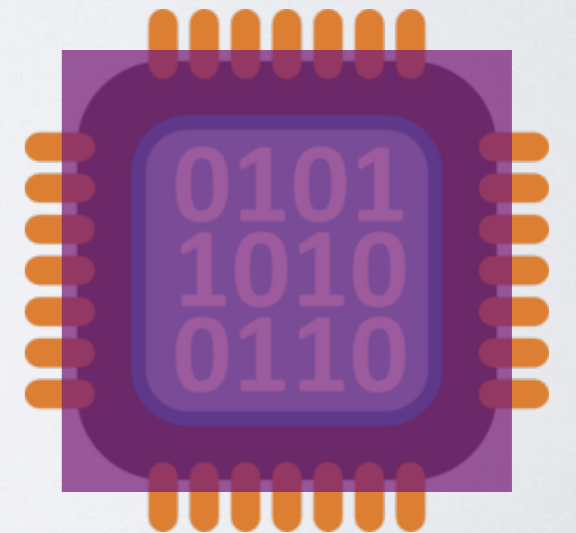
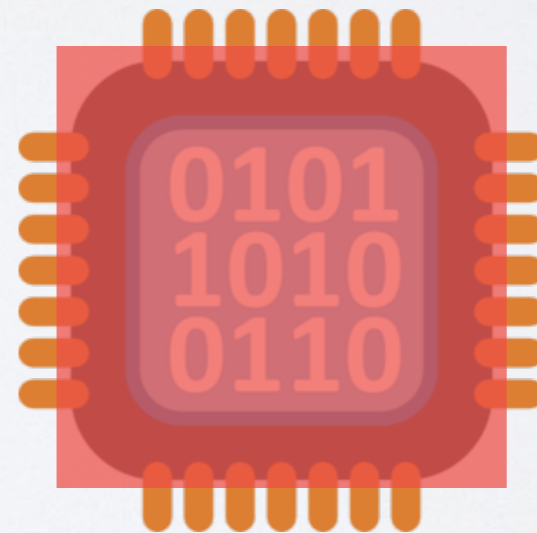
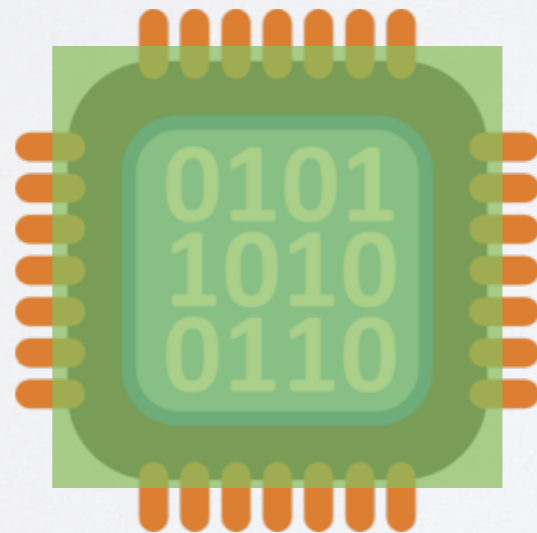
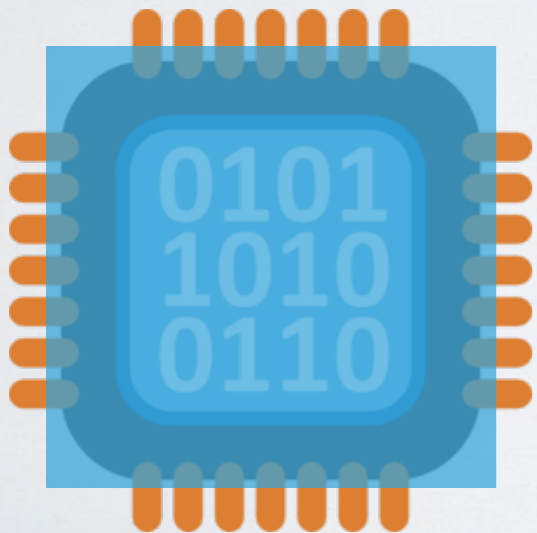




thread creation is **PREDICTABLE**



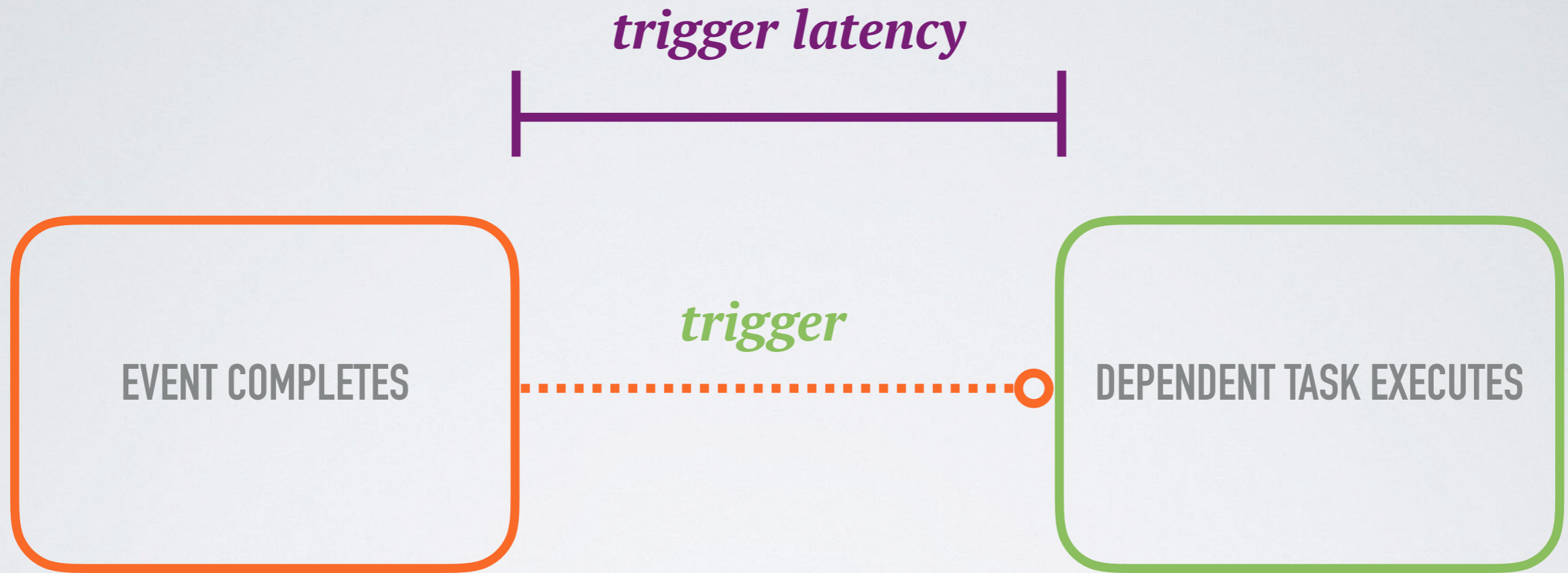
*binding to physical
processor is*
GUARANTEED



RUNTIME *control over*
physical **CPU**

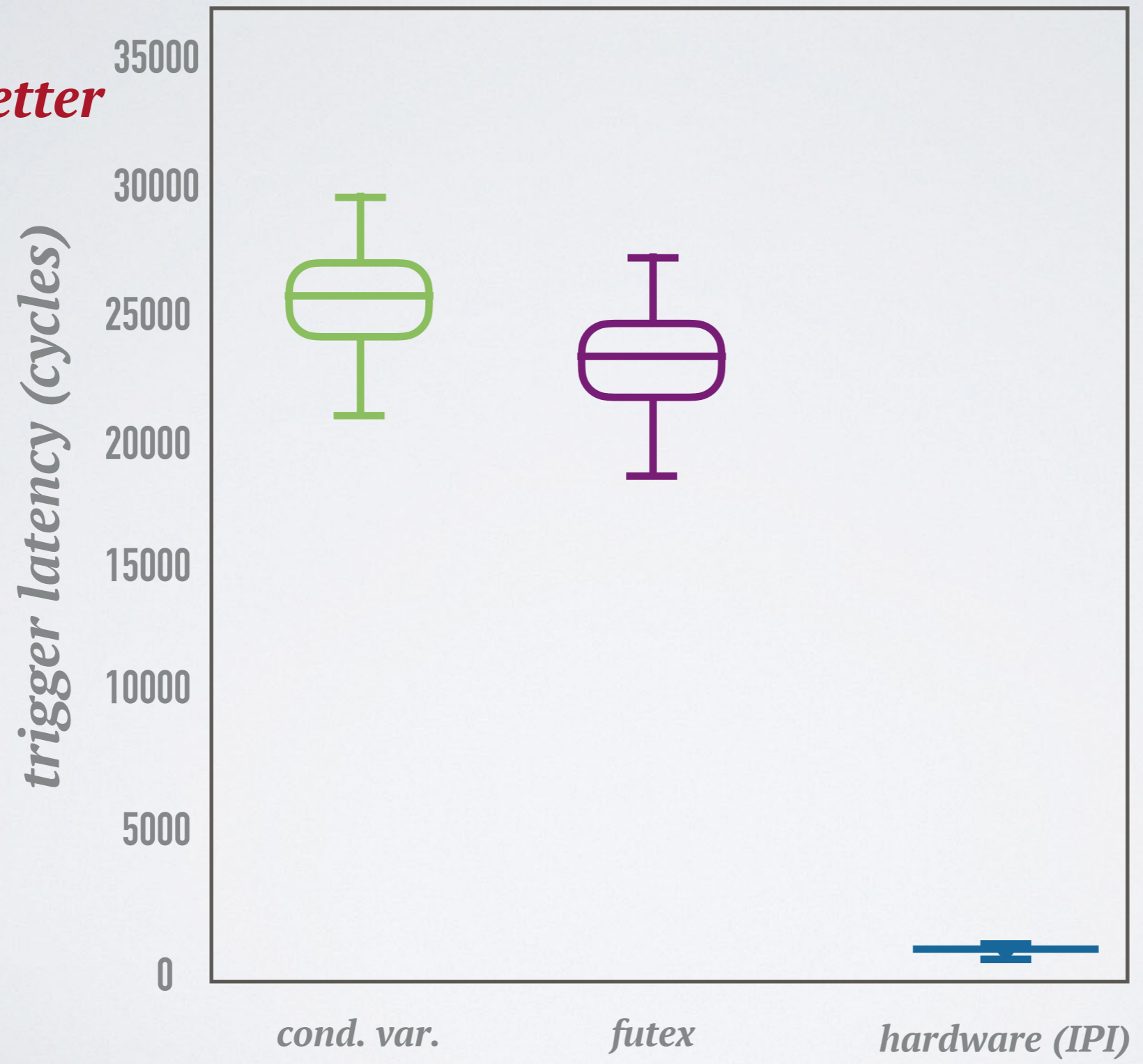
interrupts
&
events

*asynchronous
notifications*

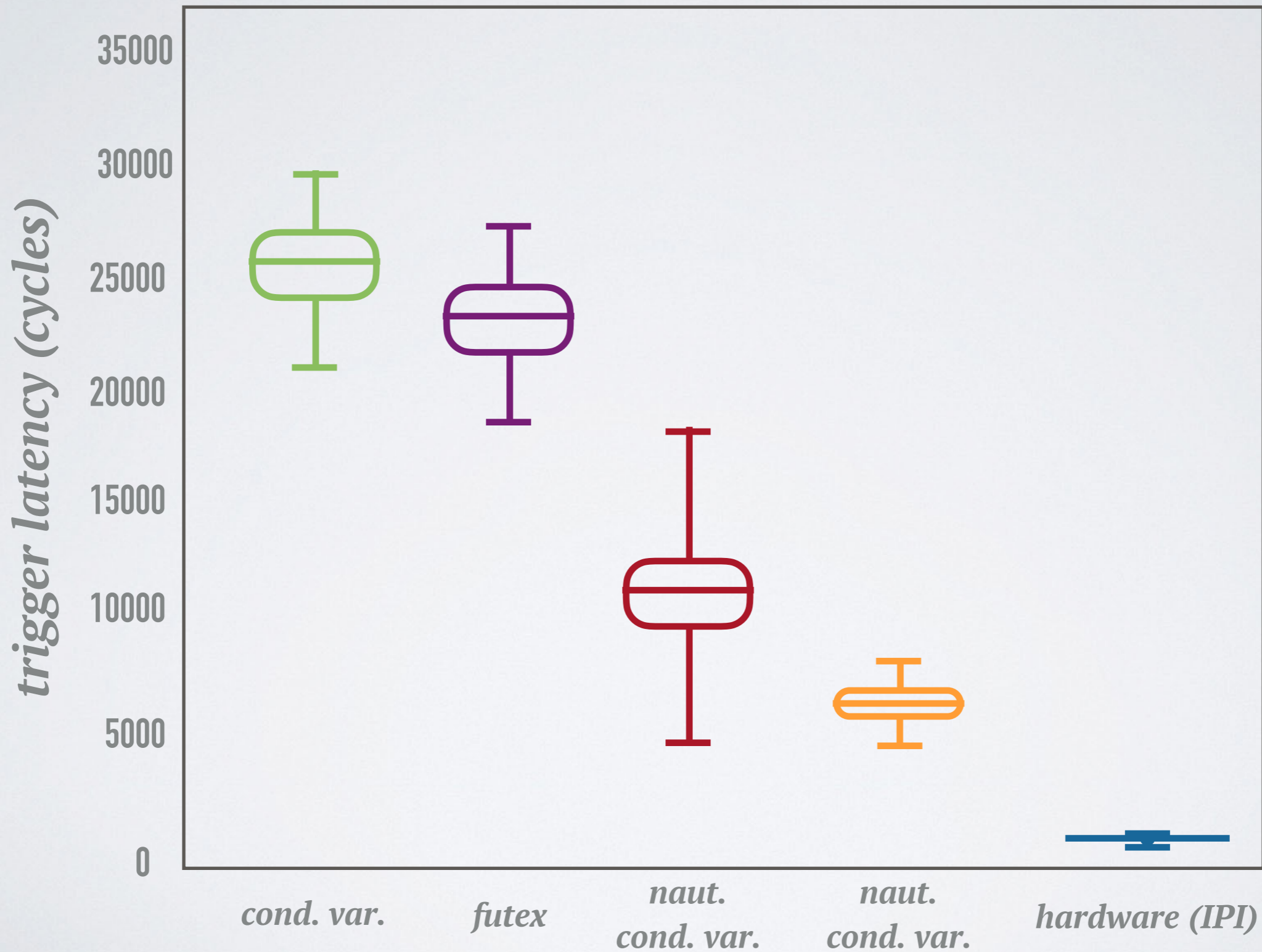


user-mode software events are SLOW

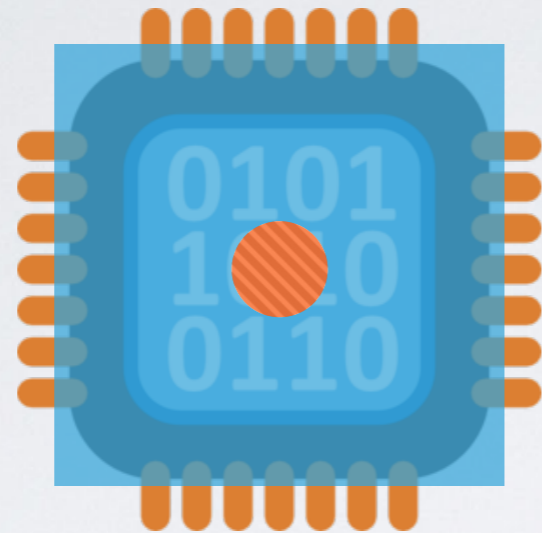
lower is better



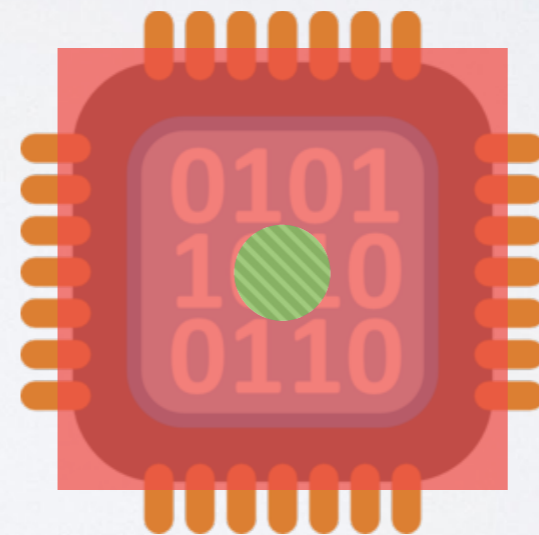
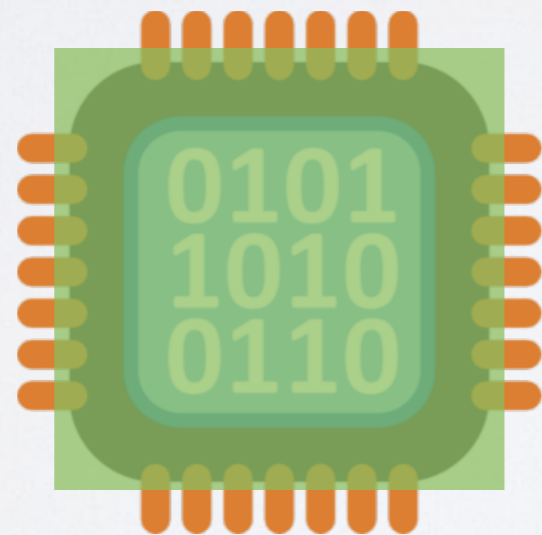
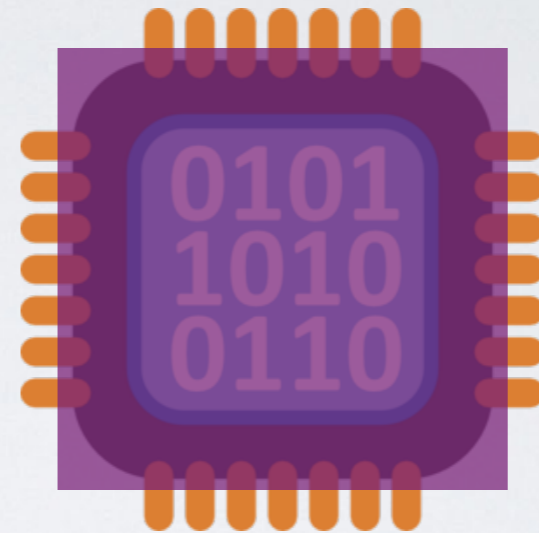
*nautilus events triggers are **FASTER***

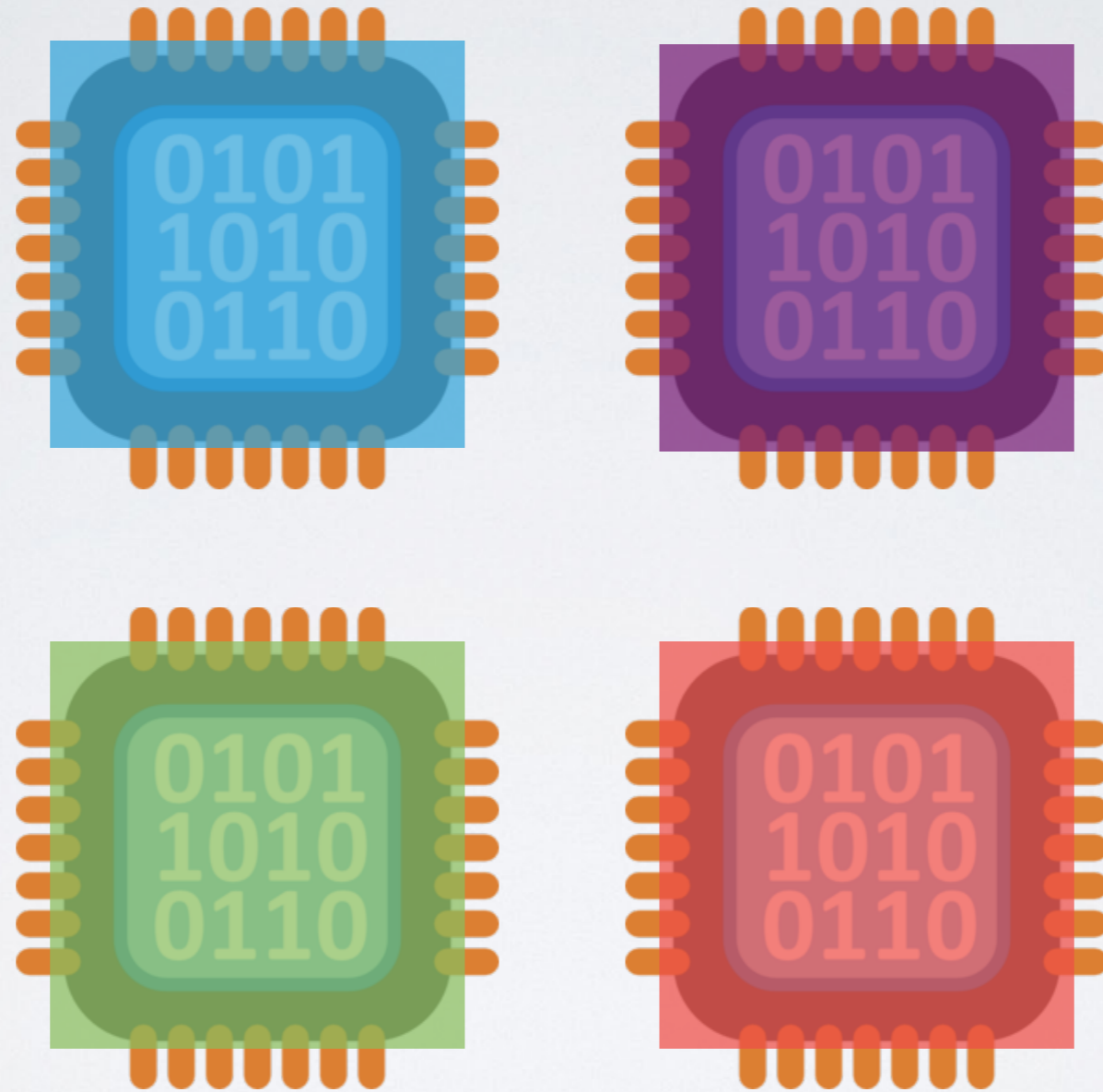


handle_event()

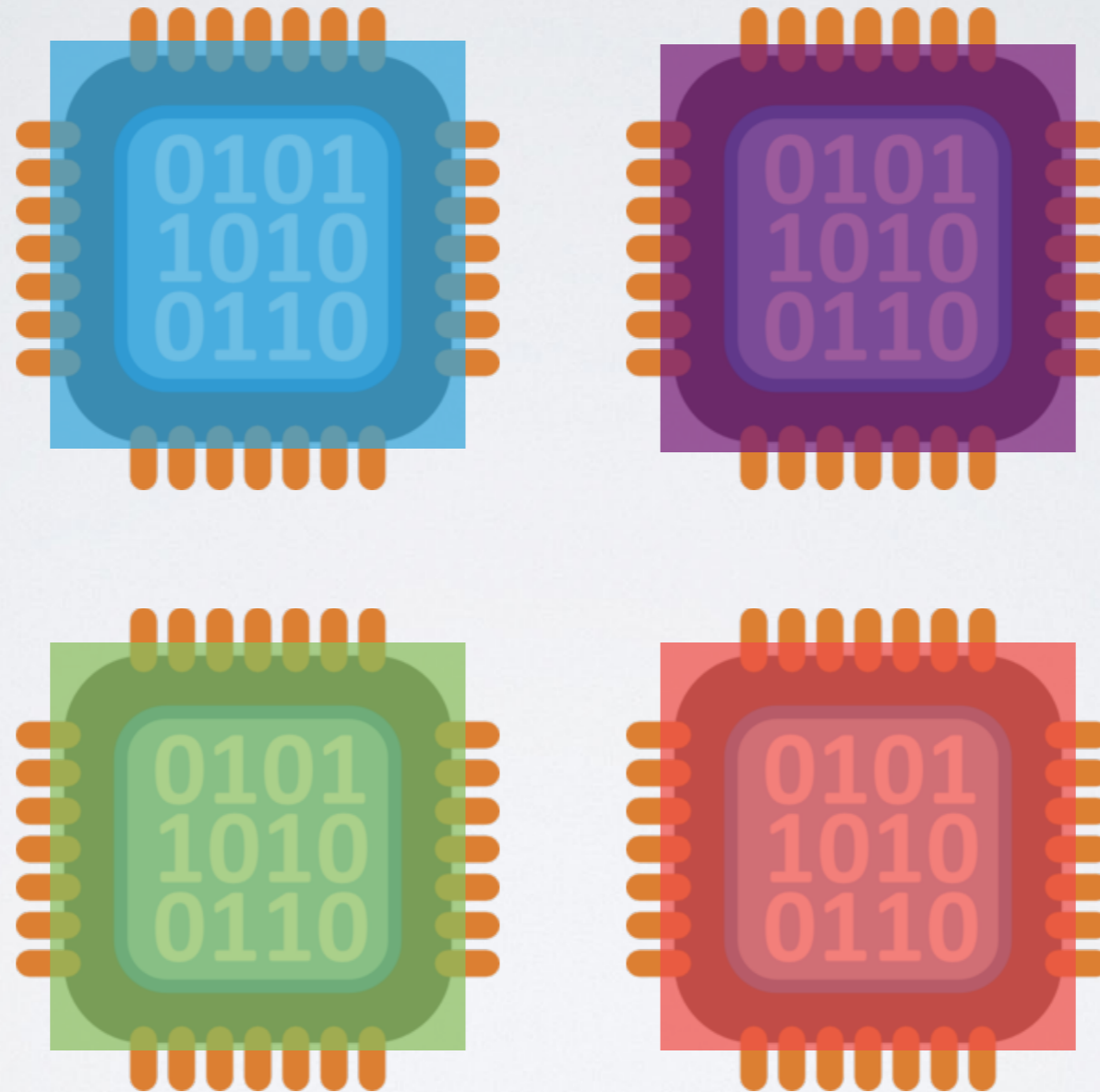


handle_event()





*asynchronous, IPI-based
execution*



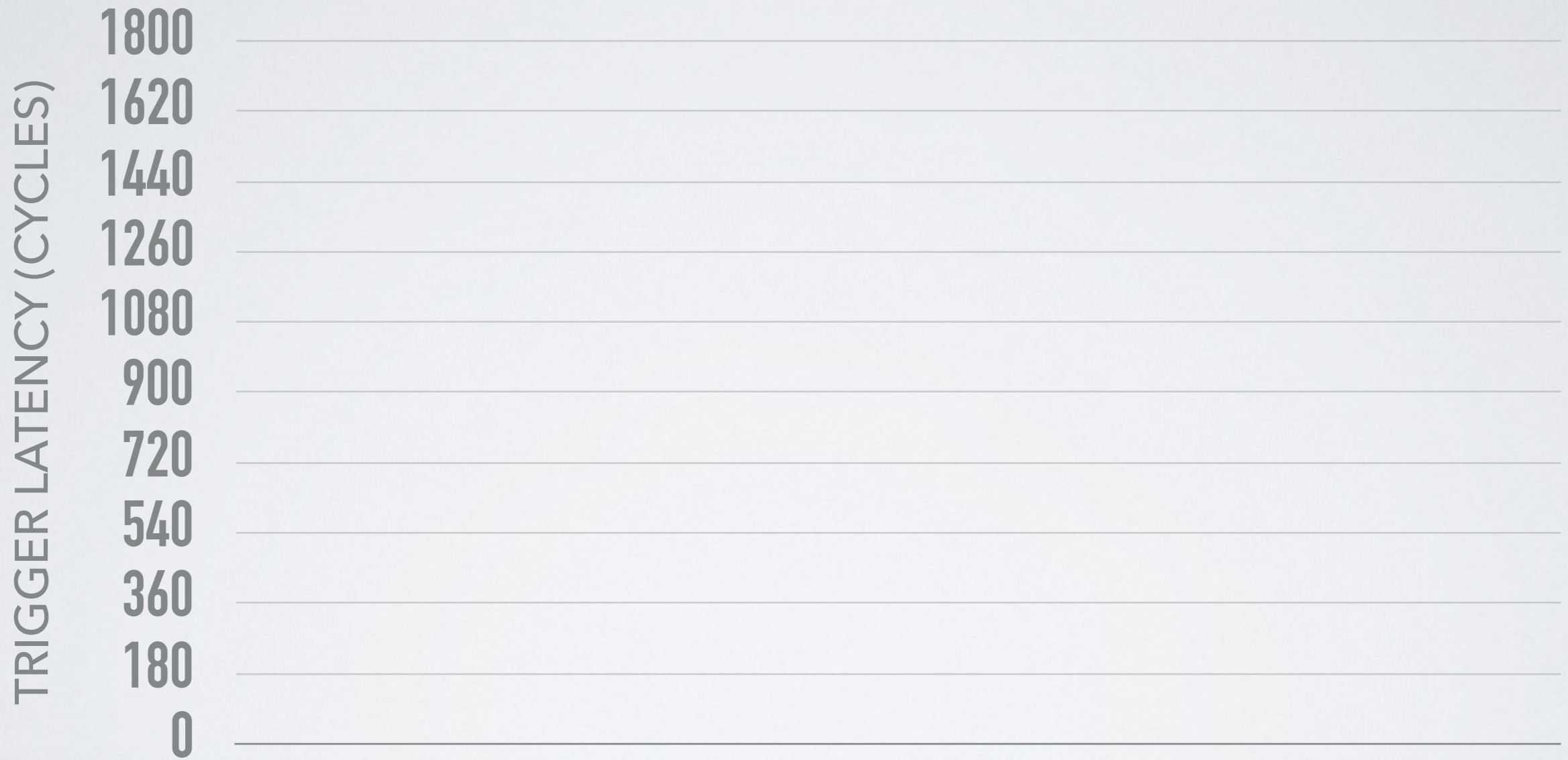
**NOT POSSIBLE IN LINUX
USERSPACE**

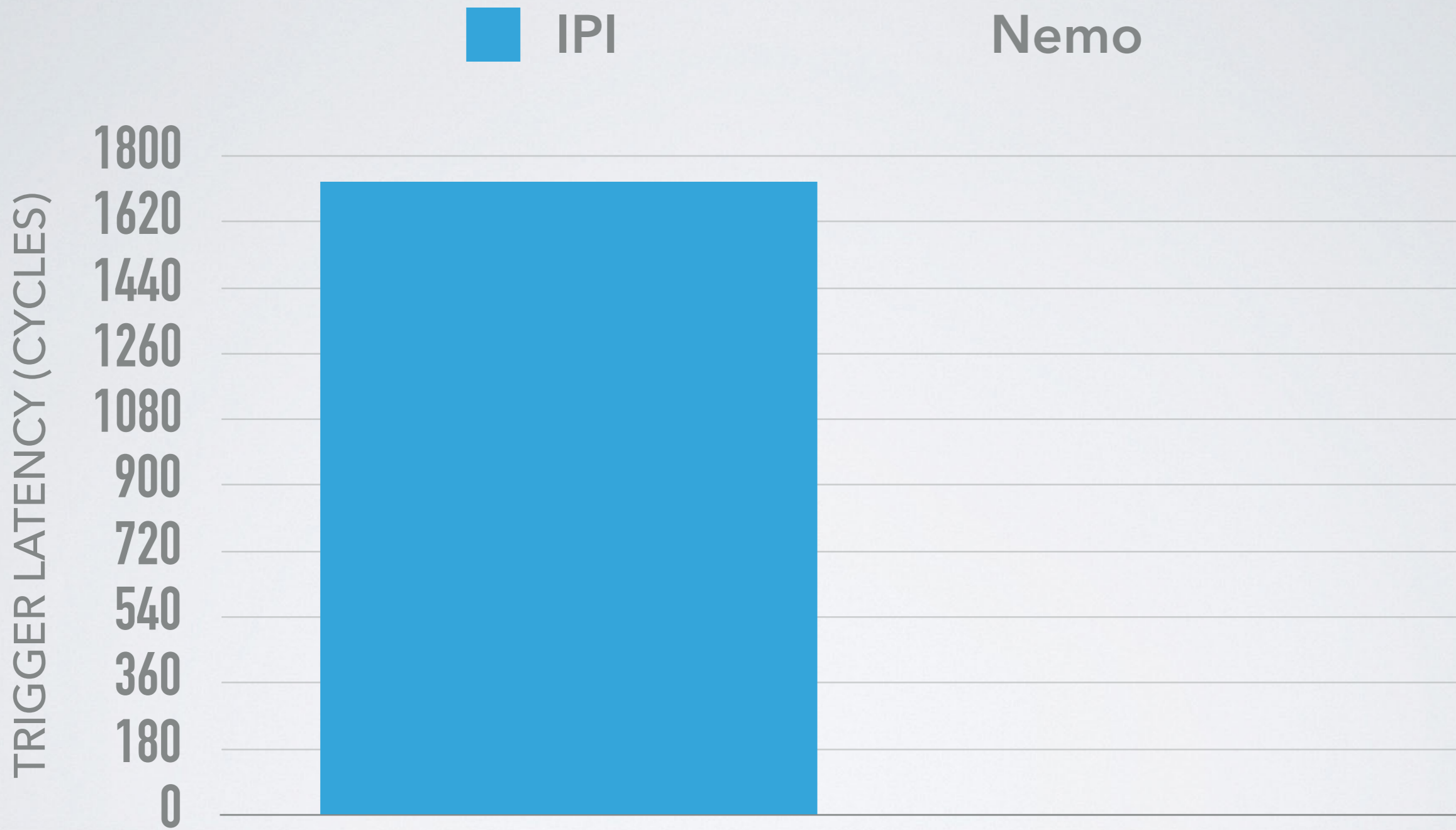
lower is better

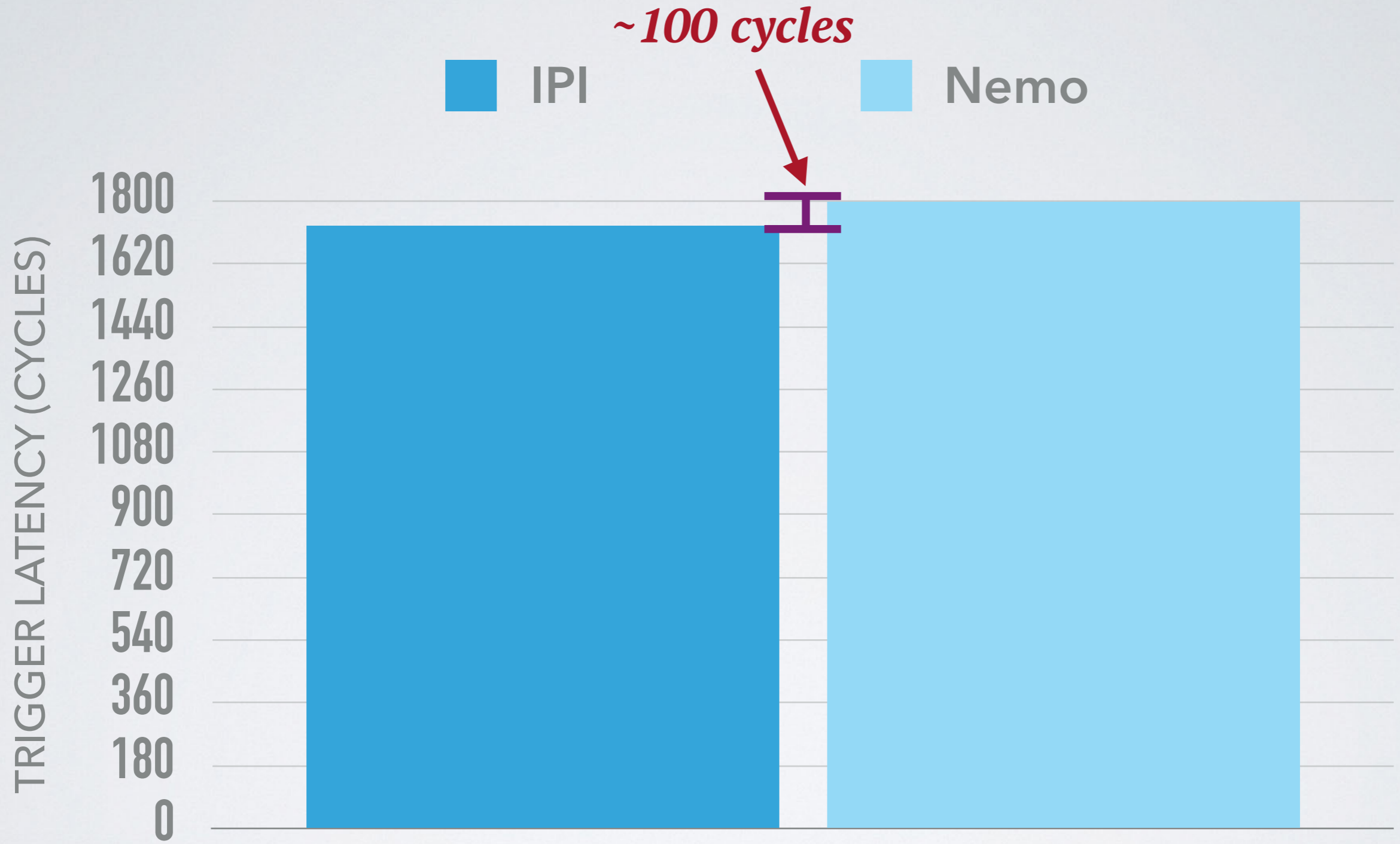


IPI

Nemo







this gives us an
incremental path
for creating HRTs

***start out with a working
HRT system***

***then pull functionality
into the HRT for hot spots***

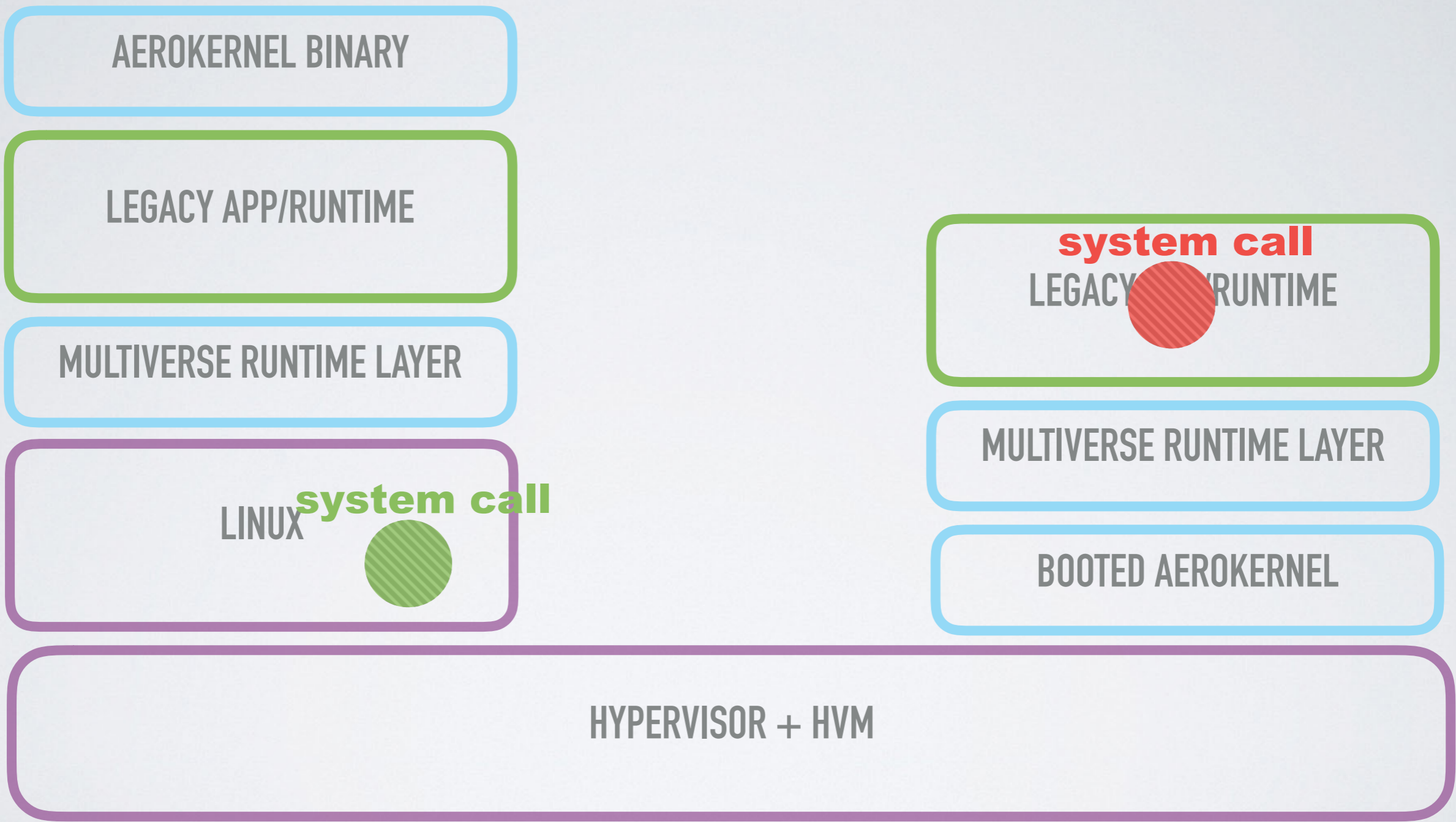
RACKET

***most widely used Scheme
implementation***

downloaded ~300 times/day

complex runtime with JIT

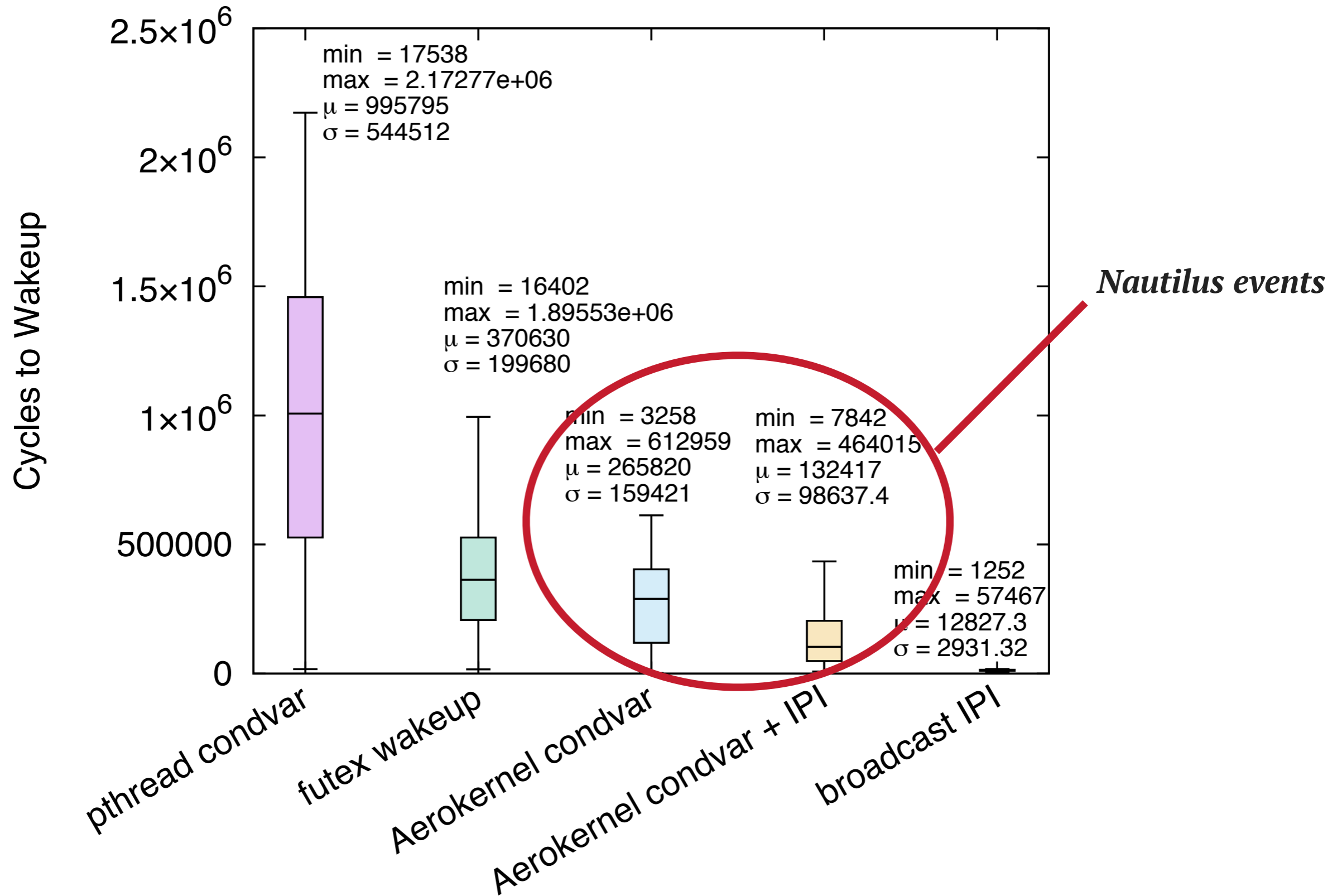
ROS/HRT communication



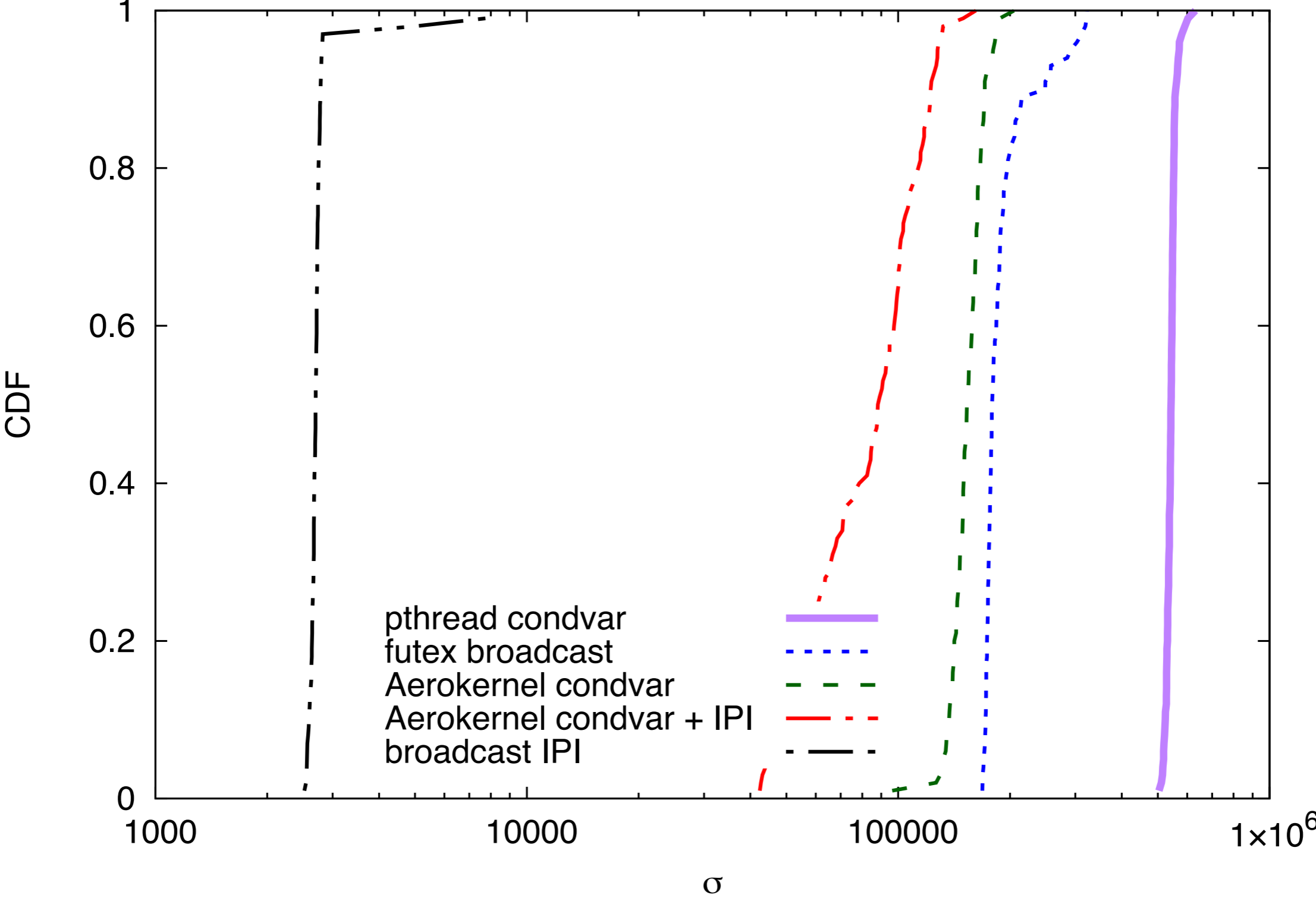
```
# 1S
bench-write.out      go          mracket-GOLD
binary-tree-2.rkt   intsum-native
bytes                ism
collects            isn
doall.sh            lgn-hpcg
doruns.sh           lgo
fannkuch-redux.rkt lost+found
fasta-3.rkt         lpm
fasta.rkt           lpn
g                   mandelbrot-2.rkt
# █
```

}

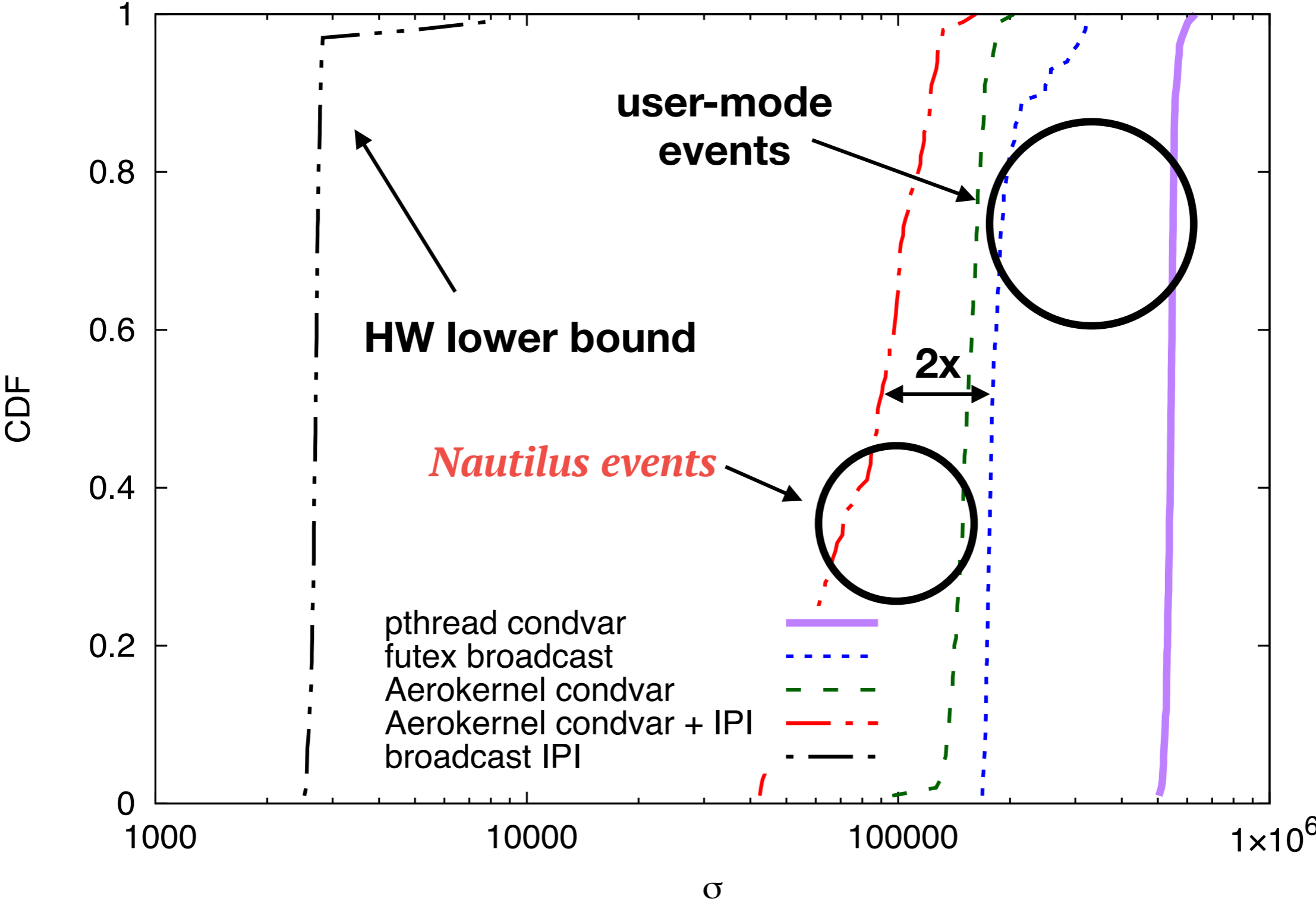
Broadcast wakeups on x64



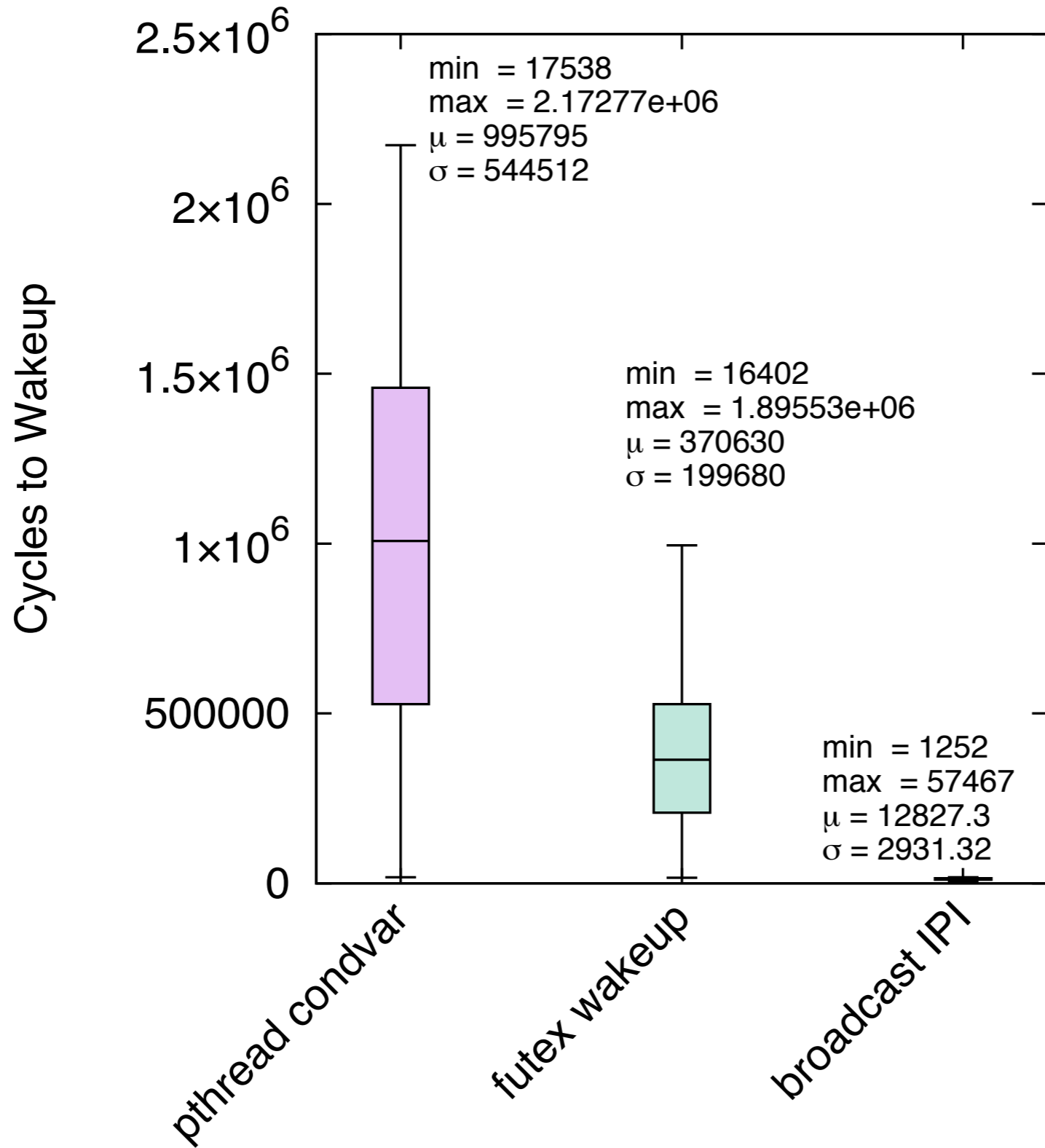
Wakeup deviation on x64



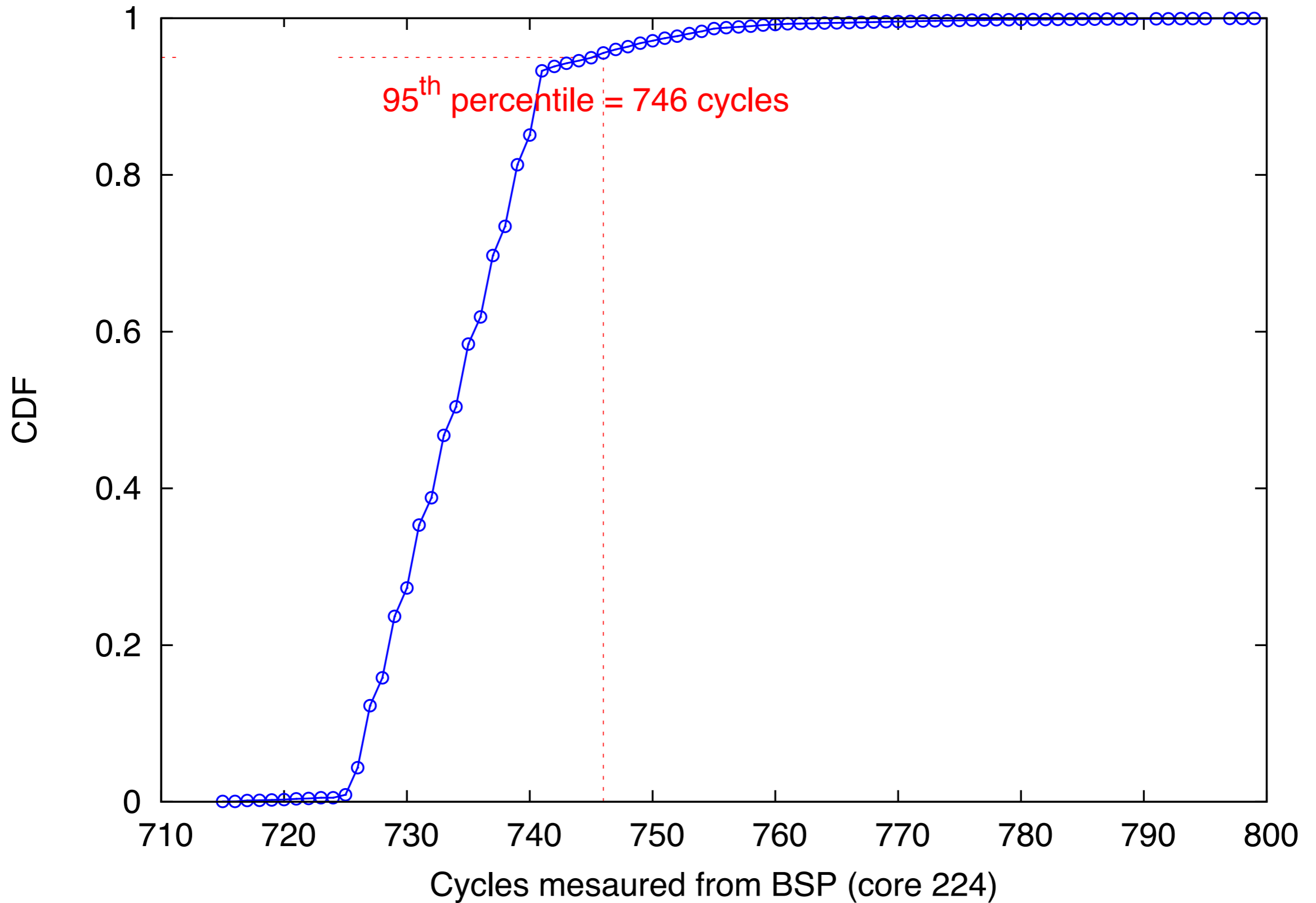
Wakeup deviation on x64



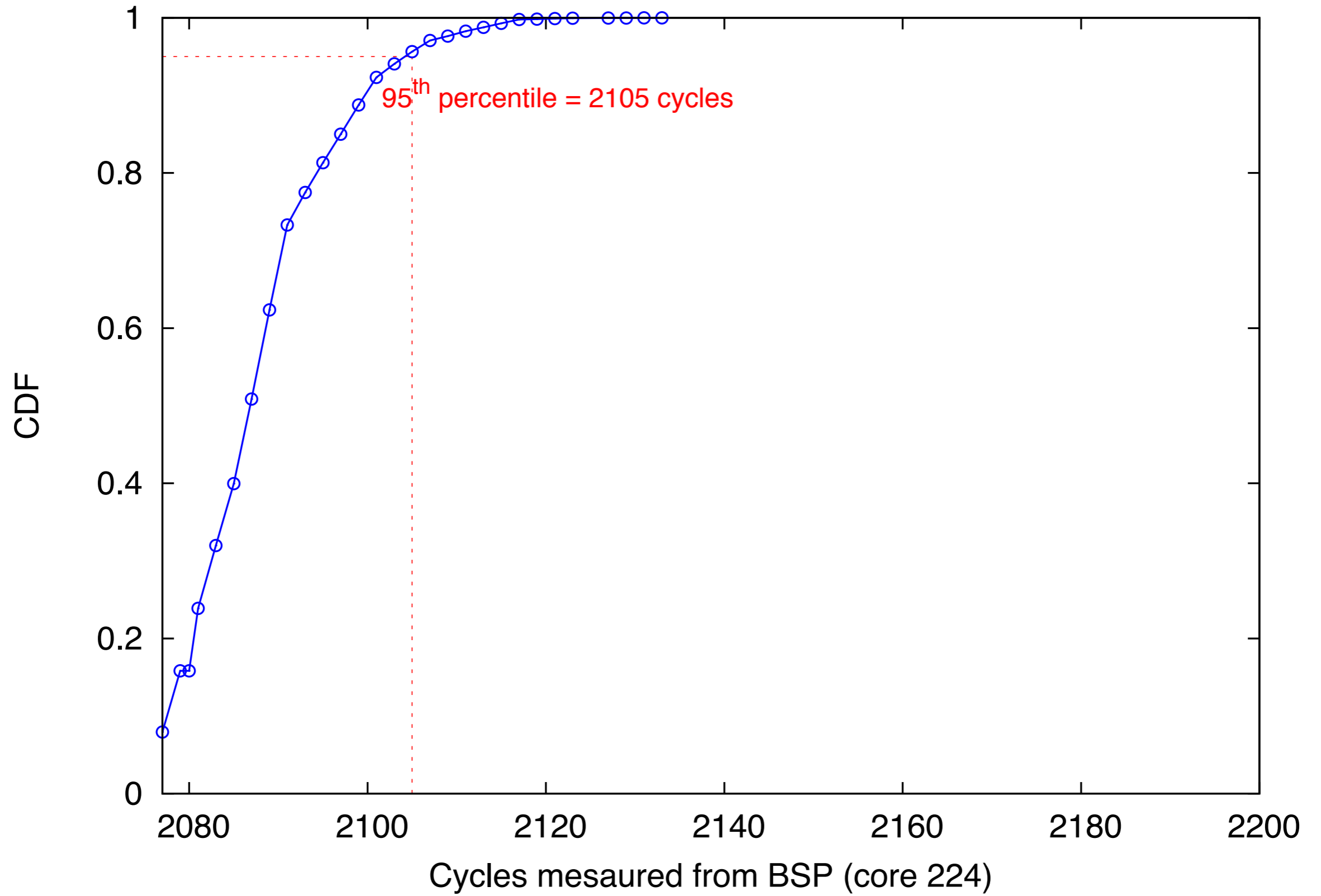
Broadcast wakeups on x64



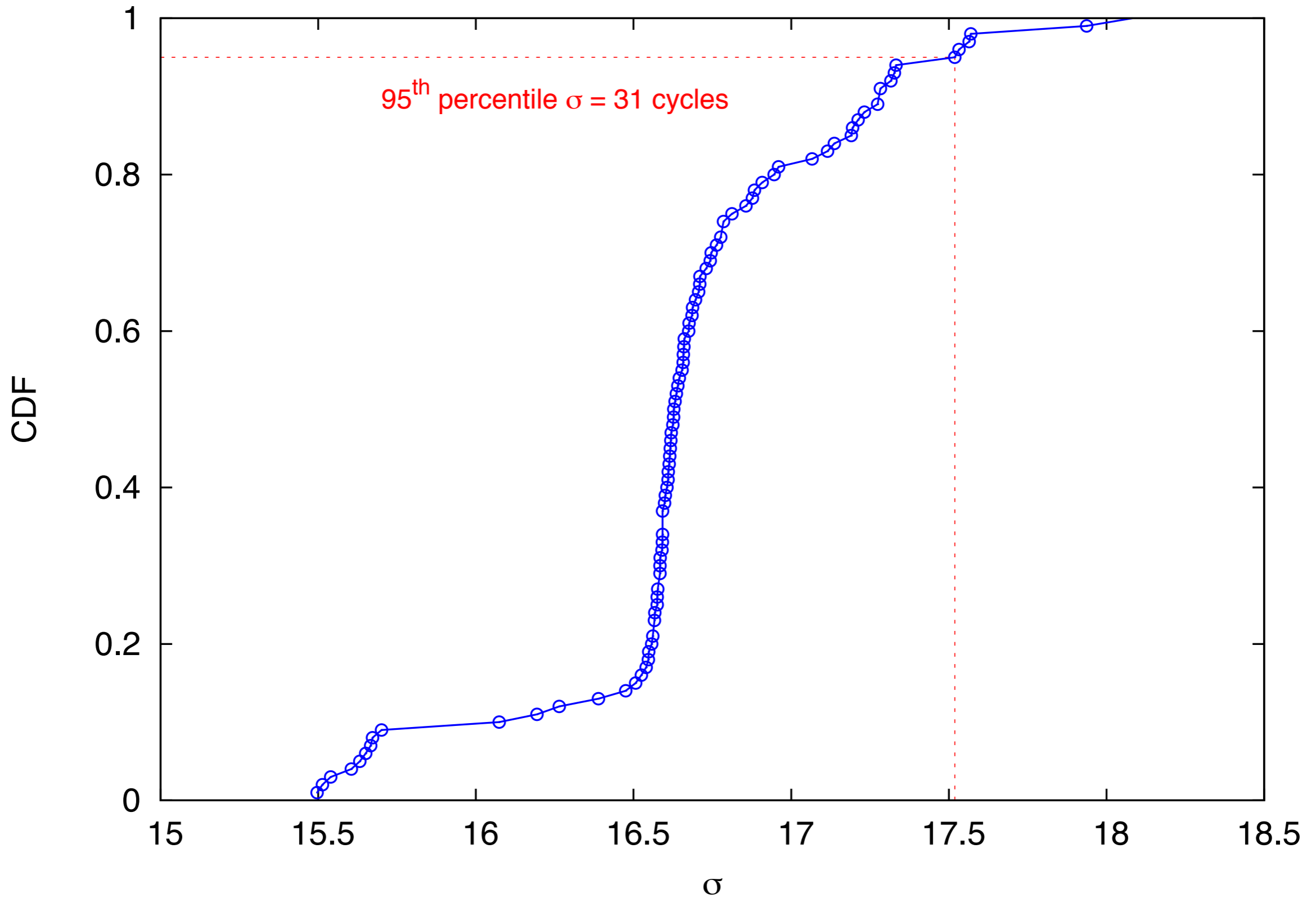
Unicast IPs on phi



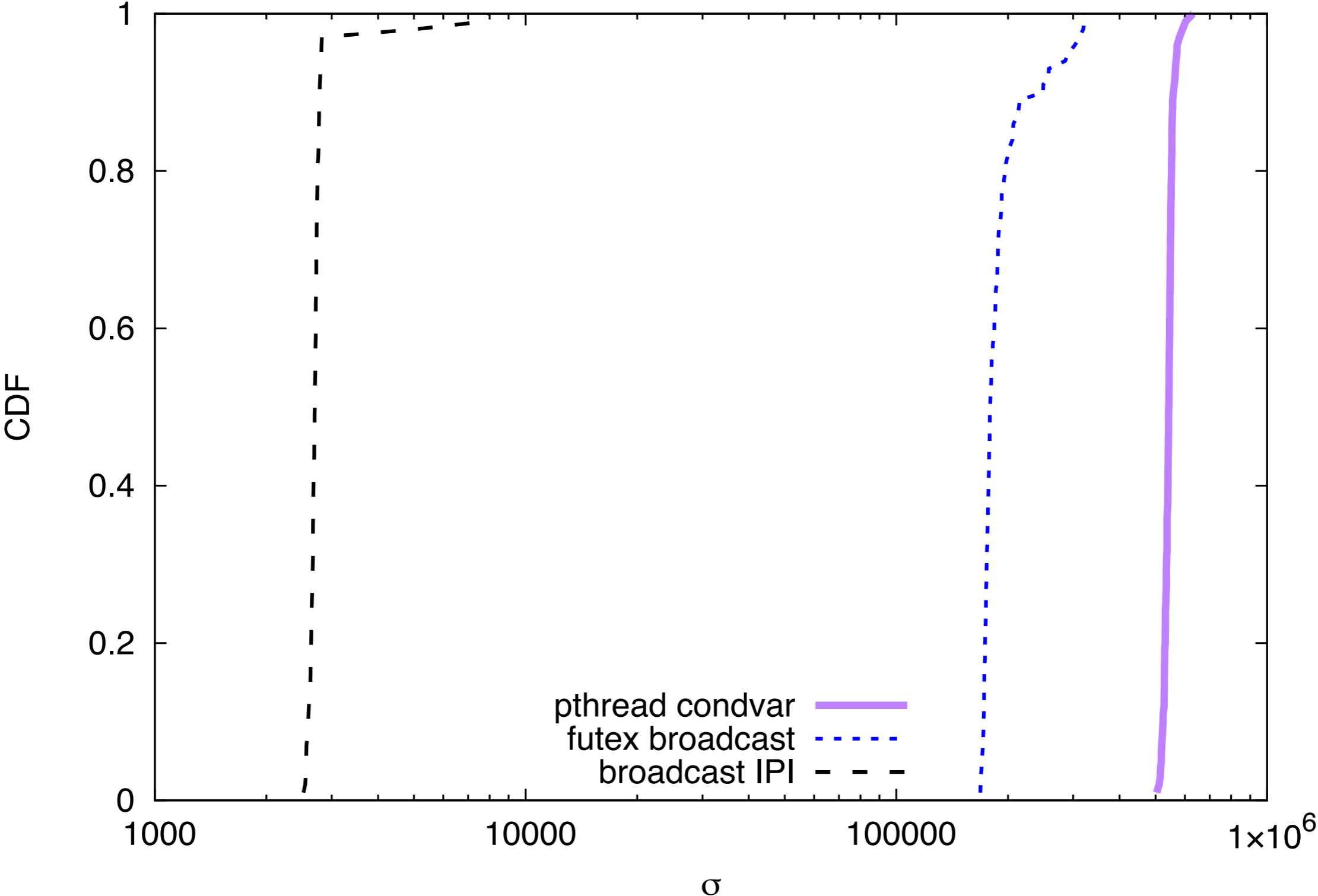
Roundtrip IPs on phi



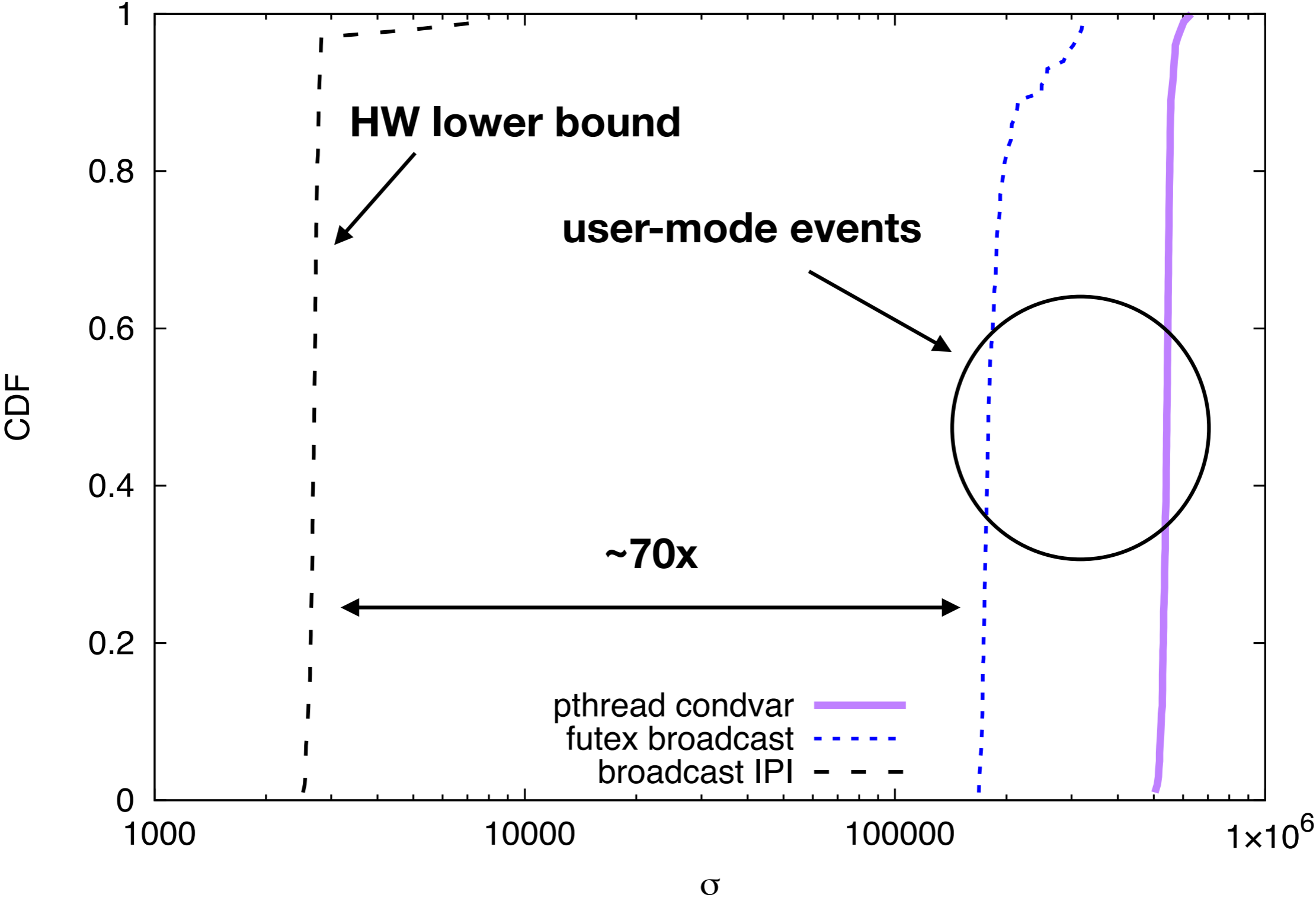
Multicast IPs on phi



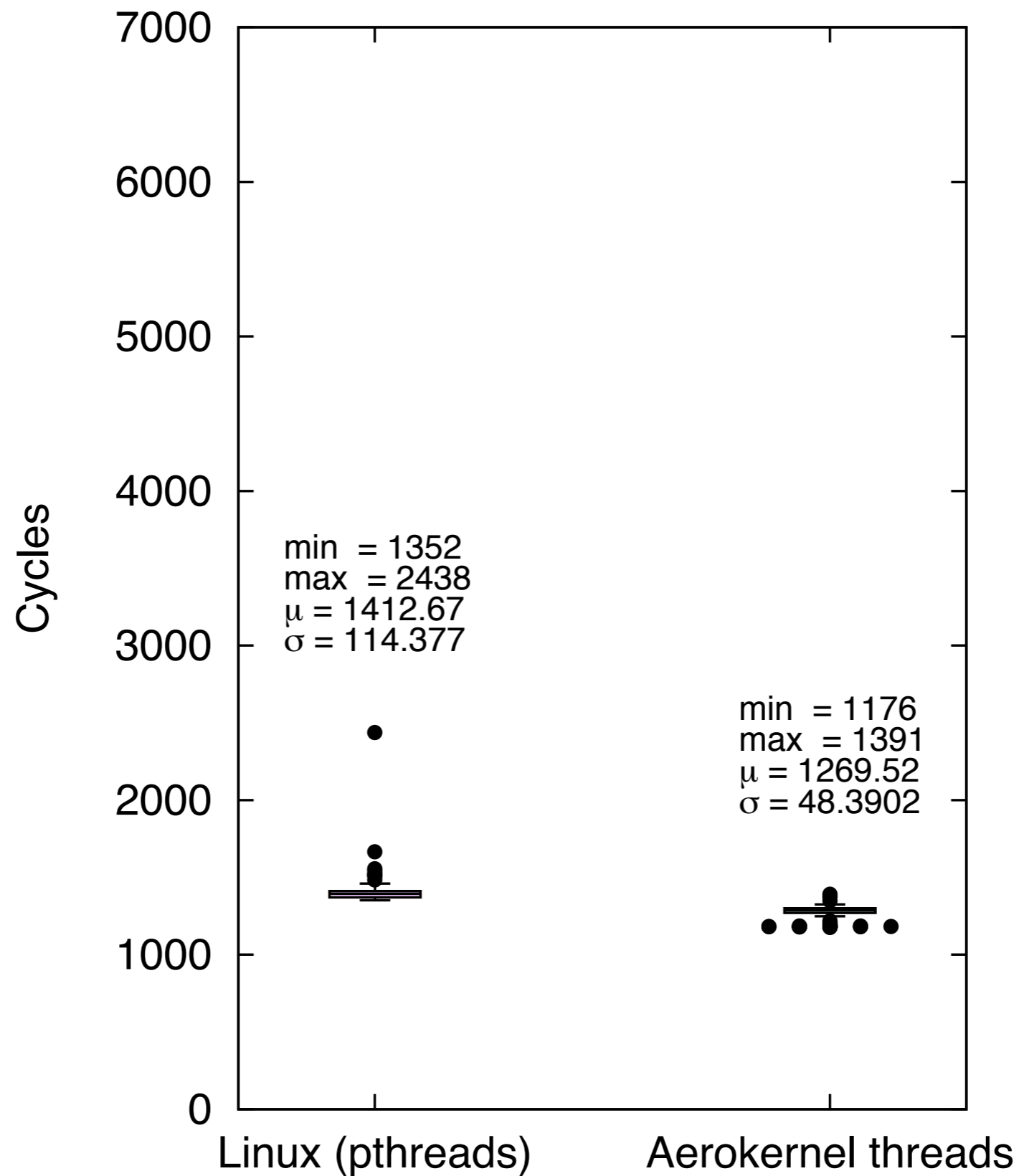
Wakeup deviation on x64



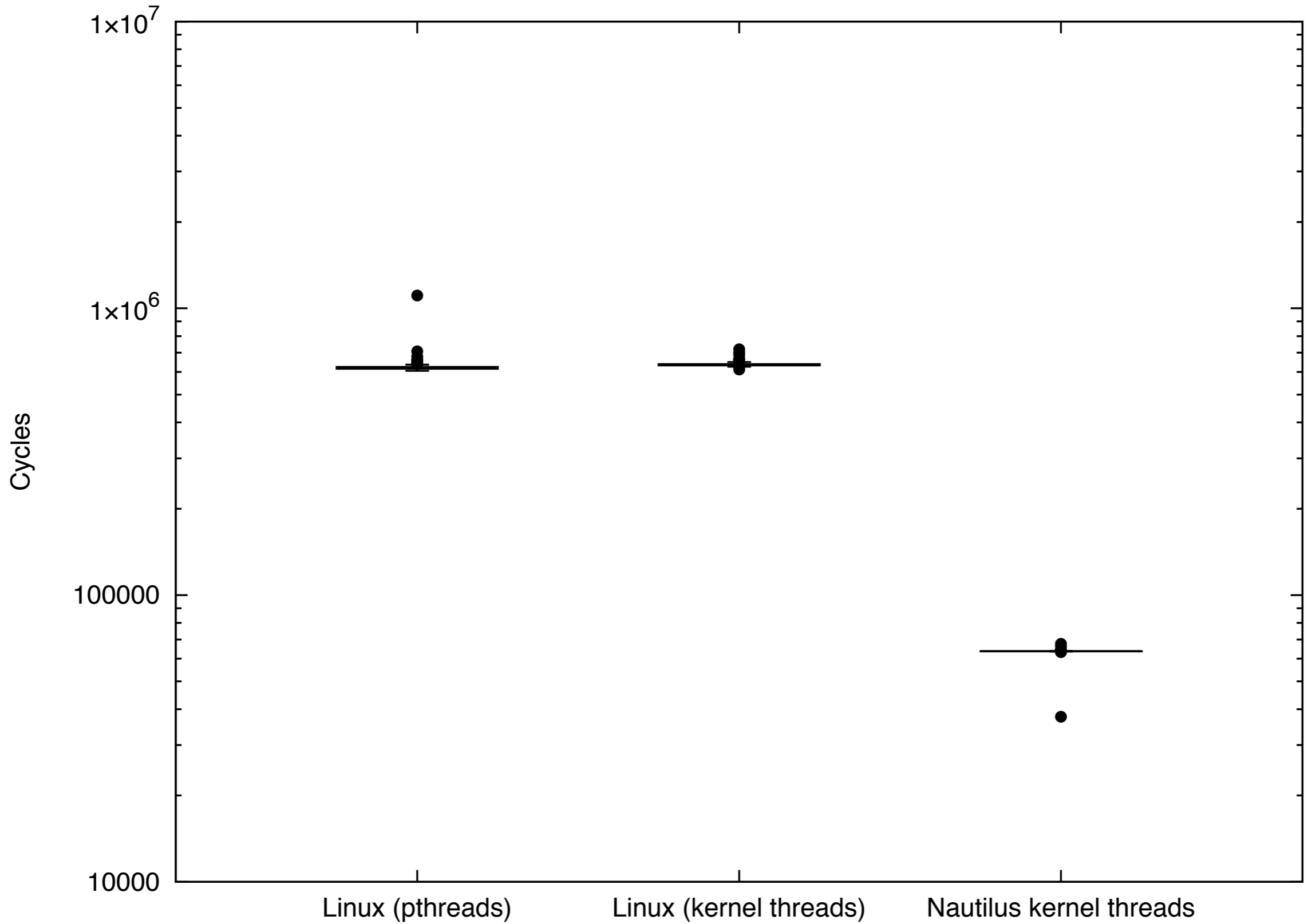
Wakeup deviation on x64



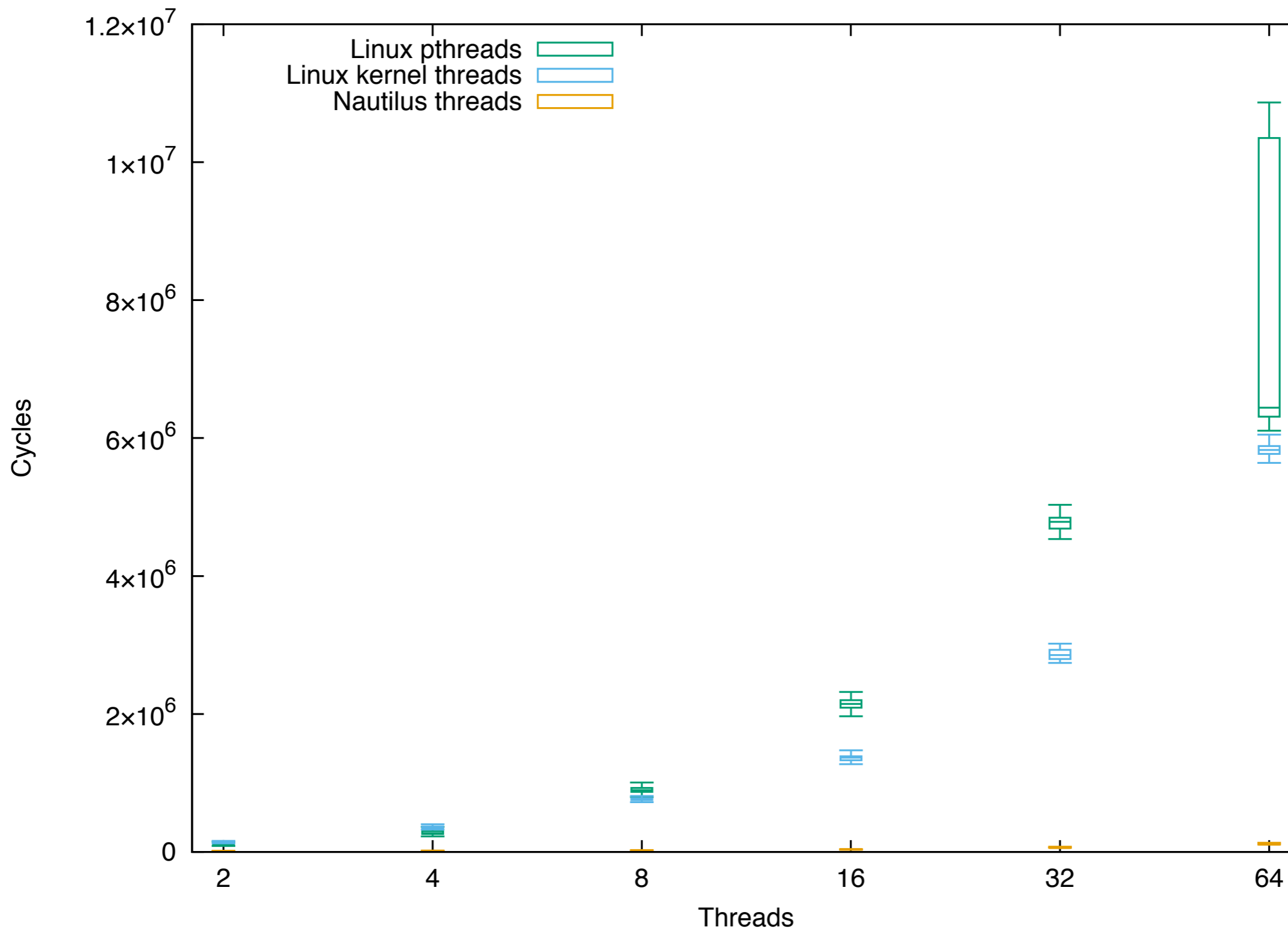
thread context switch on x64



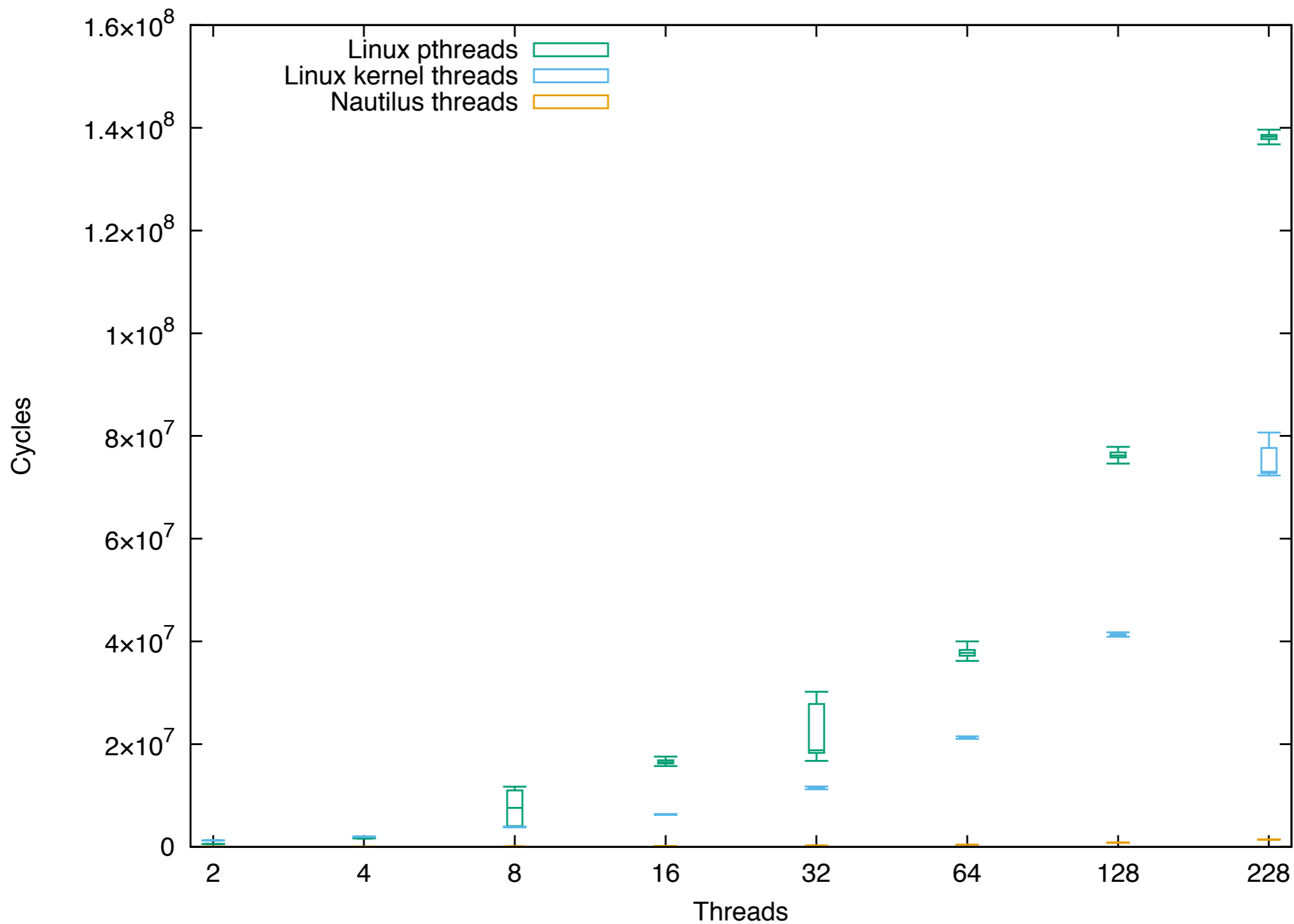
thread create + launch on phi



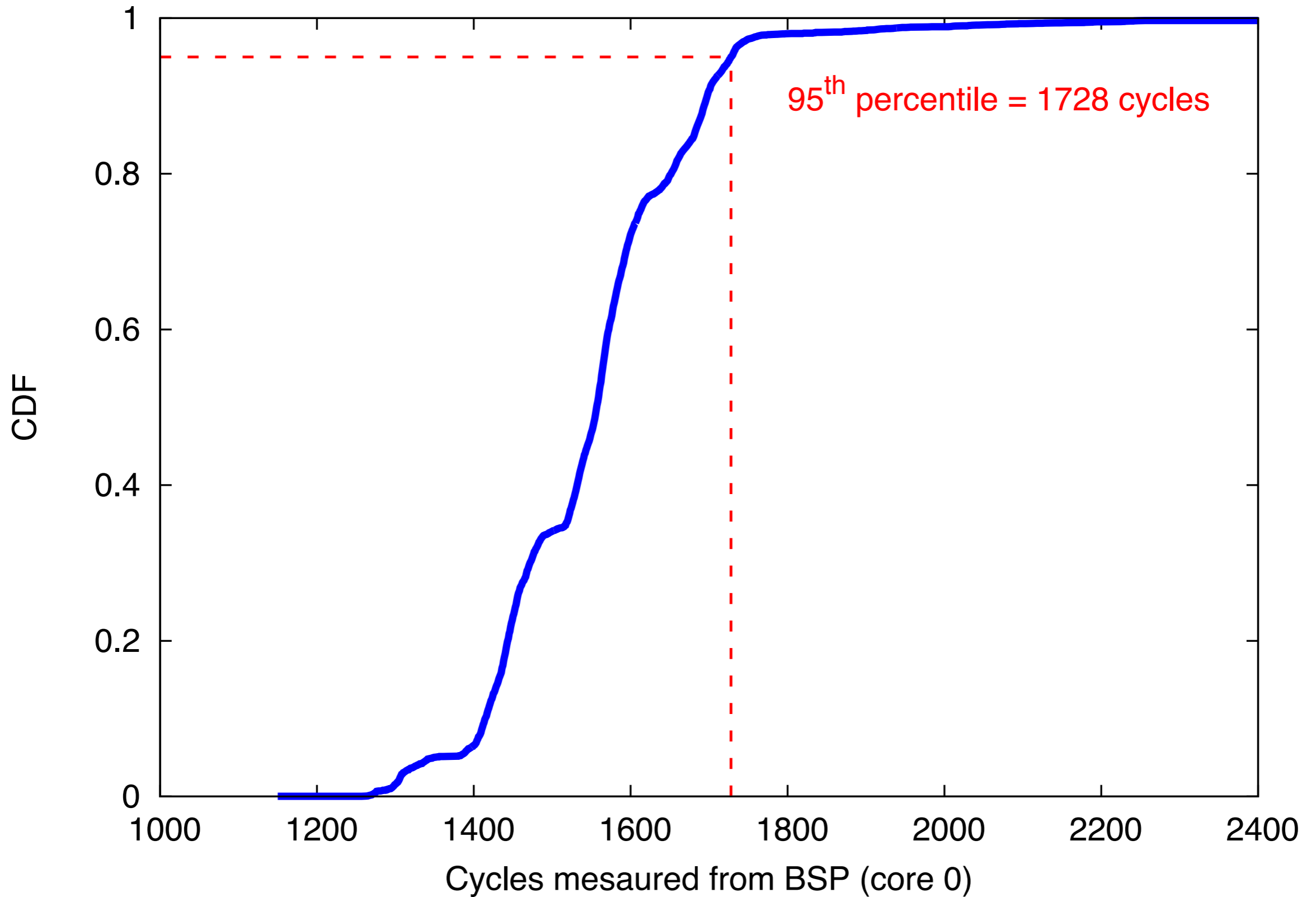
thread create + launch (many threads) on x64



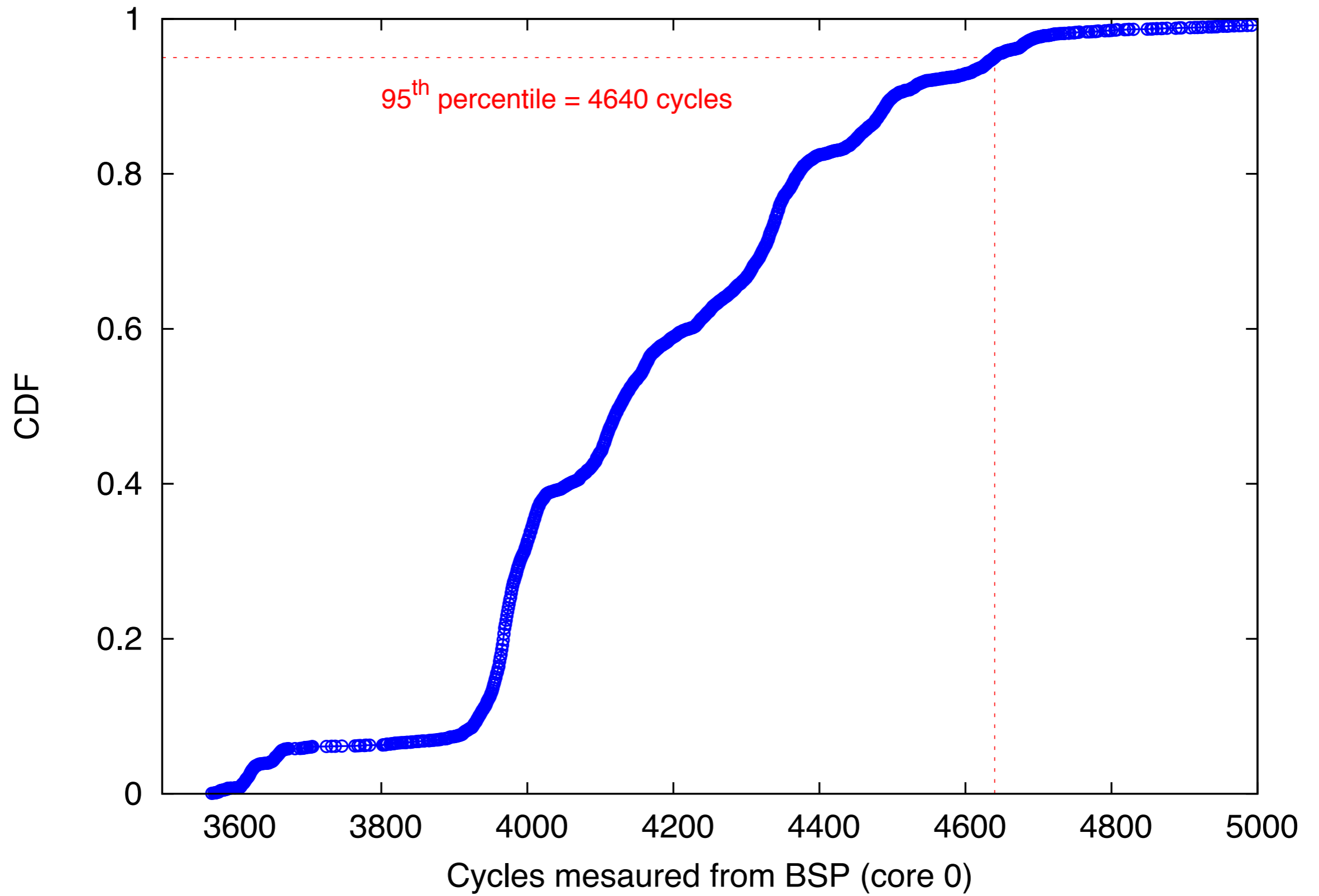
thread create + launch (many threads) on phi



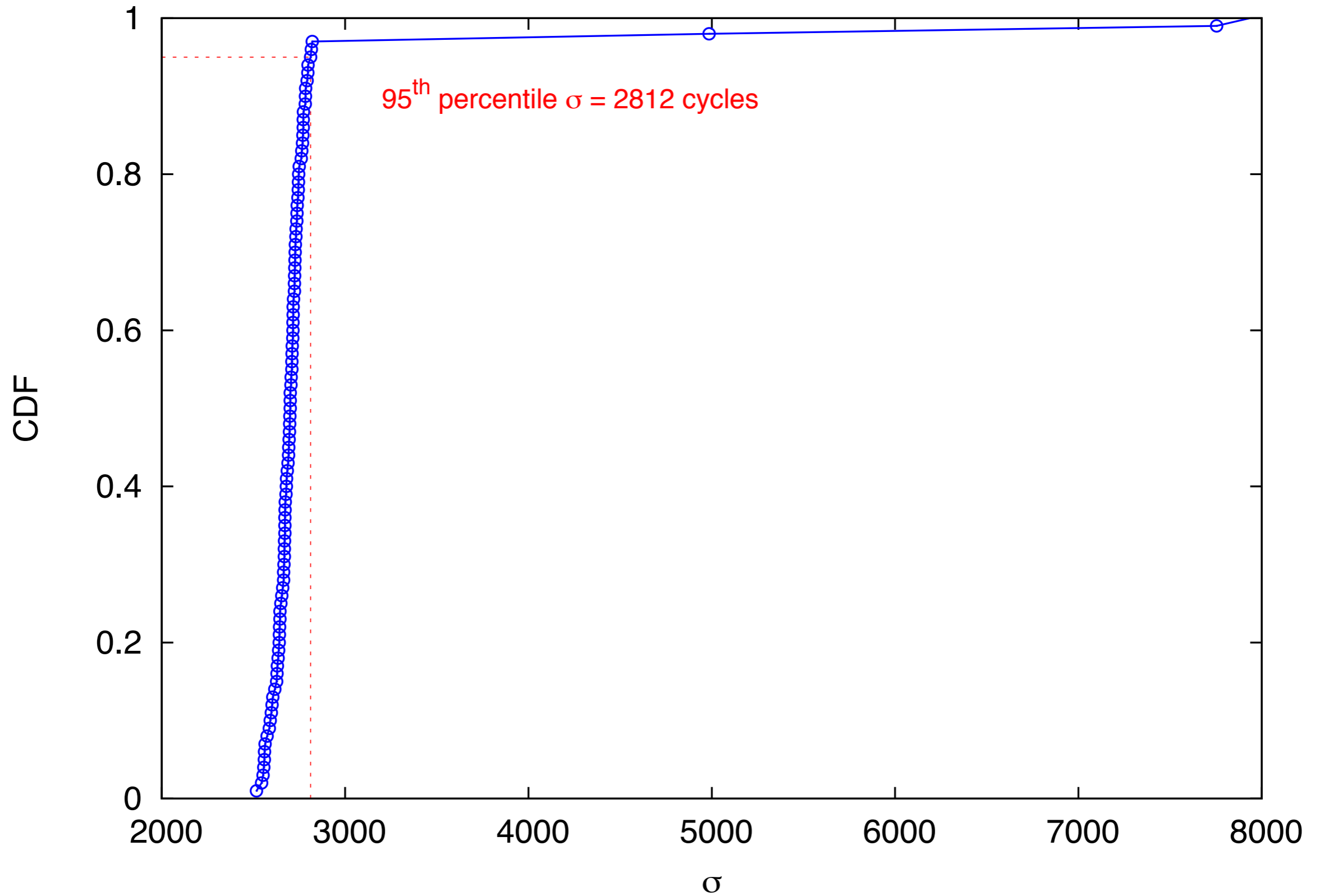
Unicast IPs on x64



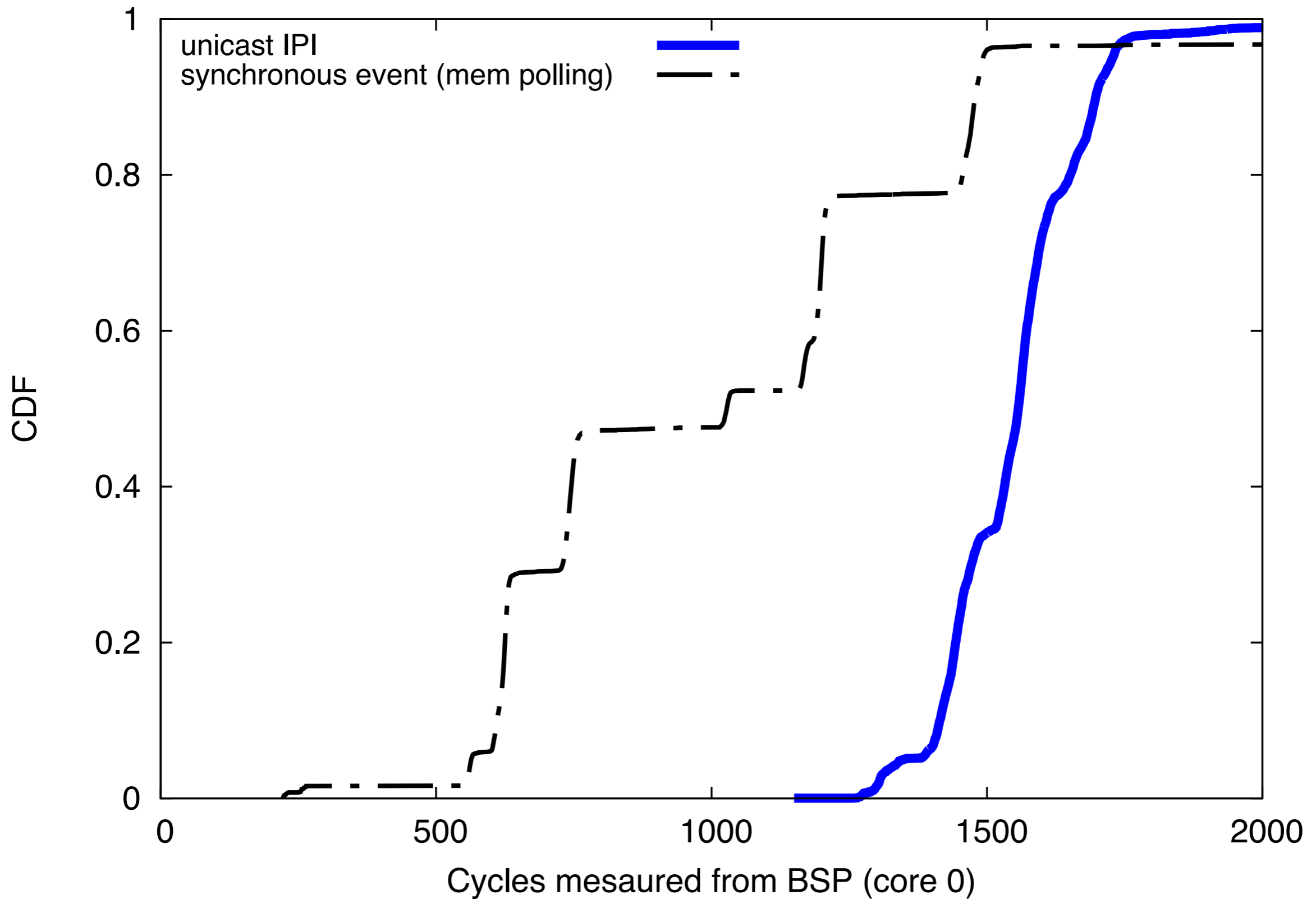
Roundtrip IPs on x64



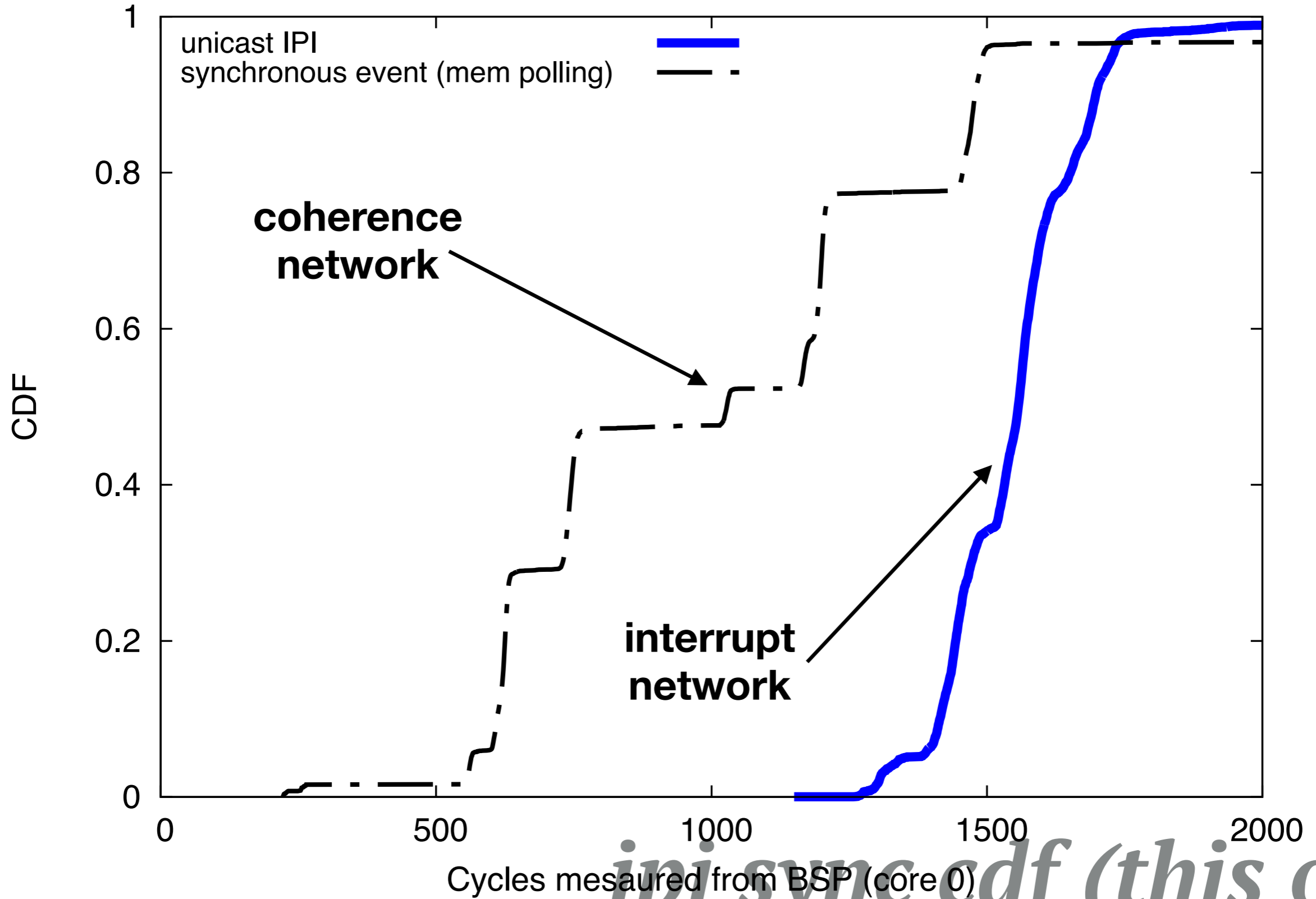
Multicast IPs on x64



Unicast IPI vs memory polling

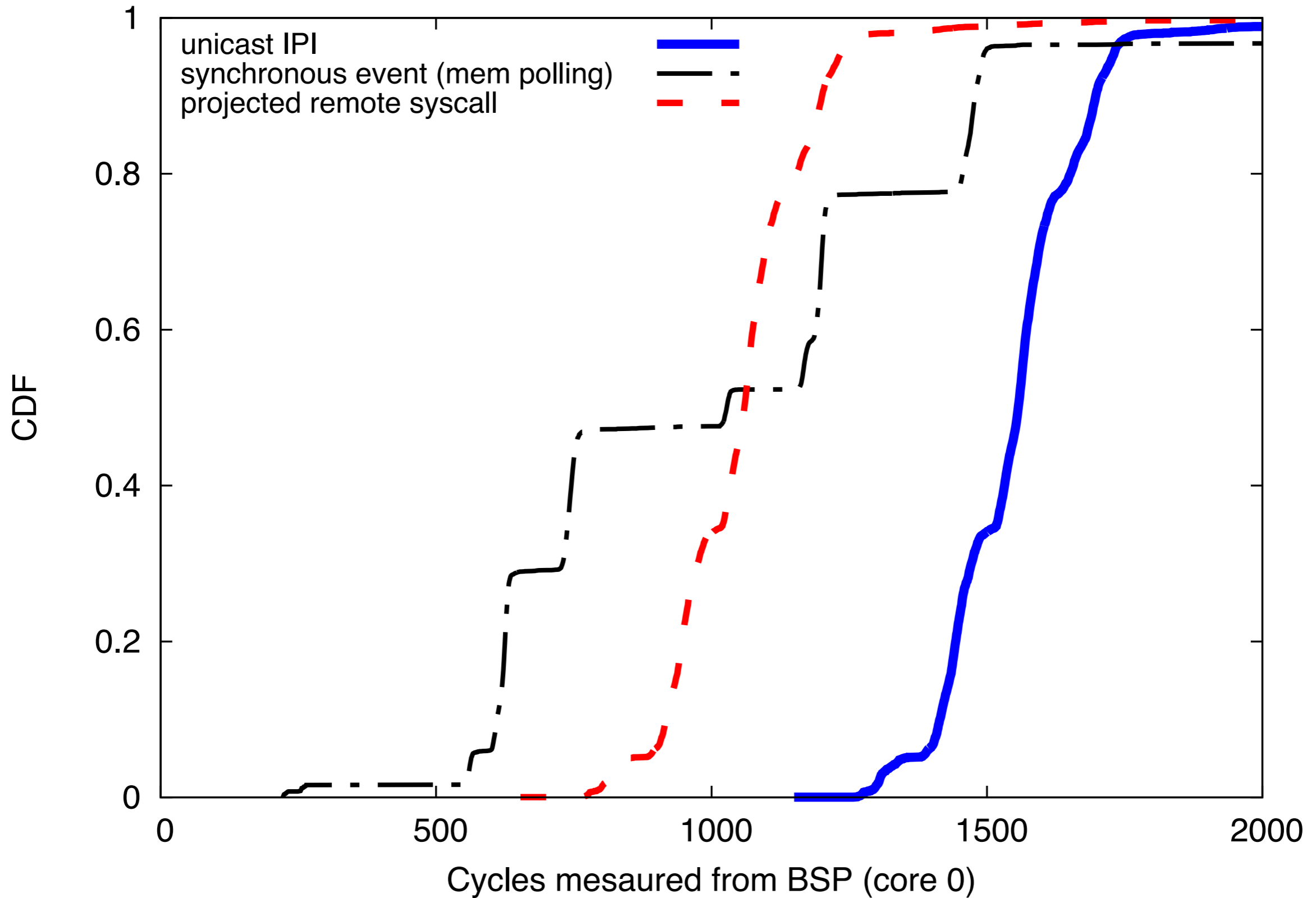


Unicast IPI vs memory polling

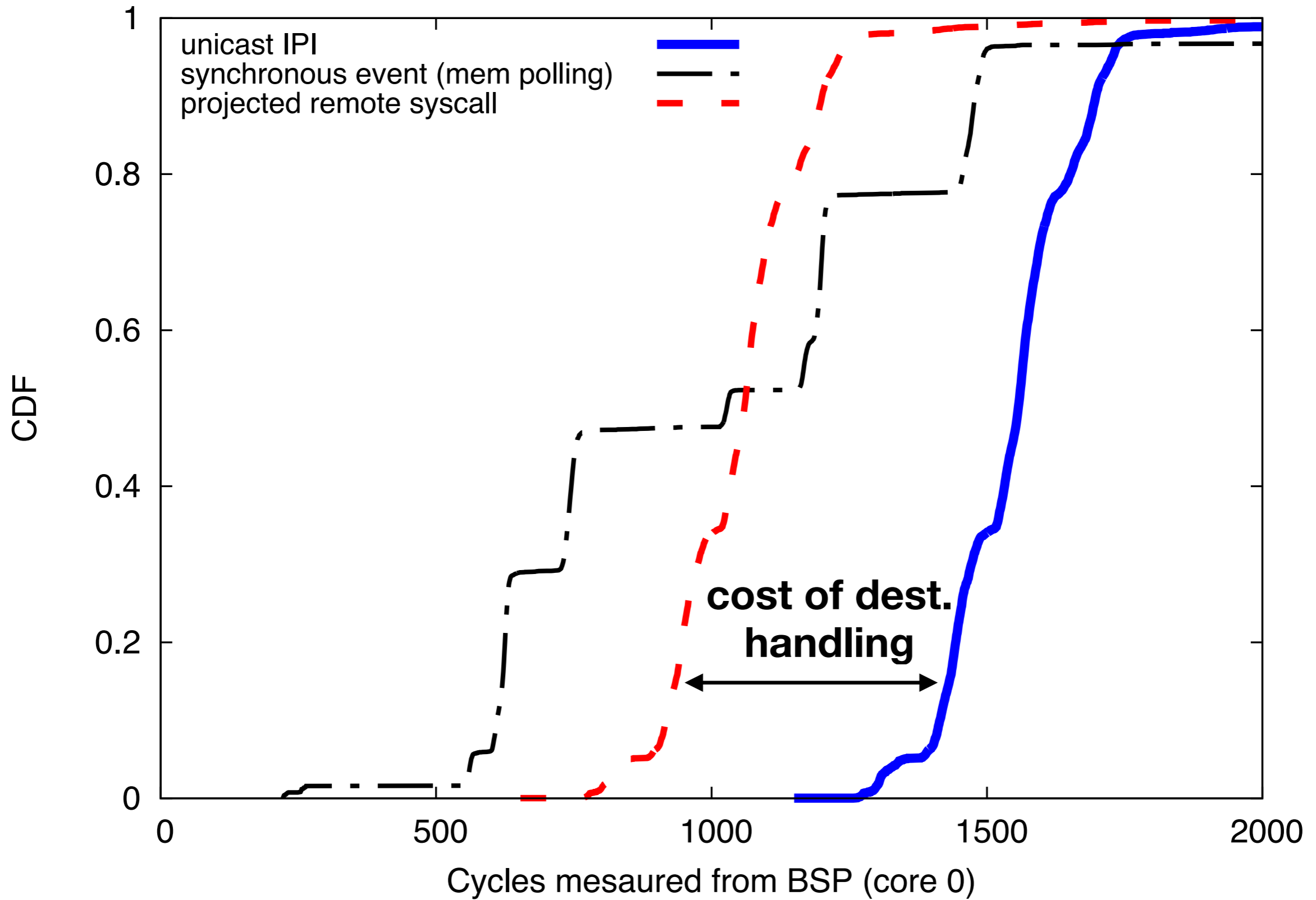


ipi sync cdf (this one annotations)

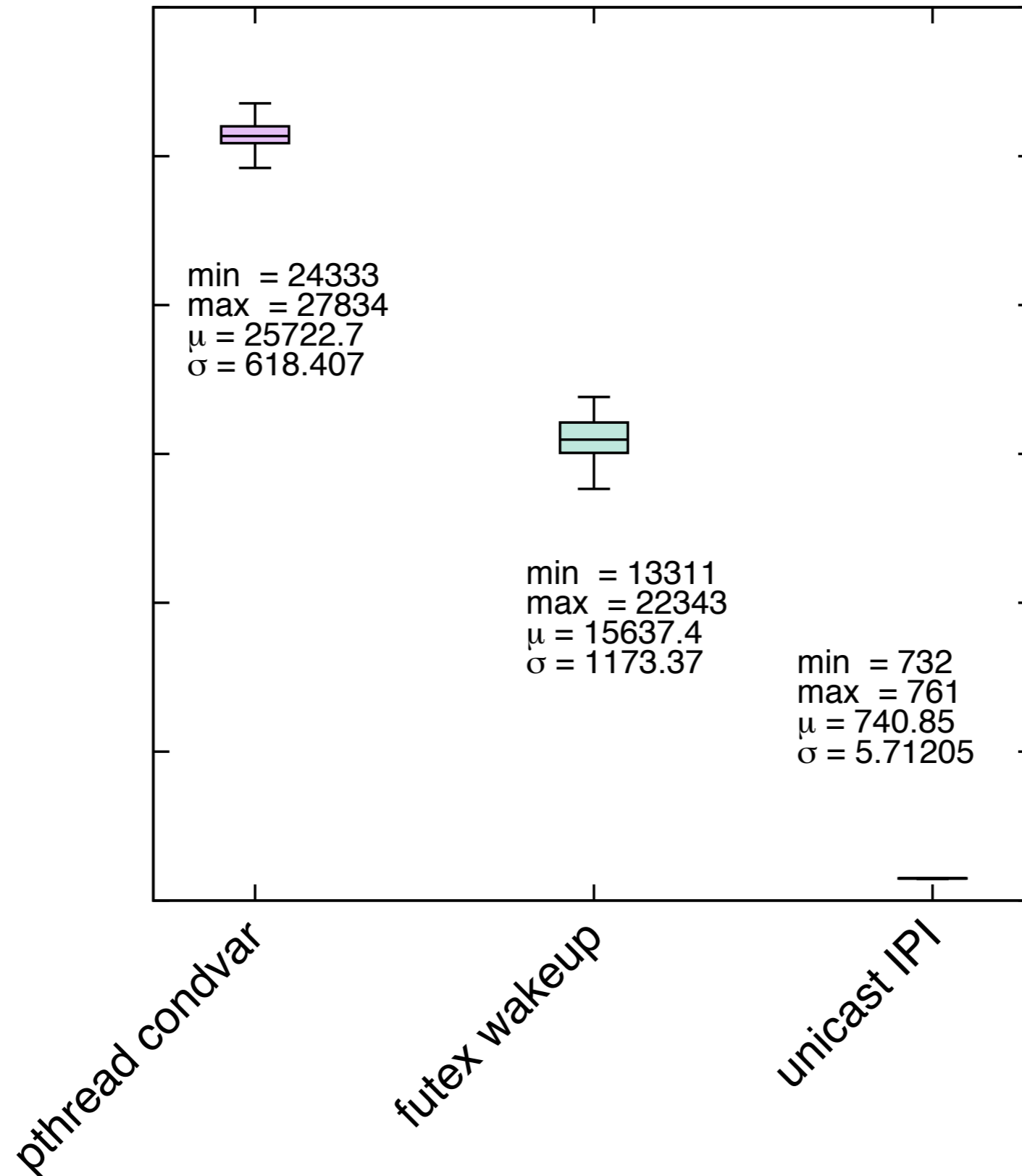
Unicast IPI vs memory polling



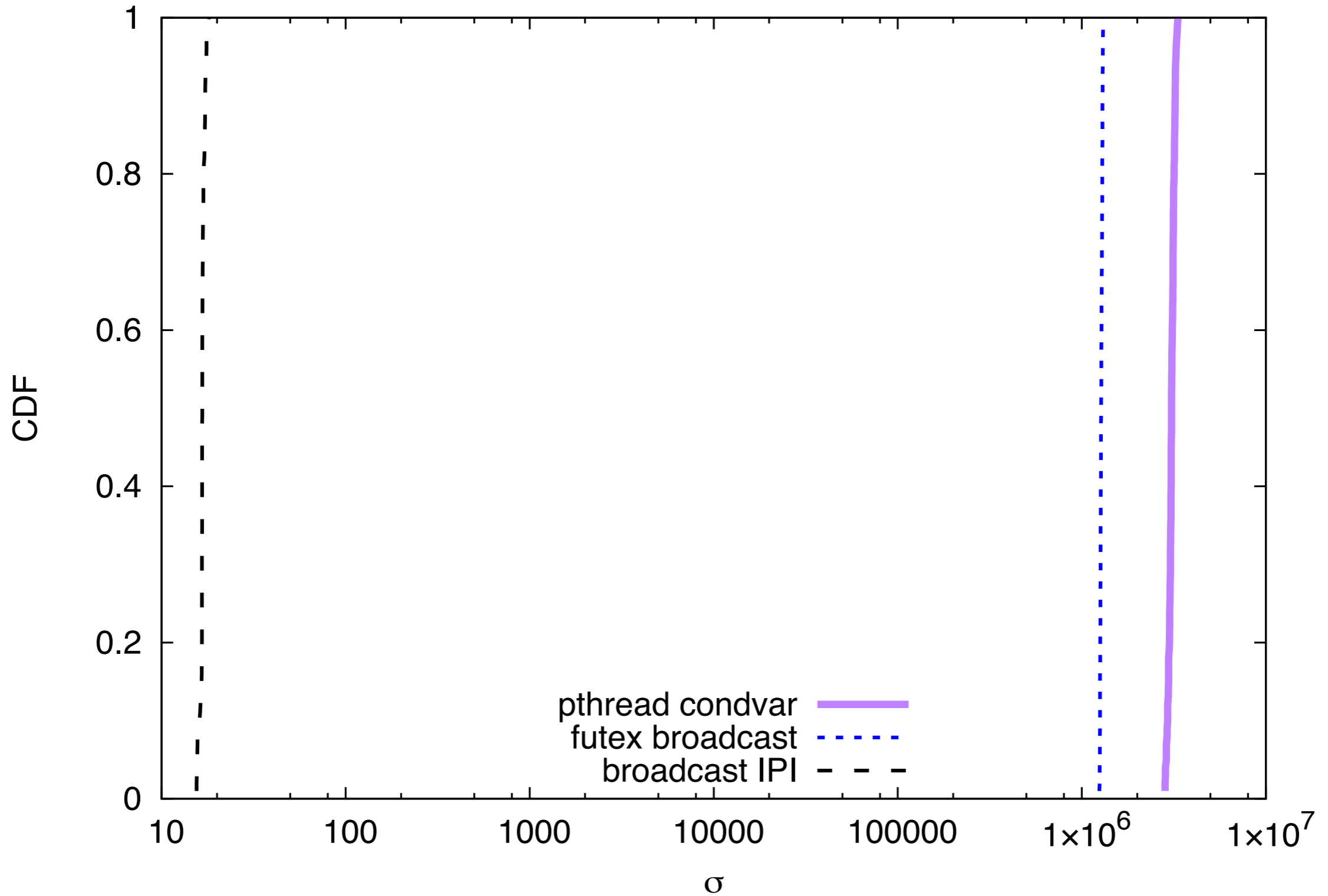
Unicast IPI vs memory polling



Single wakeup on phi

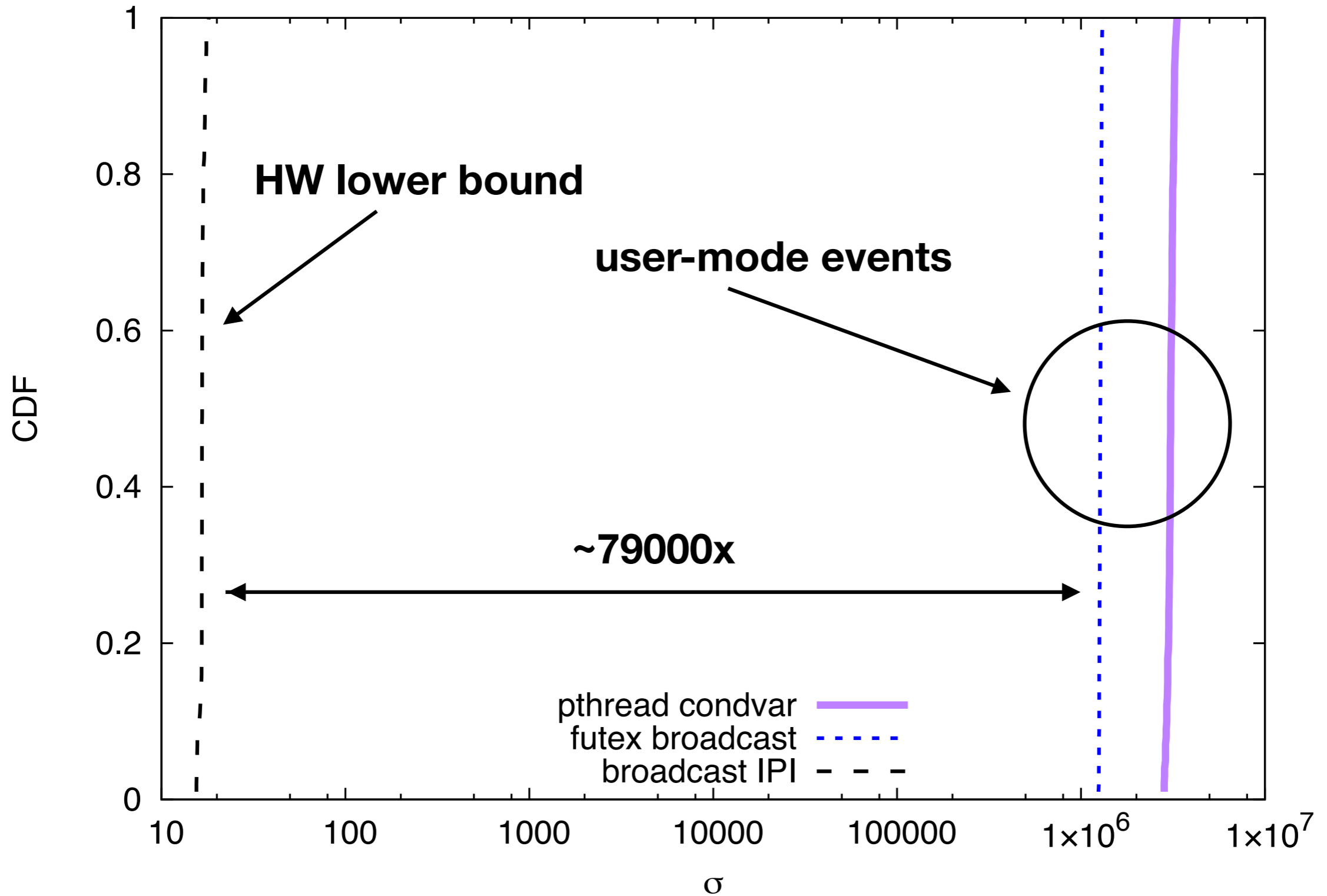


Wakeup deviation on phi

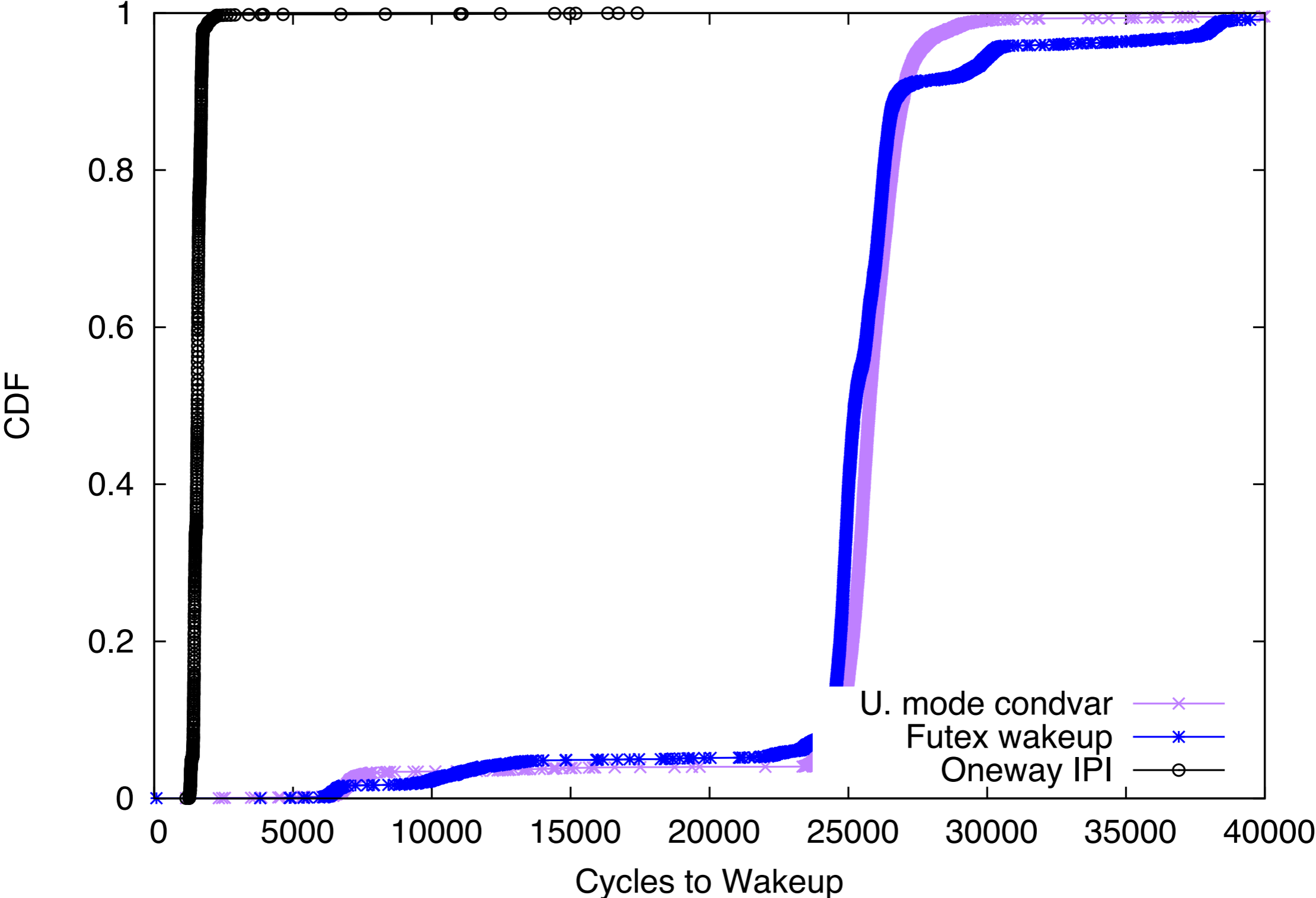


wakeup phi many

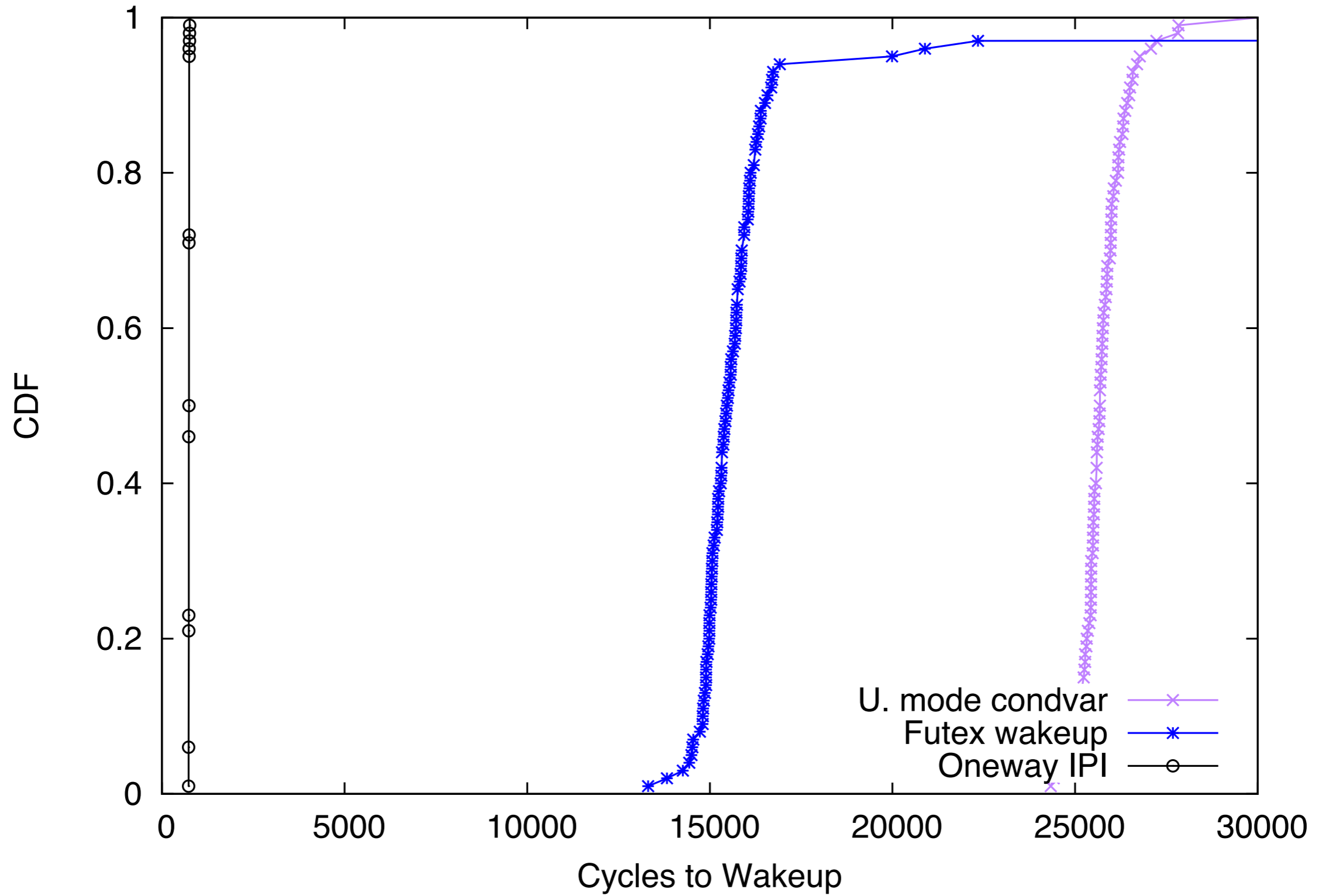
Wakeup deviation on phi



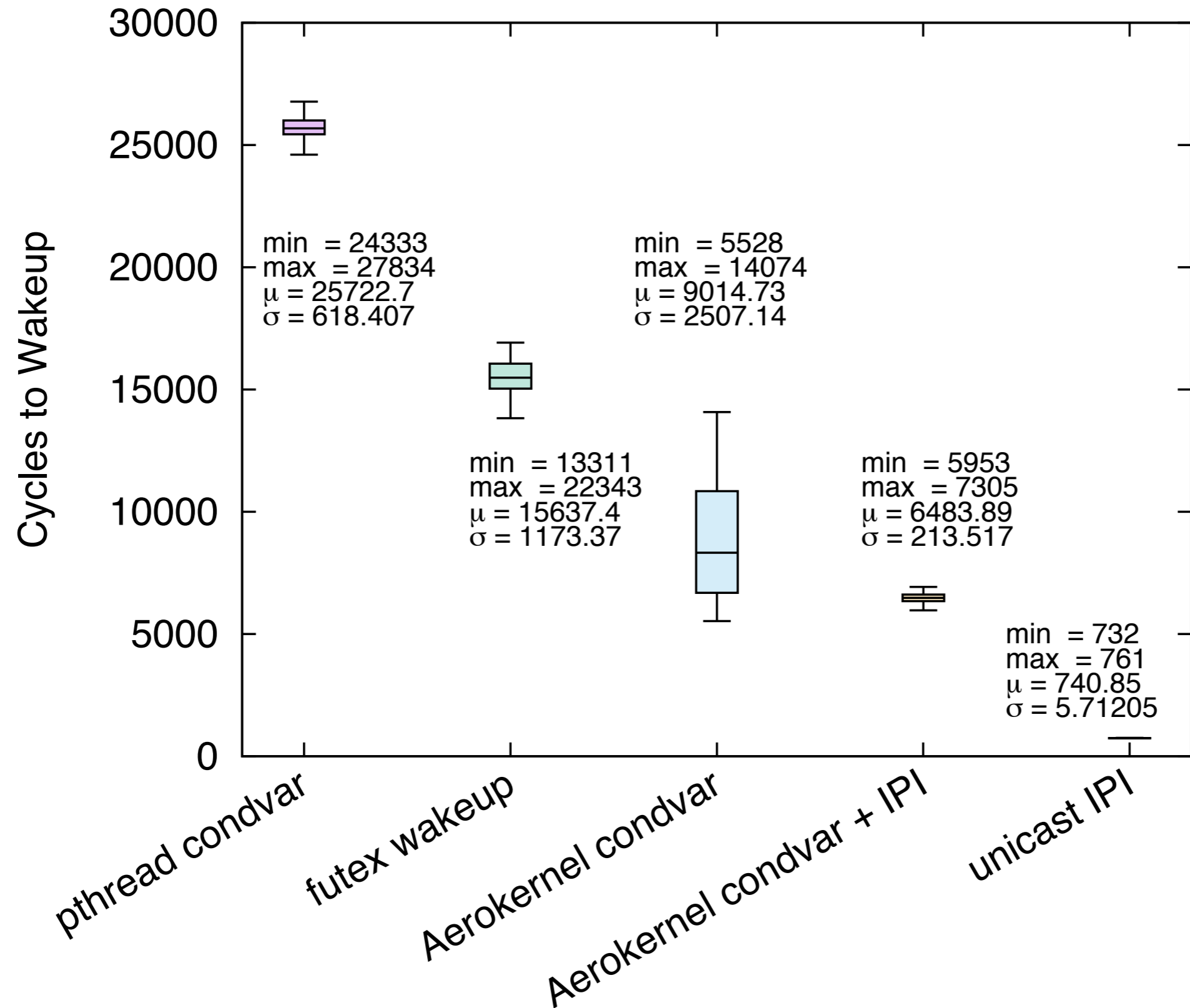
Single wakeup on x64 (CDF)



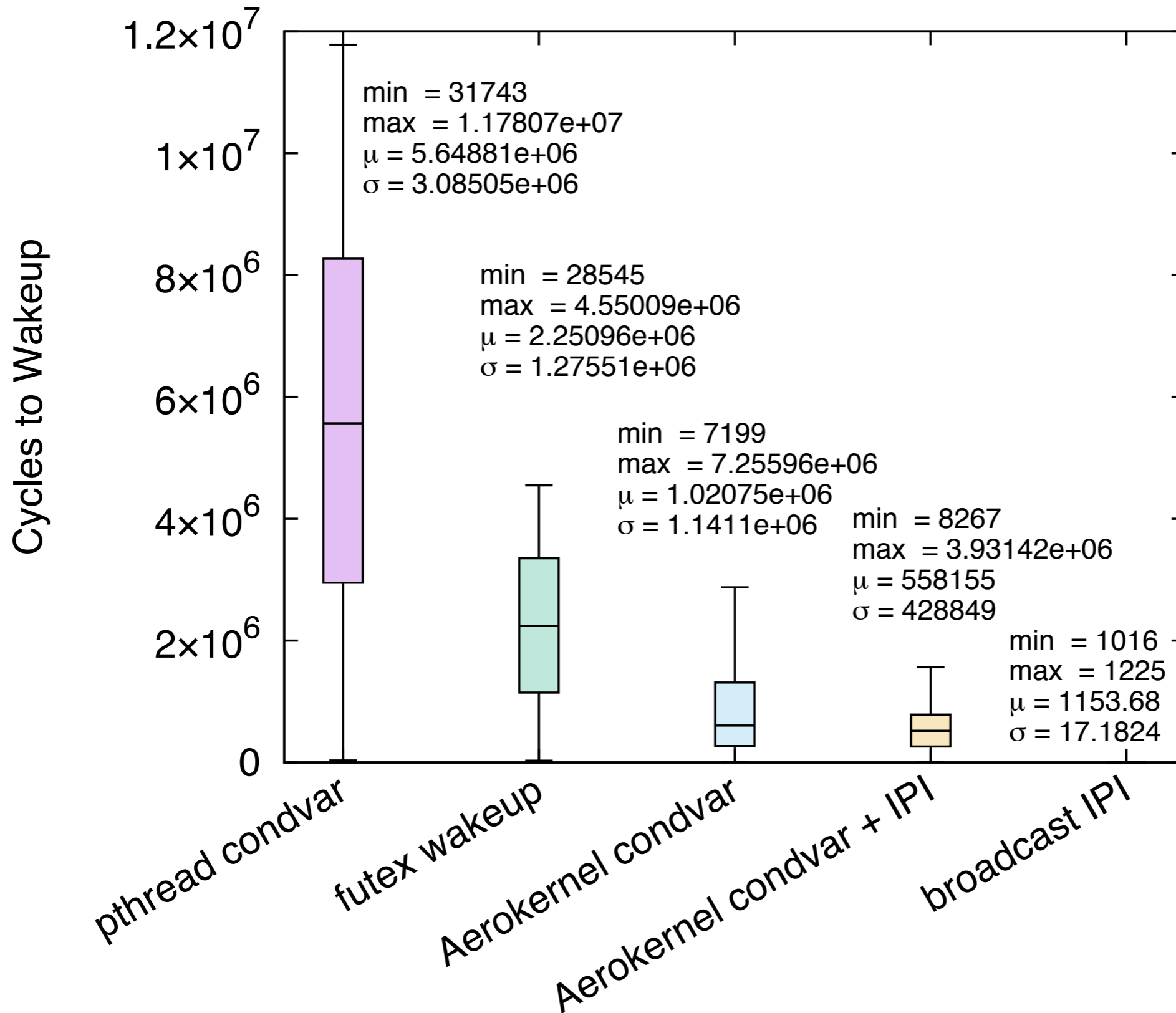
Single wakeup on phi (CDF)



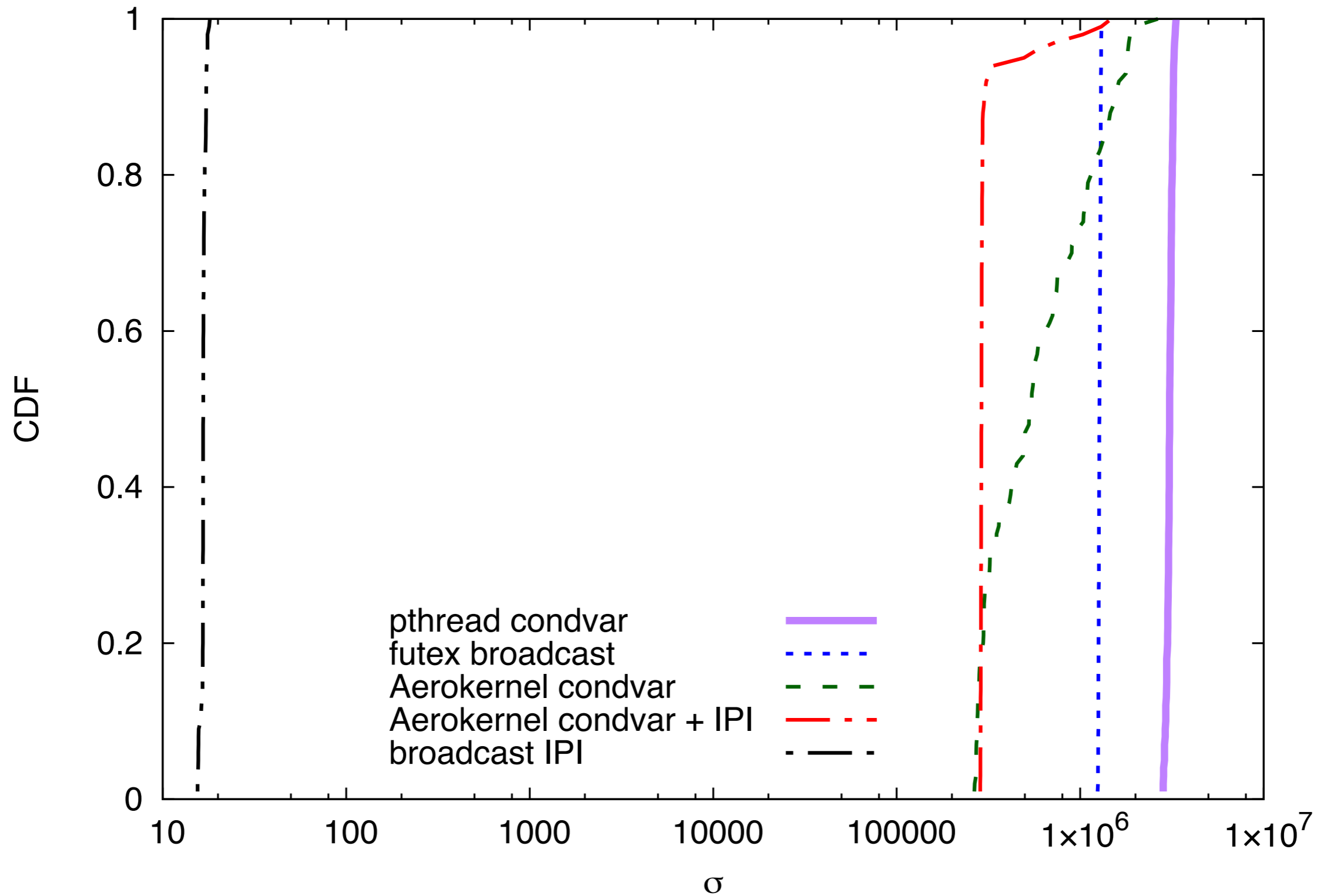
Single wakeup on phi



Many wakeups on phi

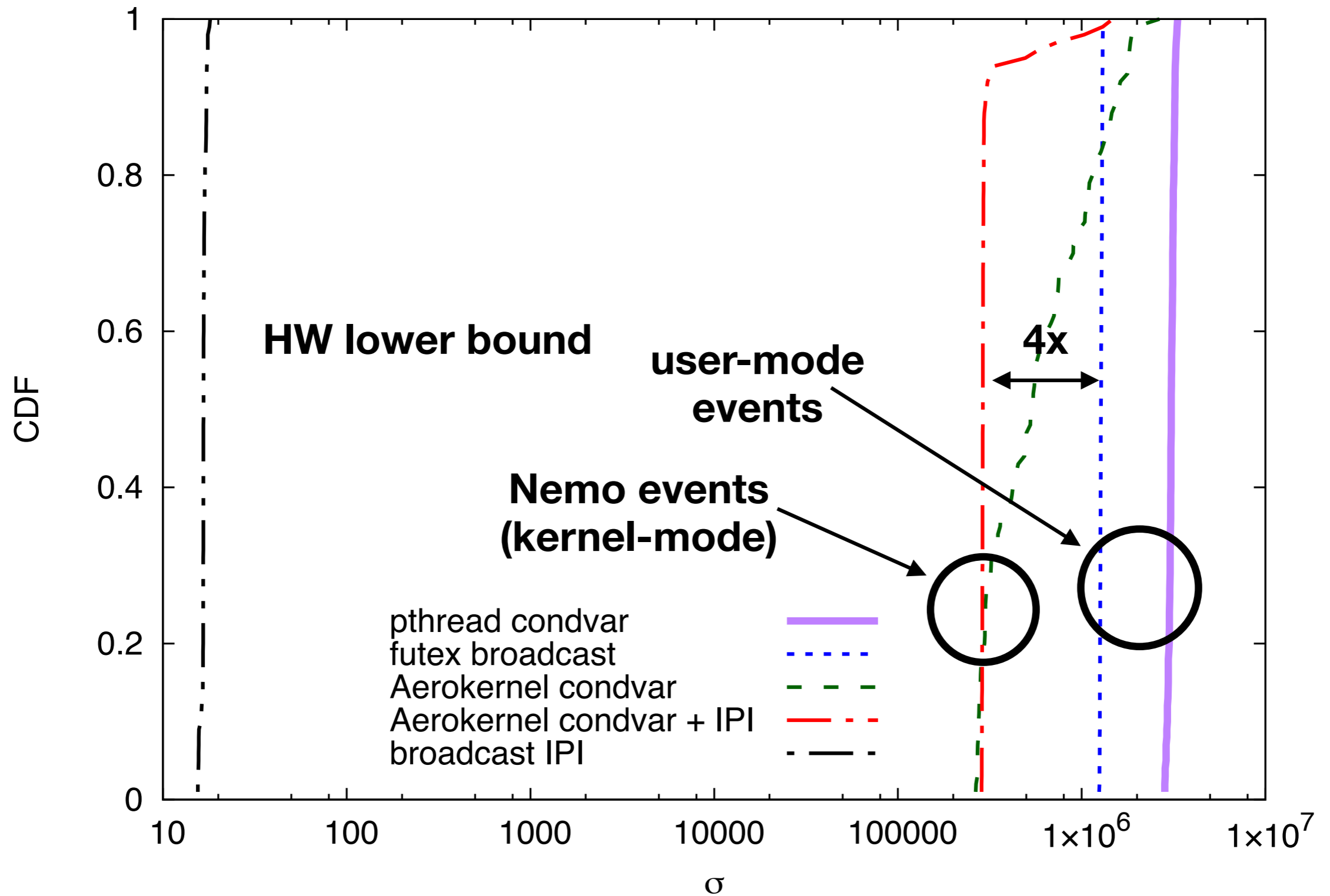


Wakeup deviation on phi

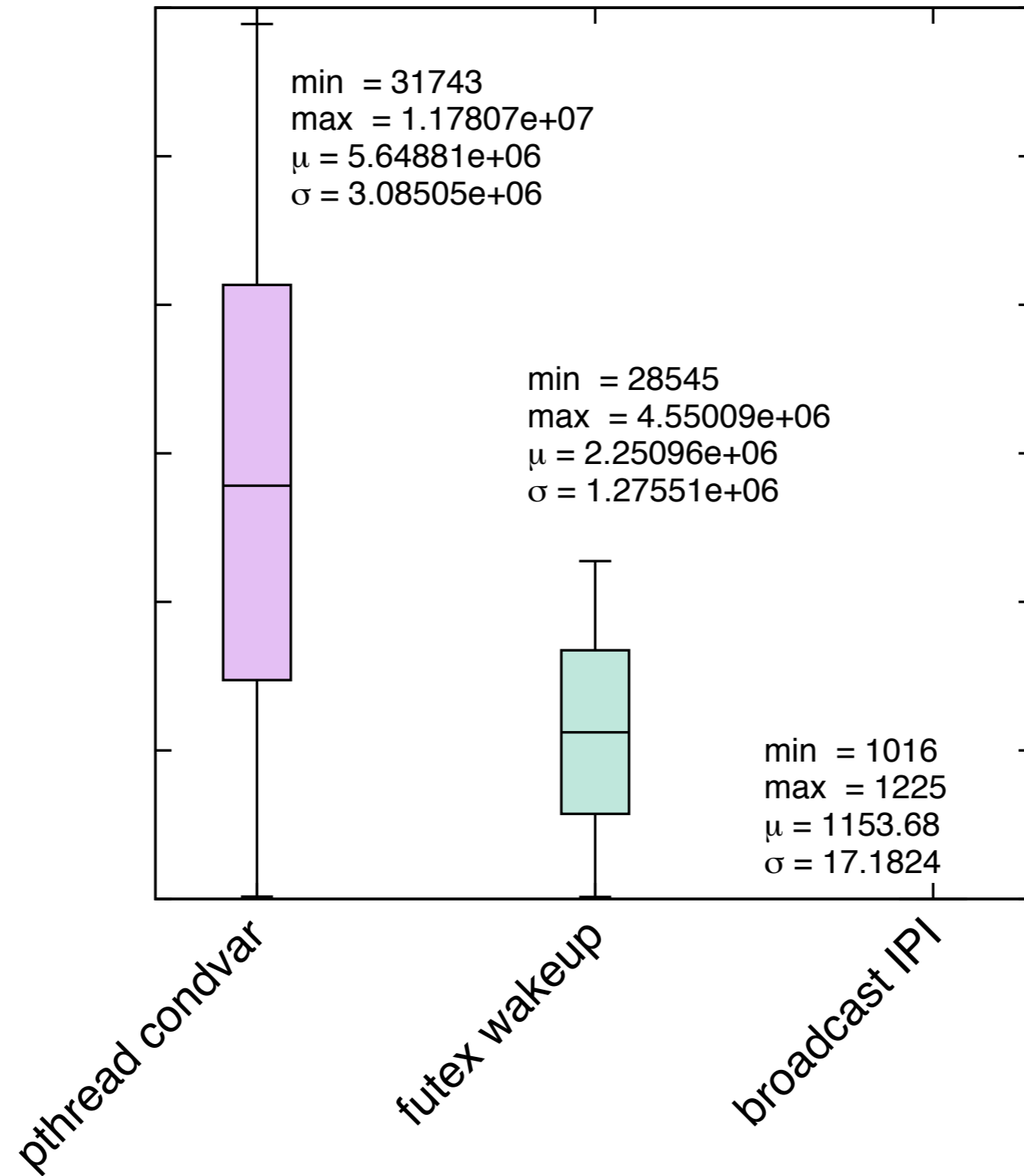


wakeup phi many wk

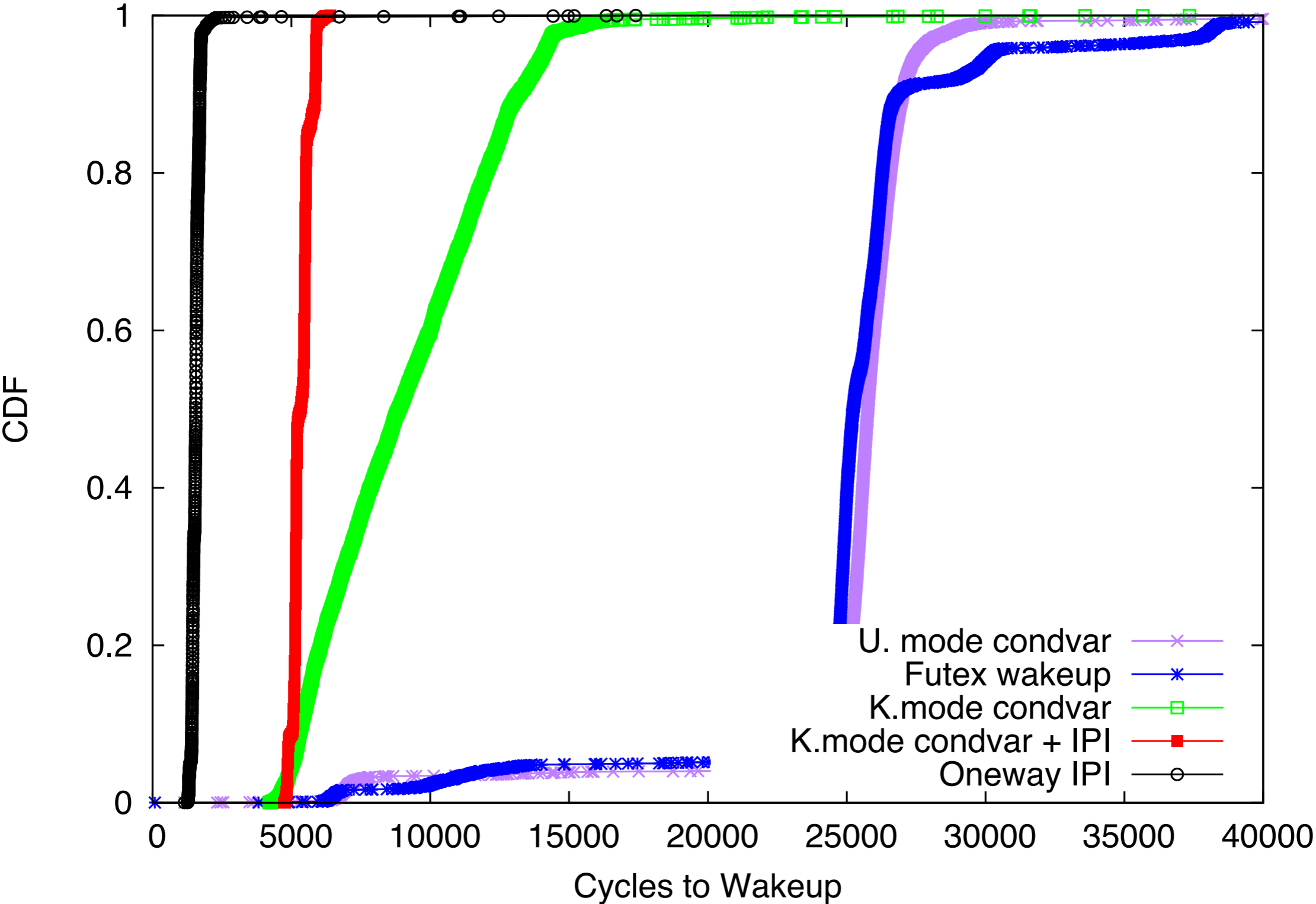
Wakeup deviation on phi



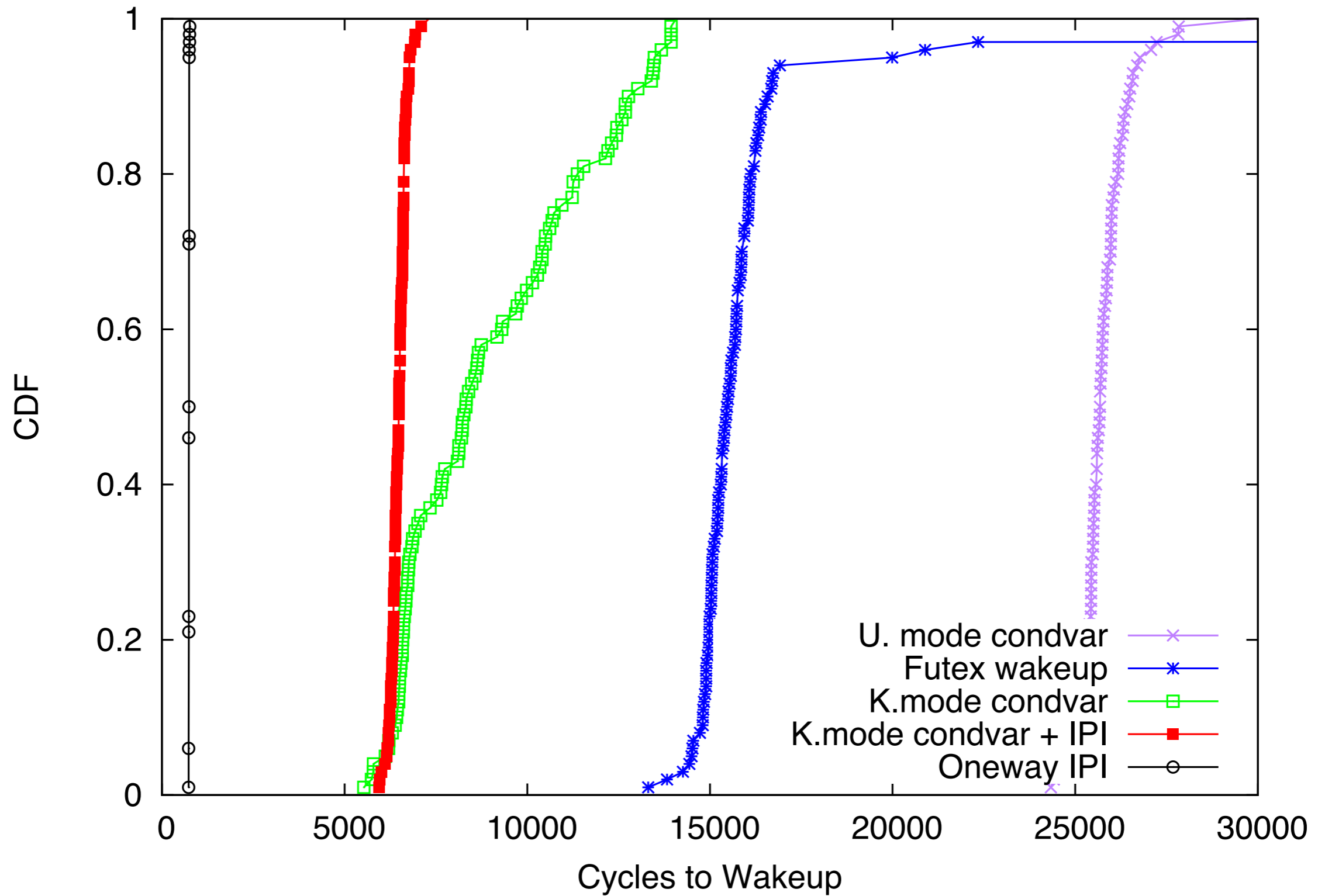
Broadcast wakeups on phi



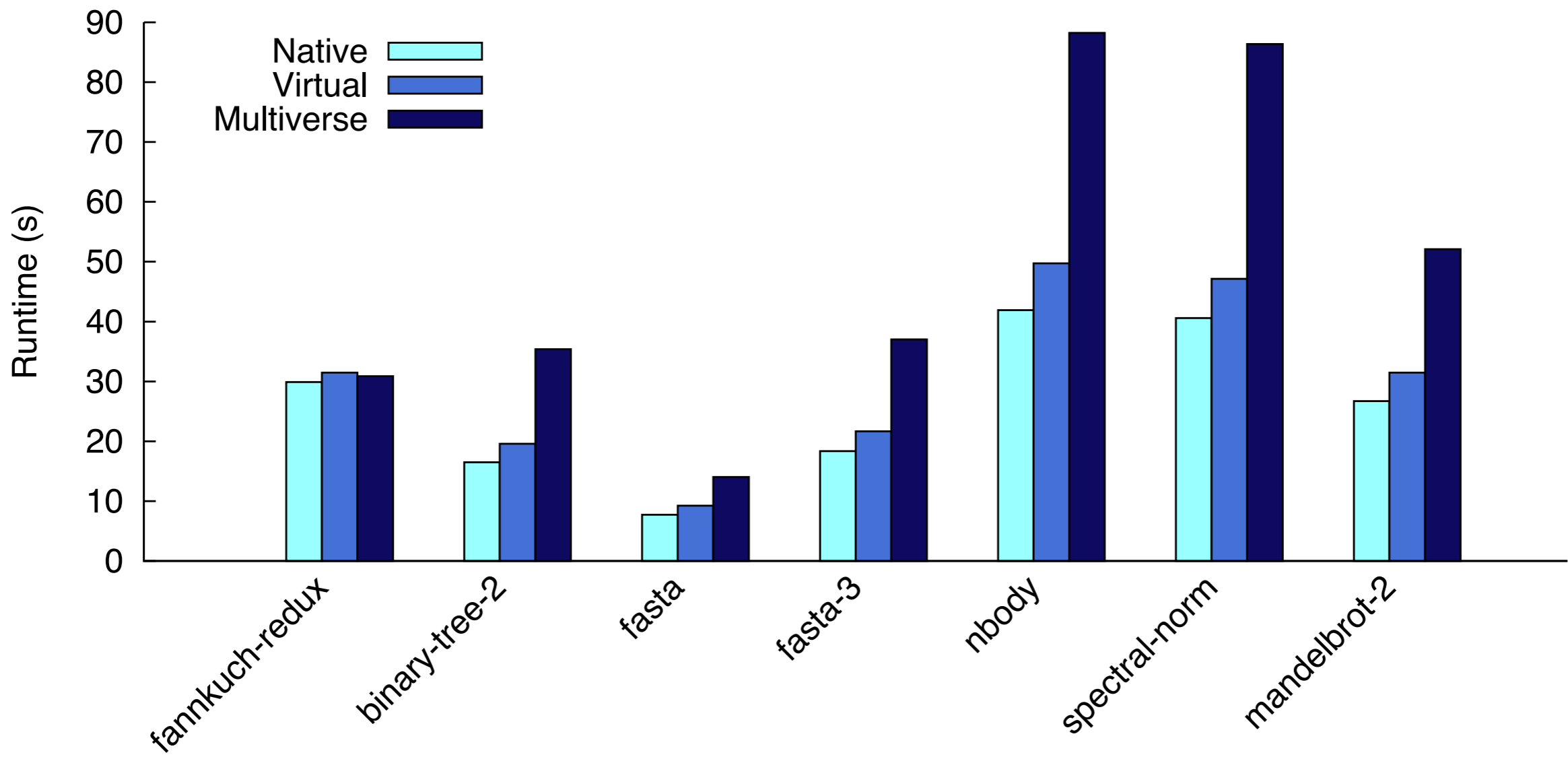
Single wakeup on x64 (CDF)



Single wakeup on phi (CDF)



racket multiverse overheads



system call overheads

