





Tapestry

Revealing Wait-For Dependencies Between Application Threads

Tomáš Faltín^{1,2}, Himadri Chhaya-Shailesh², Julia Lawall², Jean-Pierre Lozi²

¹Charles University, ²Inria Paris



Synchronization bottlenecks

Multithreaded
applications face
performance issues
due to
synchronization
causing delays and
inefficiencies.



Limitations of traditional profiling

focus on function or hardware slowdowns but *miss* inter-thread *dependencies* impacting performance.



Challenges of busy-waiting

Busy-waiting
threads spin in
userspace, hiding
actual wait times as
continuous execution
in traces.



Full picture visualization



Synchronization bottlenecks

Multithreaded
applications face
performance issues
due to
synchronization
causing delays and
inefficiencies.



Limitations of traditional profiling

focus on function or hardware slowdowns but *miss* inter-thread *dependencies* impacting performance.



Challenges of busy-waiting

Busy-waiting
threads spin in
userspace, hiding
actual wait times as
continuous execution
in traces.



Full picture visualization



Synchronization bottlenecks

Multithreaded
applications face
performance issues
due to
synchronization
causing delays and
inefficiencies.



Limitations of traditional profiling

Conventional tools
focus on function or
hardware slowdowns
but miss inter-thread
dependencies
impacting
performance.



Challenges of busy-waiting

Busy-waiting
threads spin in
userspace, hiding
actual wait times as
continuous execution
in traces.



Full picture visualization



Synchronization bottlenecks

Multithreaded
applications face
performance issues
due to
synchronization
causing delays and
inefficiencies.



Limitations of traditional profiling

focus on function or hardware slowdowns but *miss* inter-thread *dependencies* impacting performance.



Challenges of busy-waiting

Busy-waiting
threads spin in
userspace, hiding
actual wait times as
continuous execution
in traces.



Full picture visualization



Synchronization bottlenecks

Multithreaded
applications face
performance issues
due to
synchronization
causing delays and
inefficiencies.



Limitations of traditional profiling

focus on function or hardware slowdowns but *miss* inter-thread *dependencies* impacting performance.



Challenges of busy-waiting

Busy-waiting
threads spin in
userspace, hiding
actual wait times as
continuous execution
in traces.



Full picture visualization

Synchronization

Busy-waiting examples

```
/* Barrier */
if (atomic_inc(thr_count) < n_thr)</pre>
  while (*barrier_end != true)
    PAUSE;
else *barrier_end = true;
/* Simple lock */
while (atomic_test_and_set(lock, 0, 1)
       != 0)
  PAUSE;
/* Ticket lock */
tkt = atomic_inc(max_tkt);
while (*cur_tkt < tkt) PAUSE;</pre>
```

Blocking example

```
futex(*addr, WAIT, val, ...)
```

Synchronization

Busy-waiting examples

```
/* Barrier */
if (atomic_inc(thr_count) < n_thr)</pre>
  while (*barrier end != true)
    PAUSE;
else Cbarrièr end = true;
/* Simple lock */
while (atomic_test_and_set(lock, 0, 1)
       != 0)
  PAUSE;
/* Ticket lock */
tkt = atomic_inc(max_tkt);
while (*cur_tkt) < tkt) PAUSE;</pre>
```

Blocking example

```
futex *addr, WAIT, val, ...)

?
Wait-for Dependency
```

Synchronization

Busy-waiting examples

Blocking example

```
/* Barrier */
if (atomic inc(thr count) < n thr)</pre>
                                       futex(*addr, WAIT, val, ...)
 while (*barrier end)!= true)
   PA
                 How to unify the concept of
else 🔇
           dependencies for synchronization?
/* Sim
while
       ! = 0)
 PAUSE;
/* Ticket lock */
tkt = atomic_inc(max_tkt);
while (*cur_tkt) < tkt) PAUSE;</pre>
```

Busy-waiting examples

```
/* Barrier */
if (atomic_inc(thr_count) < n_thr)</pre>
  while (*barrier_end != true)
    PAUSE;
else *barrier end = true;
/* Simple lock */
while (atomic_test_and_set(lock, 0, 1)
       != 0)
  PAUSE;
/* Ticket lock */
tkt = atomic_inc(max_tkt);
while (*cur_tkt < tkt) PAUSE;</pre>
```

Blocking example

```
futex(*addr, WAIT, val, ...)
```

Identify common patterns

Busy-waiting examples

```
/* Barrier */
if (atomic_inc(thr_count) < n_thr)</pre>
  while (*barrier end!= true)
    PAUSE;
else *barrier end = true;
/* Simple lock */
while (atomic_test_and_set(lock)
       != 0)
  PAUSE;
/* Ticket lock */
tkt = atomic_inc(max_tkt);
while (*cur_tkt) < tkt) PAUSE;</pre>
```

Blocking example

```
futex(*addr) WAIT, val, ...)
```

The shared variable memory address

```
<addr, ???, ???>
```

Busy-waiting examples

```
/* Barrier */
if (atomic_inc(thr_count) < n_thr)</pre>
  while (*barrier end (!=) true)
    PAUSE;
else *barrier end = true;
/* Simple lock */
while (atomic_test_and_set(lock) 0, 1)
  PAUSE;
/* Ticket lock */
tkt = atomic_inc(max_tkt);
while (*cur tkt(<)tkt) PAUSE;
```

```
Blocking example
```

blocks if the futex word equals to the supplied val

```
futex(*addr) WAIT, val, ...)
```

A comparison operator



Busy-waiting examples

```
/* Barrier */
if (atomic inc(thr count) < n thr)</pre>
  while (*barrier end !=) true
    PAUSE;
else *barrier end = true;
/* Simple lock */
while (atomic_test_and_set(lock)
  PAUSE;
/* Ticket lock */
tkt = atomic_inc(max_tkt);
while (*cur tkt < )tkt
```

blocks if the futex word Blocking example equals to the supplied val

```
futex(*addr)
             WAIT,
```

A constant value to compared against







Busy-waiting exa <barrier_end, ==, true>

```
/* Barrier */
if (atomic inc(thr count) < n thr)</pre>
  while (*barrier end !=) true
    PAUSE;
else *barrier end = true;
                           <lock, ==, 0>
/* Simple lock */
while (atomic_test_and_set(lock)
  PAUSE;
                        <cur_tkt, ==, 0>
/* Ticket lock */
tkt = atomic_inc(max_tkt);
while (*cur tkt)< tkt
```

Blocking example blocks if the futex word equals to the supplied val

```
futex(*addr WAIT, val) ...)

<addr, ==, val>
```



blocks if the futex word Busy-waiting exa <barrier_end, ==, true> Blocking example equals to the supplied val Barrier * if (atomic inc(thr count) < n thr)</pre> futex(*addr) WAIT, while (*barrier end != true) PA How to capture ordinary inter-thread waitfor dependencies? /* Sim while !=)0) A dependency triple: PAUSE; <cur_tkt, ==, 0> /* Ticket lock */ tkt = atomic_inc(max tkt); while (*cur tkt)< tk

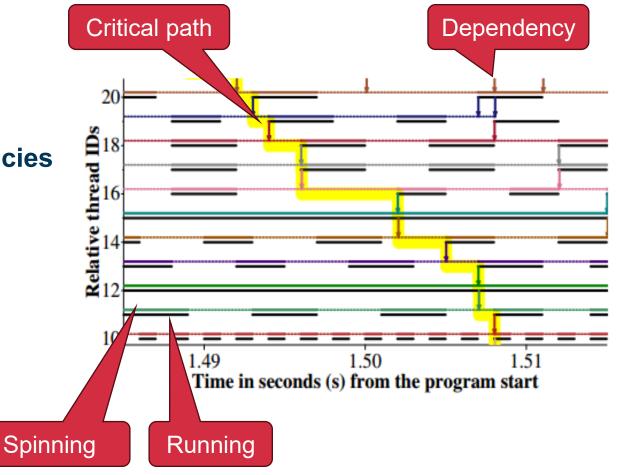
blocks if the futex word Busy-waiting exa <barrier_end, ==, true> Blocking example equals to the supplied val Barrier */ if (atomic inc(thr count) < n thr)</pre> futex(*addr) WAIT, while (*barrier end != true) PA else * **Tapestry** /* Sim while ! =) 0)A dependency triple: PAUSE; ∢addr), (⊳< <cur_tkt, ==, 0> /* Ticket lock */ tkt = atomic_inc(max_tkt); while (*cur tkt <)tkt

Overview

A Linux tracing framework

A pipeline: Analyzer → Tracer → Viewer

Captures blocking and busy-waiting dependencies
Combines hardware watchpoints with
software breakpoints



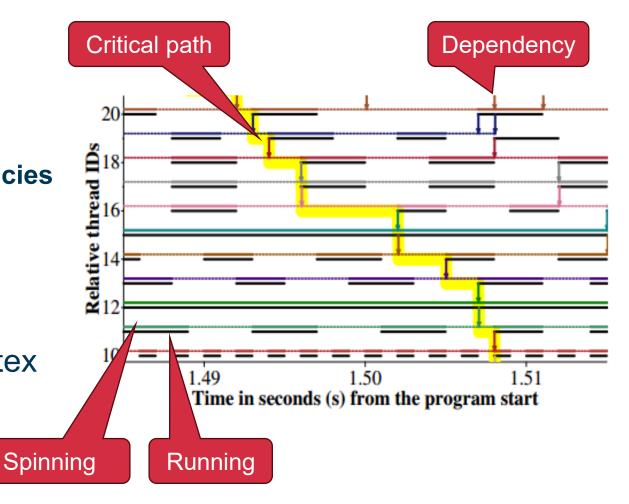
Overview

A Linux tracing framework

A pipeline: Analyzer → Tracer → Viewer

Captures blocking and busy-waiting dependencies
Combines hardware watchpoints with
software breakpoints

Blocking captured by trace-cmd via futex syscalls



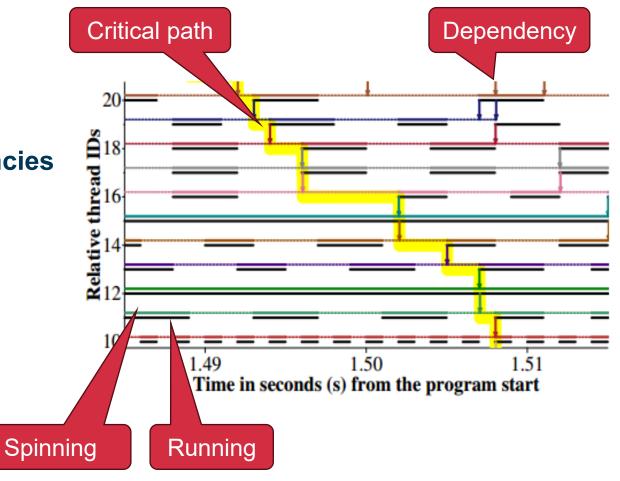
Overview

A Linux tracing framework

A pipeline: Analyzer → Tracer → Viewer

Captures blocking and busy-waiting dependencies
Combines hardware watchpoints with
software breakpoints

How to capture **busy-waiting**?



Architecture: Analyzer → Tracer → Viewer

Finding all busy-wait loops in binaries

Heuristic similar to one proposed by Jannesari & Tichy, IPDPS 2010

- The loop exit condition depends on memory load
- The conditions' value is not changed by the loop body

```
/* Barrier */
if (atomic inc(thr count) < n thr)</pre>
 while (*barrier end != 0) PAUSE;
else *barrier end = 1;
/* Simple lock */
while (atomic test and set(
    lock, 0, 1) != 0)
  PAUSE;
/* Ticket lock */
tkt = atomic inc(max tkt);
while (*cur tkt < tkt) PAUSE;
```

Architecture: Analyzer → **Tracer** → Viewer

Capturing busy-waiting synchronization

- a) Capture begin and end of busy-wait loop
- b) Capture dependency

a) Capture begin and end of busy-wait loop

Statically patch the beginning and the end of each busy-wait loop to store a tracepoint

```
while (*barrier_end != 0) PAUSE;
spinloop:
  pause
 mov (%rax), %rbx
  cmp %rbx, %rcx
  jne spinloop
```

a) Capture begin and end of busy-wait loop

Statically patch the beginning and the end of each busy-wait loop to store a tracepoint

```
store_time_begin(spinId);
while (*barrier_end != 0) PAUSE;
store time end(spinId);
call store_time_begin
spinloop:
  pause
  mov (%rax), %rbx
  cmp %rbx, %rcx
  jne spinloop
call store_time_end
```

b) Capturing dependency

Spin-variable (dependency endpoint) is found using static analysis

How to find **dependency store**?

```
if (atomic_inc(thr_count) < n_thr)</pre>
  while (*barrier_end != 0) PAUSE;
else
  *barrier end = 1;
spinloop:
                   Spin-variable
  pause
      (%rax), %rbx
  cmp %rbx, %rcx
  jne spinloop
mov %rdx, (%rdi)
```

b) Capturing dependency

```
if (atomic_inc(thr_count) < n_thr)
  while (*barrier_end != 0) PAUSE;
else
  *barrier_end = 1;</pre>
```

Solution: Use write watchpoints (as in debuggers)

How to find dependency store?

```
cmp %rbx, %rcx
jne spinloop
...
mov %rdx, (%rdi)
```

b) Capturing dependency

- (i) Statically inject code that sets a write watchpoint on the dependency variable address
- (ii) When the watchpoint is triggered:
- Store the dependency

```
if (atomic_inc(thr_count) < n_thr)</pre>
  while (*barrier_end != 0) PAUSE;
else
  *barrier end = 1;
call set_watchpoint(%rax)
spinloop:
                   Spin-variable
  pause
  mov (%rax), %rbx
  cmp %rbx, %rcx
  jne spinloop
mov %rdx, (%rdi)
```

b) Capturing dependency

(i) Statically inject code that sets a write

watchp

addres

```
if (atomic_inc(thr_count) < n_thr)
  while (*barrier_end != 0) PAUSE;
else
  *barrier_end = 1;</pre>
```

Problem: limited number of watchpoints *E.g., four on x86*

```
cmp %rbx, %rcx
jne spinloop
...
mov %rdx, (%rdi)
```

b) Capturing dependency

(i) Statically inject code that sets a write

watchp

addres

```
if (atomic_inc(thr_count) < n_thr)
  while (*barrier_end != 0) PAUSE;
else
  *barrier_end = 1;</pre>
```

Solution: dynamically replace watchpoints with software breakpoints

```
cmp %rbx, %rcx
jne spinloop
...
mov %rdx, (%rdi)
```

b) Capturing dependency

(i) Statically inject code that sets a write

watchp

addres

```
if (atomic_inc(thr_count) < n_thr)
  while (*barrier_end != 0) PAUSE;
else
  *barrier_end = 1;</pre>
```

Solution: dynamically replace watchpoints with software breakpoints

```
cmp %rbx, %r
jne spinloop

code addresses

mov %rdx, (%rdi)
```

b) Capturing dependency

- (i) Statically inject code that sets a write watchpoint on the dependency variable address
- (ii) When the watchpoint is **triggered**:
- 1. Store the dependency
- 2. Disable the watchpoint
- 3. Set a breakpoint on the store code address (%rip) for future dependencies

```
on_trigger(watch_id):
    store_dep(<%rdx,==,0>);
    disable_watch(watch_id);
    set_break(%rip,&store_dep);
```

```
if (atomic_inc(thr_count) < n_thr)</pre>
  while (*barrier_end != 0) PAUSE;
else
  *barrier_end = 1;
call set_watchpoint(%rax)
spinloop:
  pause
  mov (%rax),
              %rbx
  cmp %rbx, %rcx
  jne spinloop
mov %rdx, (%rdi)
```

Architecture: Analyzer → Tracer → **Viewer**

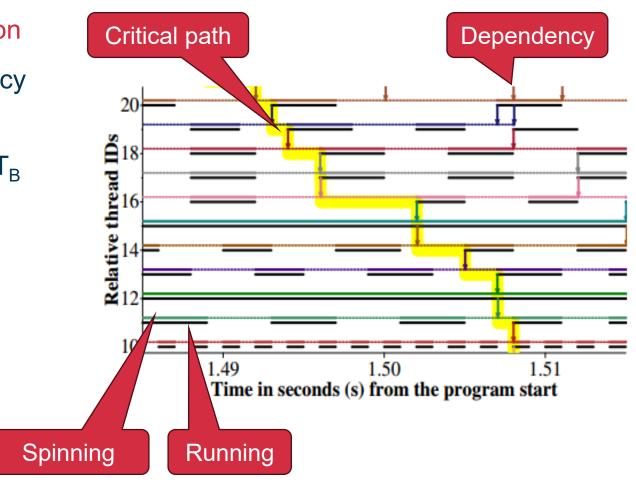
Post-processing of raw traces and visualization

Connects events with the same dependency

 $T_A \rightarrow T_B$: A thread T_A resolves a dependency T_B was waiting on

Forms a dependency graph

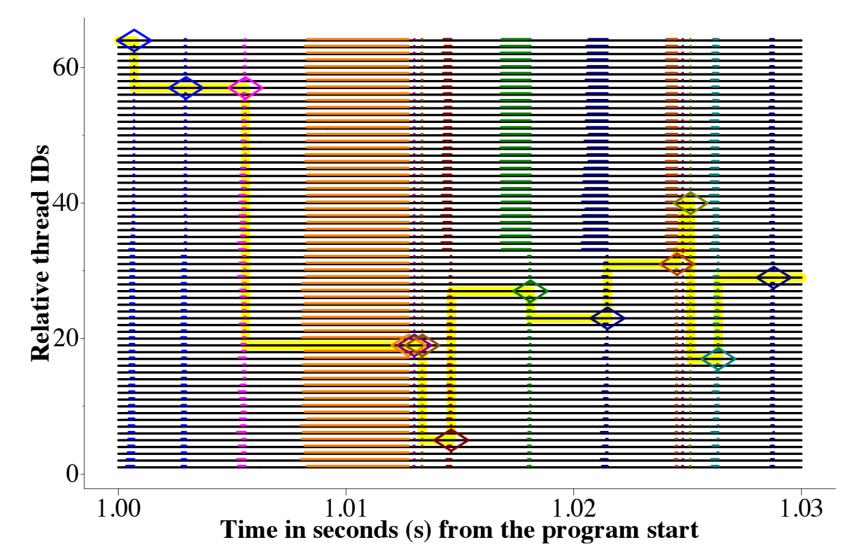
Critical path is the longest path in the dependency graph



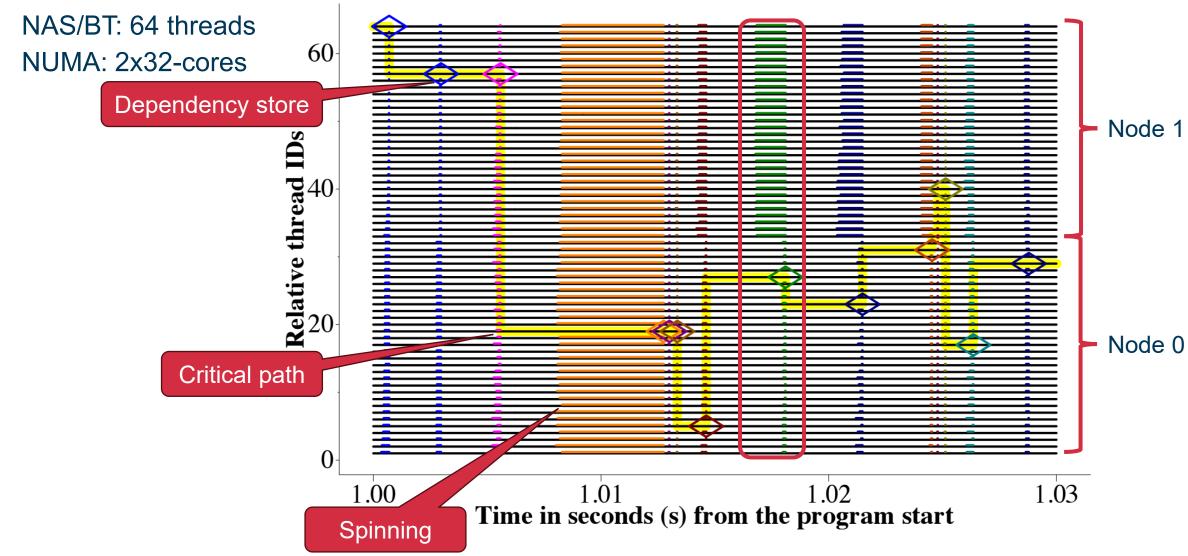
Observing NUMA Effect of Spinning Barrier

NAS/BT: 64 threads

NUMA: 2x32-cores



Observing NUMA Effect of Spinning Barrier



Diagnosing Pathological Busy-Waiting

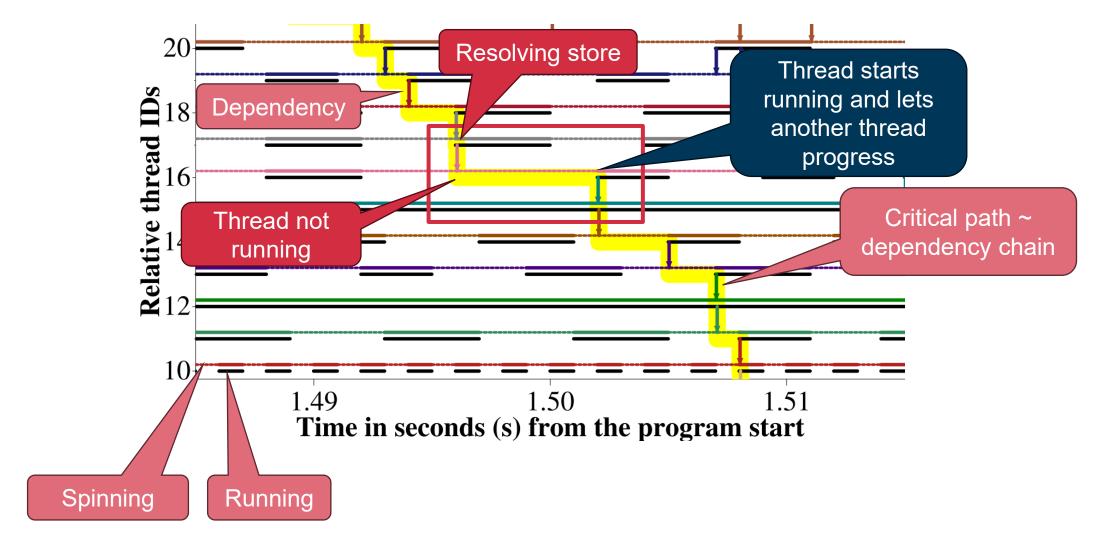
Running two NAS/LU in parallel OpenMP set to **blocking**

Why is the parallel execution of LU so **slow**?

Bench.	Single	Parallel
BT	2.67	5.72
CG	0.09	0.20
EP	0.25	0.52
FT	0.18	0.35
IS	0.02	0.05
LU	2.22	100.09
MG	0.14	0.26
SP	1.58	3.86
UA	3.58	7.23

(a) Single vs. Parallel.

Diagnosing Pathological Busy-Waiting



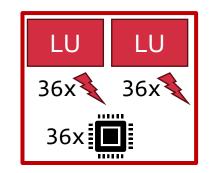
Running two NAS/LU in parallel:

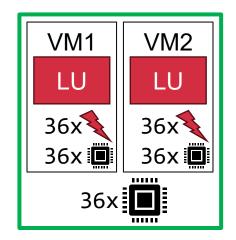
- BM: instances run on bare metal.
- VM: each instance runs in a separate VM

OpenMP set to default

 Busy-waiting for a fixed number of iterations, then blocking

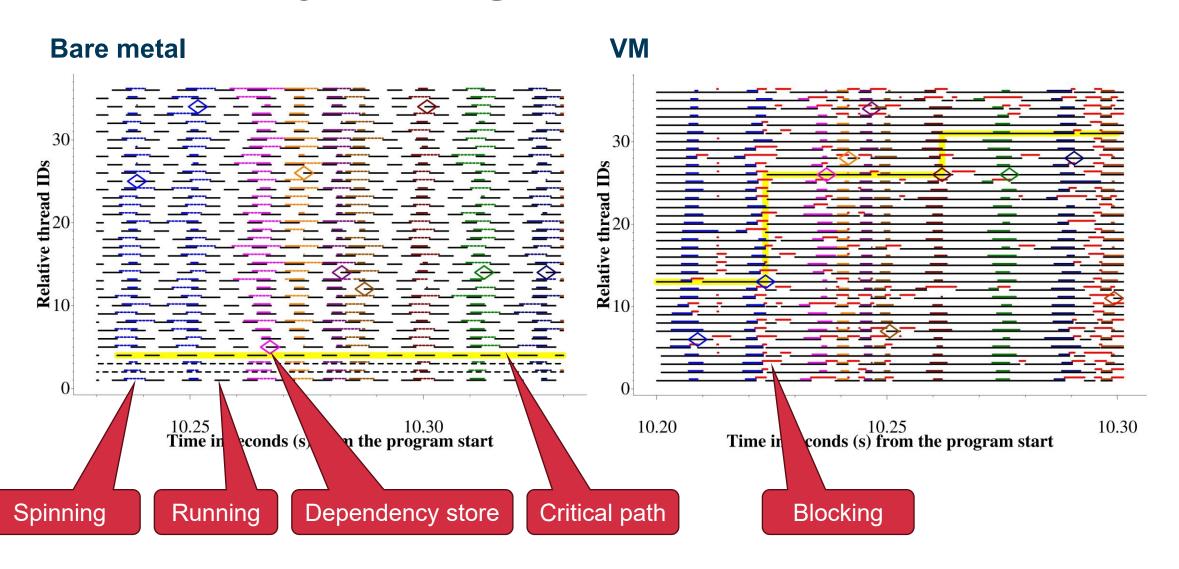
Why is **running in VM faster** than on bare metal?

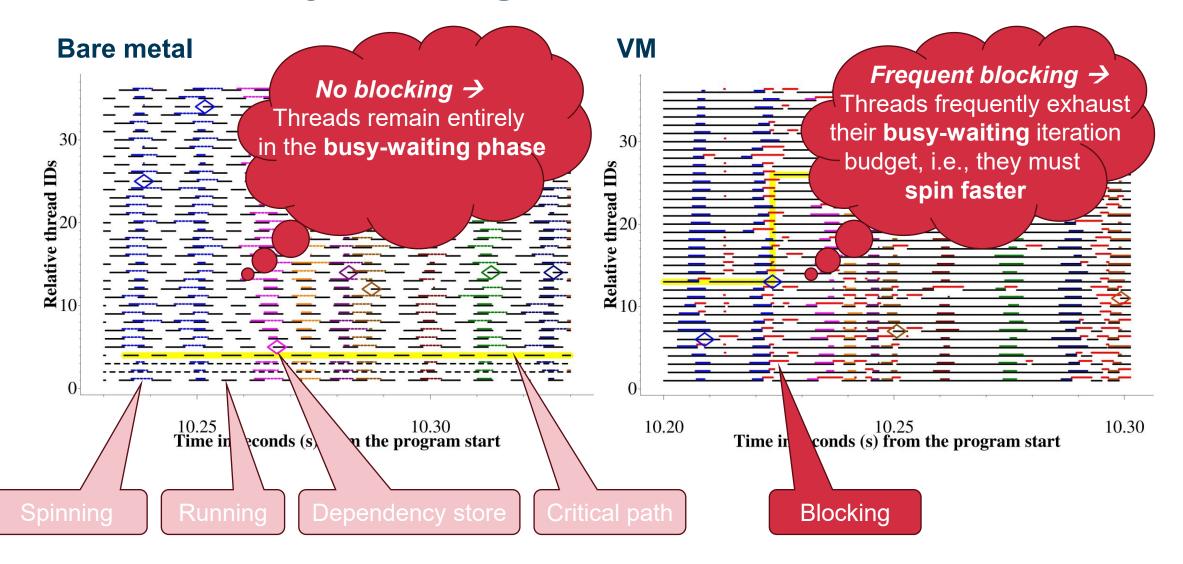


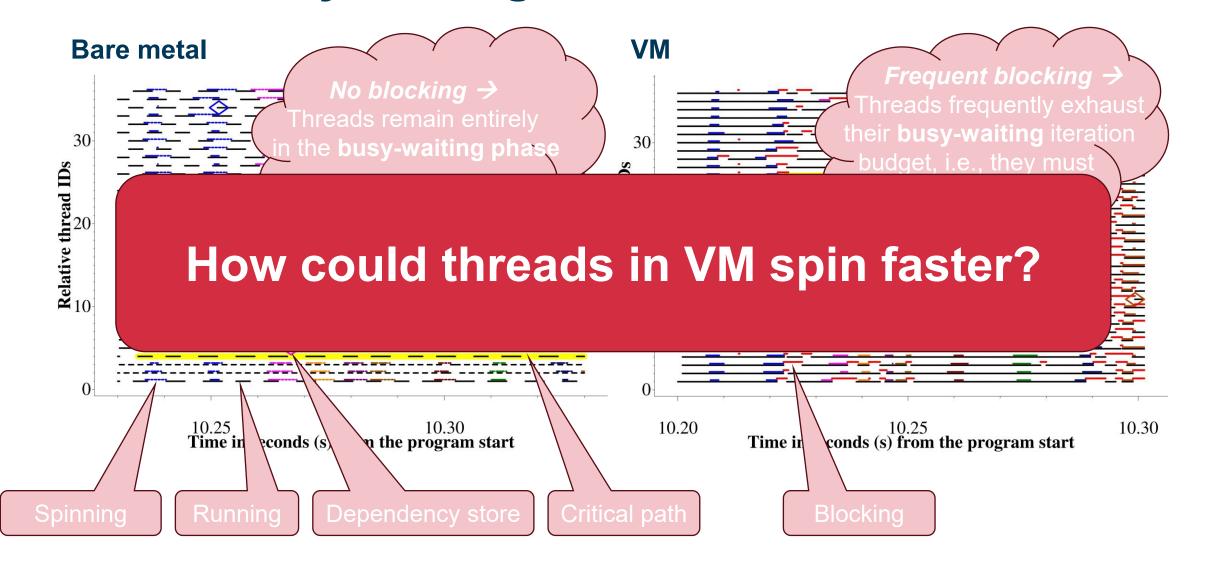


Bench.	BM	VM		
BT	13.22	11.94		
CG	8.44	5.02		
EP	0.53	0.45		
FT	0.57	0.54		
IS	0.23	0.16		
LU	190.75	138.44		
MG	2.51	1.6		
SP	31.9	22.57		
UA	300.2	185.7		

(b) Bare Metal (BM) vs. VM.







Pause-Loop Exiting (PLE)

A hardware feature of Intel VMX

PLE triggers VM exits for spinloops to help the hypervisor mitigate guest-kernel lock holder/waiter preemption issue

"PLE VM-execution control is **ignored in userspace**" [Intel manual, §36.7.3]

```
// userland test:
for (i=0; i < 109; i++) PAUSE;</pre>
```

Config.	Rt (s)	f (GHz)	10 ⁹ Instr.		10 ⁹ PAUSEs	
	200 (0)		H	\boldsymbol{G}	H	G
H	14.12	3.89	5.0	0	1.0	0
G PLE	6.19	3.89	0	5.0	0	0
G no PLE	14.17	3.89	0	5.0	0	1.0

Table 2. Impact of PLE on a userland busy-wait loop with 10^9 iterations (Rt=runtime, f=CPU frequency, H=host, G=guest).

Pause-Loop Exiting (PLE)

A hardware feature of Intel VMX

PLE triggers VM exits for spinloops to help the hypervisor mitigate guest-kernel lock holder/waiter preemption issue

"PLE VM-execution control is **ignored in** userspace" [Intel manual, §36.7.3]

And yet, with PLE enabled, user PAUSE instructions are not even retired on Skylake and Cascade Lake machines

```
// userland test:
for (i=0; i < 109; i++) PAUSE;</pre>
```

Config.		Rt (s) f (GHz)	10 ⁹ Instr.		10 ⁹ PAUSEs		
			<i>j</i> (3222)	H	\boldsymbol{G}	H	\boldsymbol{G}
H		14.12	3.89	5.0	0	1.0	0
\boldsymbol{G}	PLE	6.19	3.89	0	5.0	0	0
\boldsymbol{G}	no PLE	14.17	3.89	0	5.0	0	1.0

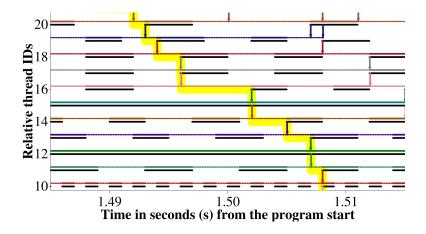
Table 2. Impact of PLE on a userland busy-wait loop with 10^9 iterations (Rt=runtime, f=CPU frequency, H=host, G=guest).

A **dependency model** that unifies the concept of inter-thread wait-for dependencies for blocking and busy-waiting synchronization

A **dependency model** that unifies the concept of inter-thread wait-for dependencies for blocking and busy-waiting synchronization

Tapestry: A **tracing framework** to observe wait-for dependencies

Combining the use of hardware watchpoints and software breakpoints



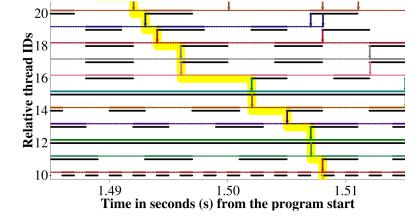
A **dependency model** that unifies the concept of inter-thread wait-for dependencies for blocking and busy-waiting synchronization

Tapestry: A tracing framework to observe wait-for dependencies

Combining the use of hardware watchpoints and software breakpoints

Investigated three anomalies using Tapestry

- NUMA effect on spin-barriers
- Pathological busy-waiting
- Faster busy-waiting in VM vs. bare metal



Found undocumented hardware effect on some Intel machines

A **dependency model** that unifies the concept of inter-thread wait-for dependencies for blocking and busy-waiting synchronization

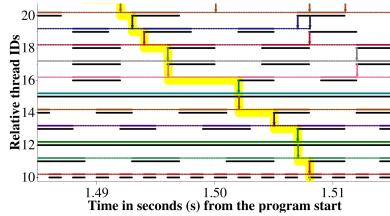
Tapestry: A tracing framework to observe wait-for dependencies

Combining the use of hardware watchpoints and software breakpoints

Investigated three anomalies using Tapestry

- NUMA effect on spin-barriers
- Pathological busy-waiting
- Faster busy-waiting in VM vs. bare metal





Thank you!