

# Living on the Edge : Rapid Probe Toggling on x86

**Buddhika Chamith**, Bo Joel Svensson, Luke Dalessandro  
Ryan R. Newton





How to modify a memory  
word atomically?

## ***A data word***

... “32” “53” “**42**” “41”“85” ...

## ***A data word***

... “32” “53” “**43**” “41”“85” ...

# Compare and Swap (CAS)?

Asserts ***Lock*** signal on memory bus

**A data word**

*Cache line boundary*



**A data word**

*Cache line boundary*



... "32" "53" 4

3 "41" "85" ...

How to modify **an instruction**  
word atomically?

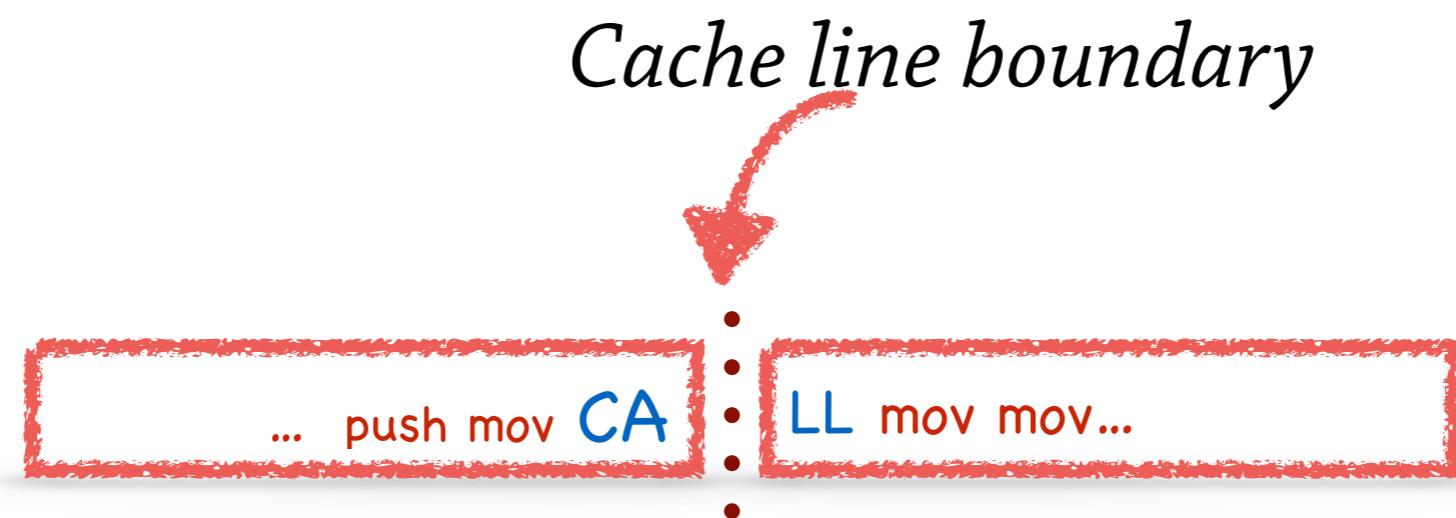
## ***An instruction word***

... push mov CALL ... push mov

## *An instruction word*

... push mov NOOP... push mov

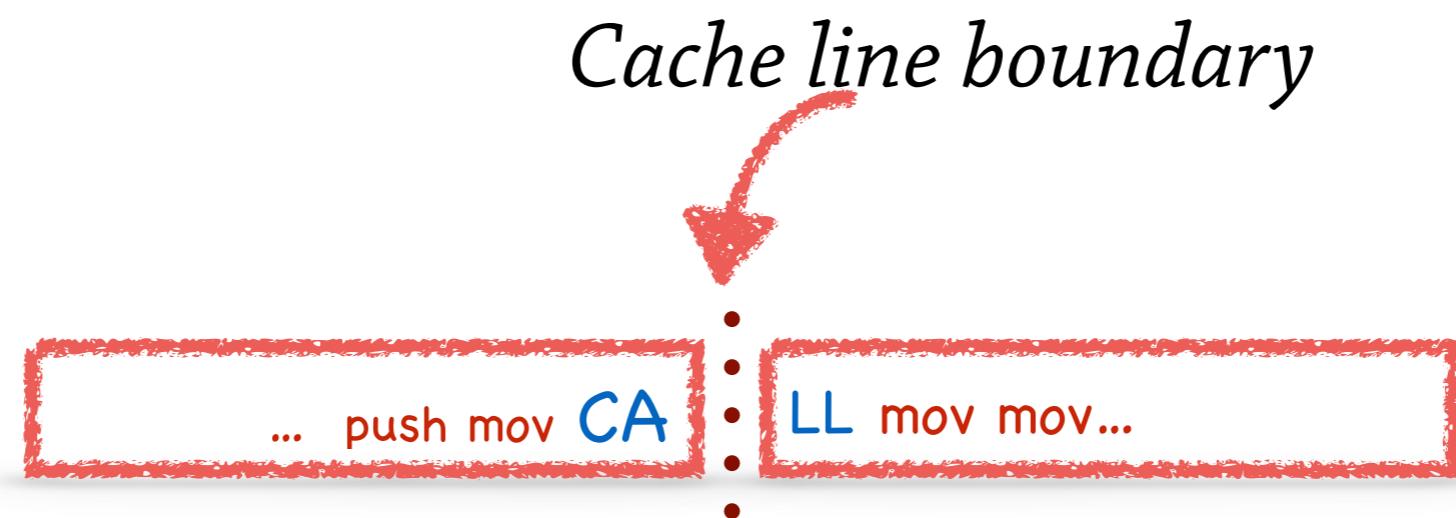
# Woes of CAS



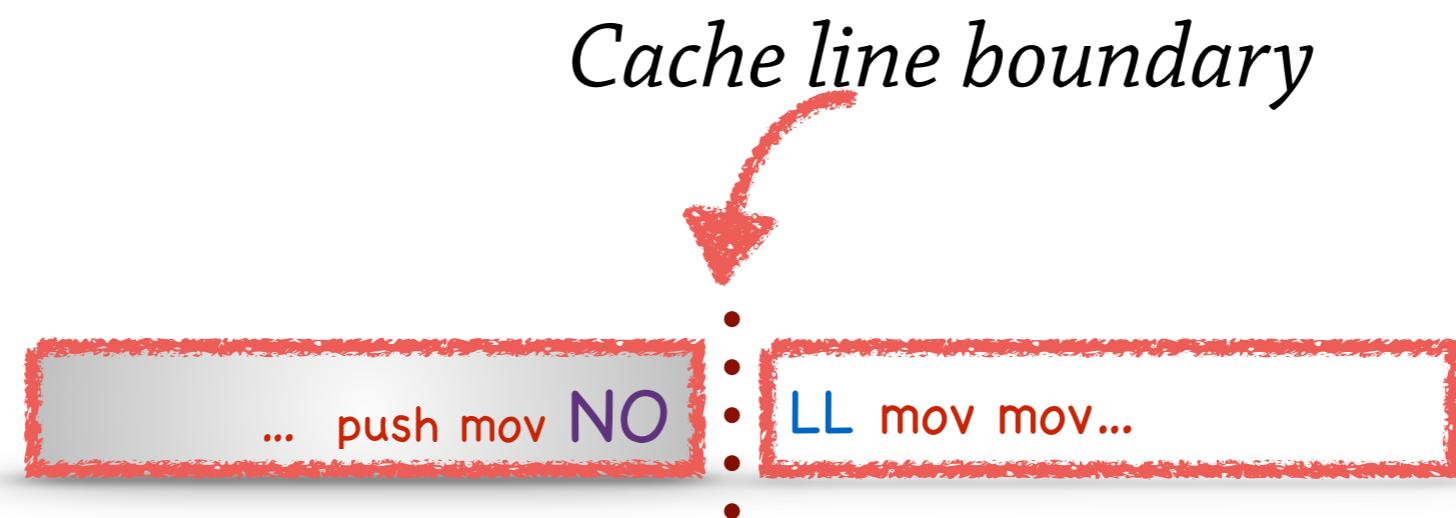
# Woes of CAS

***Locked operations*** are not atomic with respect to ***instruction fetch*** when they operate on multiple cache-lines.

# Woes of CAS



# Woes of CAS



**INCORRECT**

# Cross Modification

***Definition*** : The act of one processor writing data into the currently executing code segment of a second processor.

So how can we do **CROSS**  
**modification** safely?



We wanted to build a **lightweight profiler**.

# Lightweight Profiler

- Instrumentation based

*foo: <call instrument\_fn>*

*push*

*mov*

*...*

*<call instrument\_fn>*

*ret*

# Lightweight Profiler

- Periodic Probing

*foo: <call instrument\_fn>*

*push*

*mov*

...

*<call instrument\_fn>*

*ret*

# Lightweight Profiler

- Periodic Probing

*foo:* ***noop***

*push*

*mov*

...

***noop***

*ret*

**Cheap**

# Requirements

*fn1:*

***call***

....

***call***

*ret*

*fn2:* ***call***

...

***call***

*ret*

...

*fn3: call*

...

***call***

*ret*

*fn4: call*

...

***call***

*ret*

...

*fn5: call*

...

***call***

*ret*

...

*fnN: call*

...

***call***

*ret*

Cheap

# Requirements

*fn1:*

*call*

....

*call*

*ret*

*fn2:* ***noop***

...

***noop***

*ret*

...

*fn3: call*

...

*call*

*ret*

*fn4: noop*

...

***noop***

*ret*

...

*fn5: call*

...

*call*

*ret*

...

*fnN: noop*

...

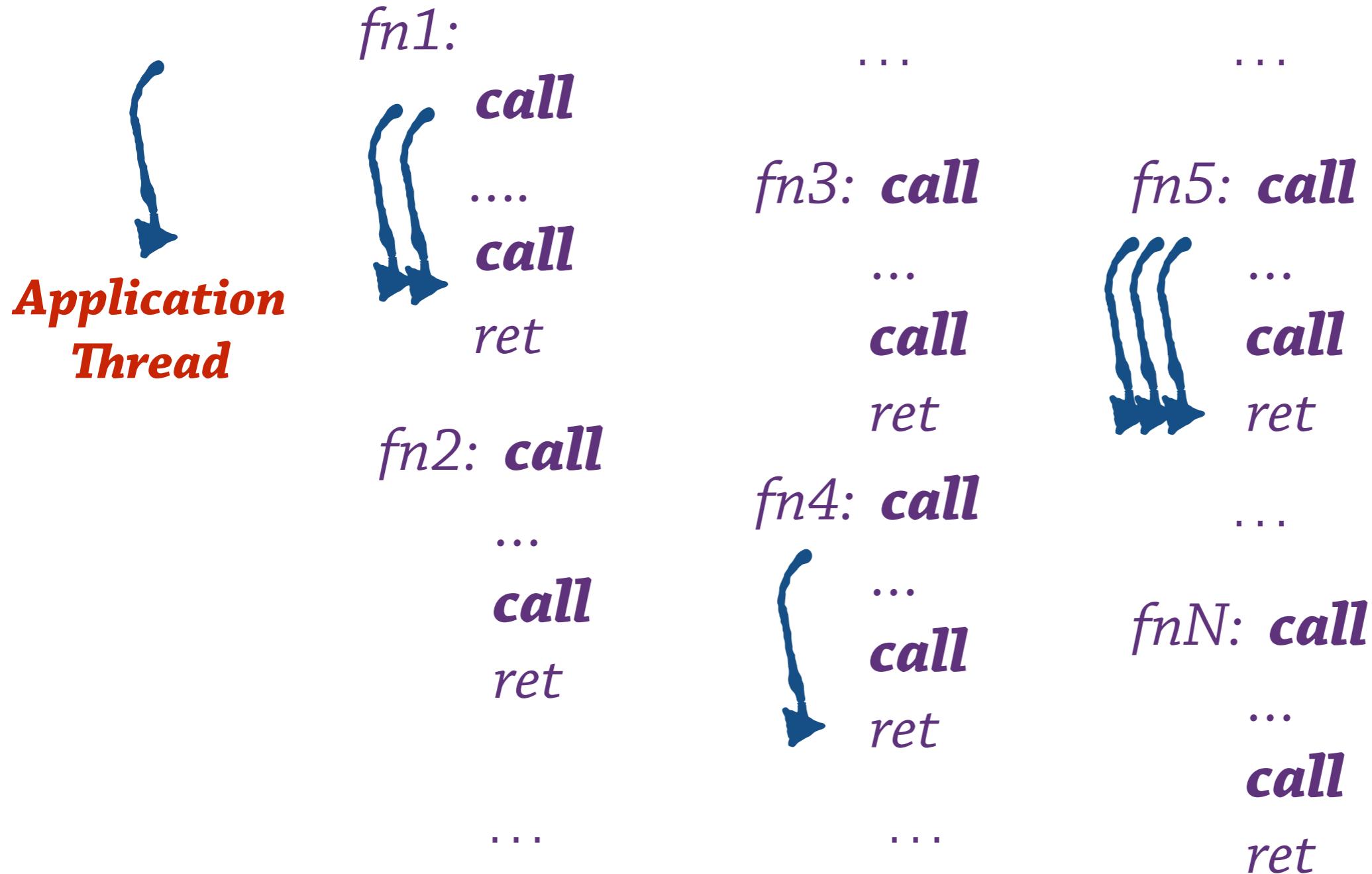
***noop***

*ret*

# Requirements

Cheap

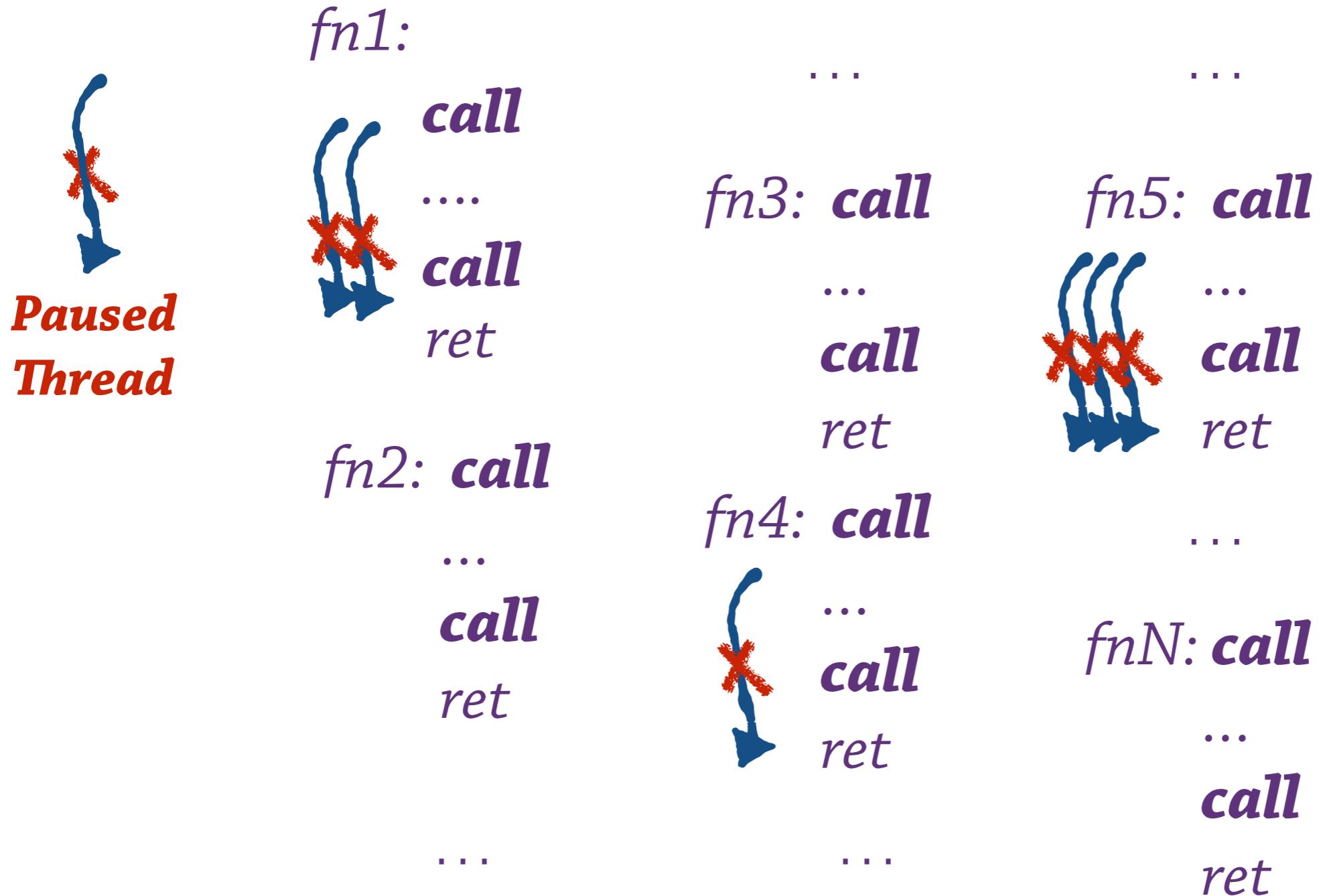
## Thread Safe



# Requirements

Cheap

## Thread Safe

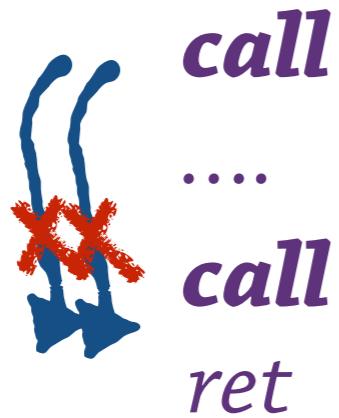


# Requirements

Cheap

**Thread Safe**

*fn1:*



*fn2: call*



...

*fn3: call*



*fn4: call*



...

*fn5: call*



*fnN: call*

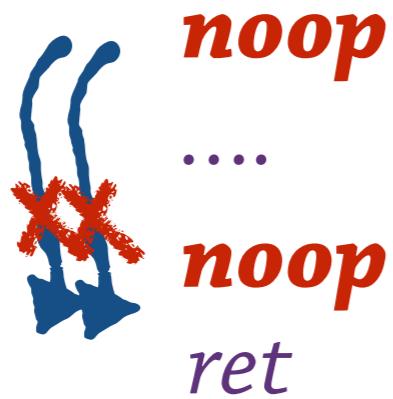


# Requirements

Cheap

**Thread Safe**

*fn1:*



*fn2: call*

...  
**call**  
**ret**

...

*fn3: call*

...  
**call**  
**ret**

*fn4: call*

Diagram illustrating the assembly code for *fn4*. It shows a sequence of instructions: a **call** instruction, followed by a **ret** instruction. The **call** instruction is highlighted with red scribbles, indicating it is cheap. The **ret** instruction is also highlighted with blue scribbles.

...

*fn5: call*

Diagram illustrating the assembly code for *fn5*. It shows a sequence of instructions: three **call** instructions, followed by a **ret** instruction. The first **call** instruction is highlighted with red scribbles, indicating it is cheap. The **ret** instruction is also highlighted with blue scribbles.

...

*fnN: call*

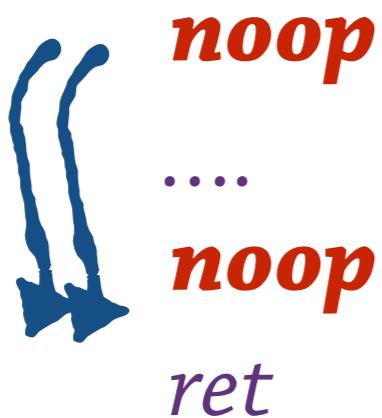
...  
**call**  
**ret**

# Requirements

Cheap

**Thread Safe**

*fn1:*



*fn2: call*

...  
*call*  
*ret*

...

*fn3: call*

...  
*call*  
*ret*

*fn4: call*

Diagram illustrating the execution flow of *fn4*. A single blue arrow points from the start of the function to a purple ***call*** label.

...

*fn5: call*

Diagram illustrating the execution flow of *fn5*. Three parallel blue arrows point from the start of the function to a purple ***call*** label.

...  
*call*  
*ret*

*fnN: call*

...  
*call*  
*ret*

# Requirements

Cheap

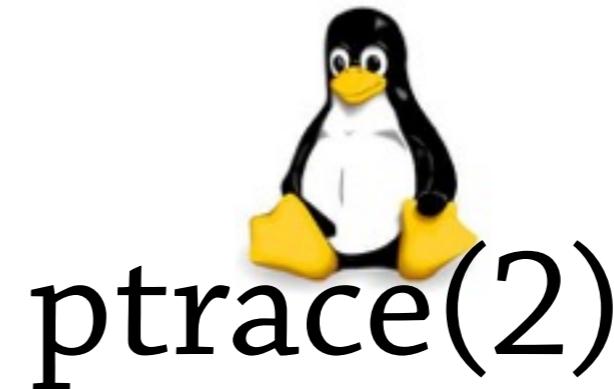
Thread Safe

**Scalable**



# Some Options : In Place Instrumentation

*Dyninst*



- ✗ *Cheap*
- ✓ *Thread safe*
- ✗ *Scalable*

- ✗ *Cheap*
- ✓ *Thread safe*
- ✓ *Scalable*

# Some Options : In Place Instrumentation

*Dyn  
inst*



ptrace(2)

- ✗ *Cheap*
- ✓ *Thread safe*
- ✗ *Scalable*

Dtrace



- ✓ *Cheap*
- ✓ *Thread safe*
- ✓ *Scalable*

# Some Options (Contd): Out of Place Translation

Valgrind



- ✗ *Cheap*
- ✓ *Thread safe*
- ✗ *Scalable*

Dynamorio



# Assumption

There exists an upper bound for change propagation and visibility from other cores when doing cross modification.

# $T_{max}$

There is an upper bound  $T_{max}$  where instruction cache/fetch & decode units are made eventually consistent after a write to code from one core/processor.

# Split Write Protocol (a.k.a. wordpatch)



# Split Write Protocol (a.k.a. wordpatch)



# Split Write Protocol (a.k.a. wordpatch)



A **Tmax** wait

# Split Write Protocol (a.k.a. wordpatch)



A **Tmax** wait

# Split Write Protocol (a.k.a. wordpatch)



Another **Tmax** wait.

# Split Write Protocol (a.k.a. wordpatch)



# Prerequisites

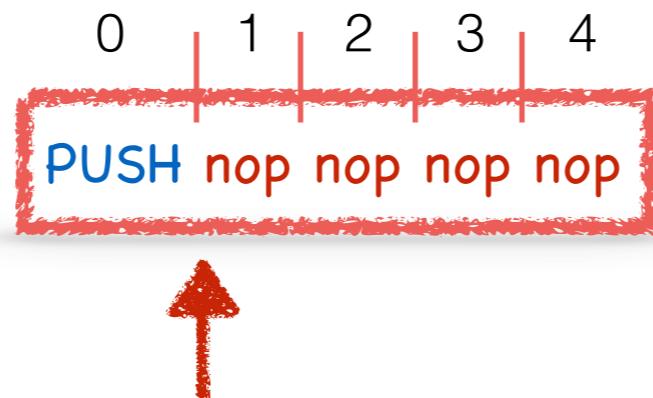
**Writable code** : Patch address must be writable.

**Layout equivalence** : Set of valid PC addresses remains invariant

**No collisions** : No addresses  $a_1$  and  $a_2$  may be concurrently modified if  $0 \leq |a_1 - a_2| \leq \text{word size}$

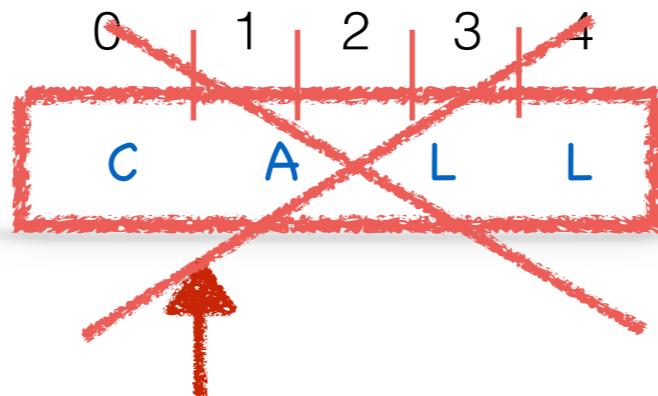
# Layout Equivalence

Layout (in bytes) = [1 | 1 | 1 | 1 | 1]



# Layout Equivalence

Layout (in bytes) = [5]



THREADn

**FAIL**

# Validating $T_{max}$

- ***A single thread toggling a call site***
- Stress test with multiple threads attempting to execute the call site.
- Many different configurations for
  - ▶ Straddle position
  - ▶ # of calling threads
- About 100 such test cases

# Validating T<sub>max</sub>

- A single thread toggling a call site
- ***Stress test with multiple threads attempting to execute the call site.***
- Many different configurations for
  - ▶ Straddle position
  - ▶ # of calling threads
- About 100 such test cases

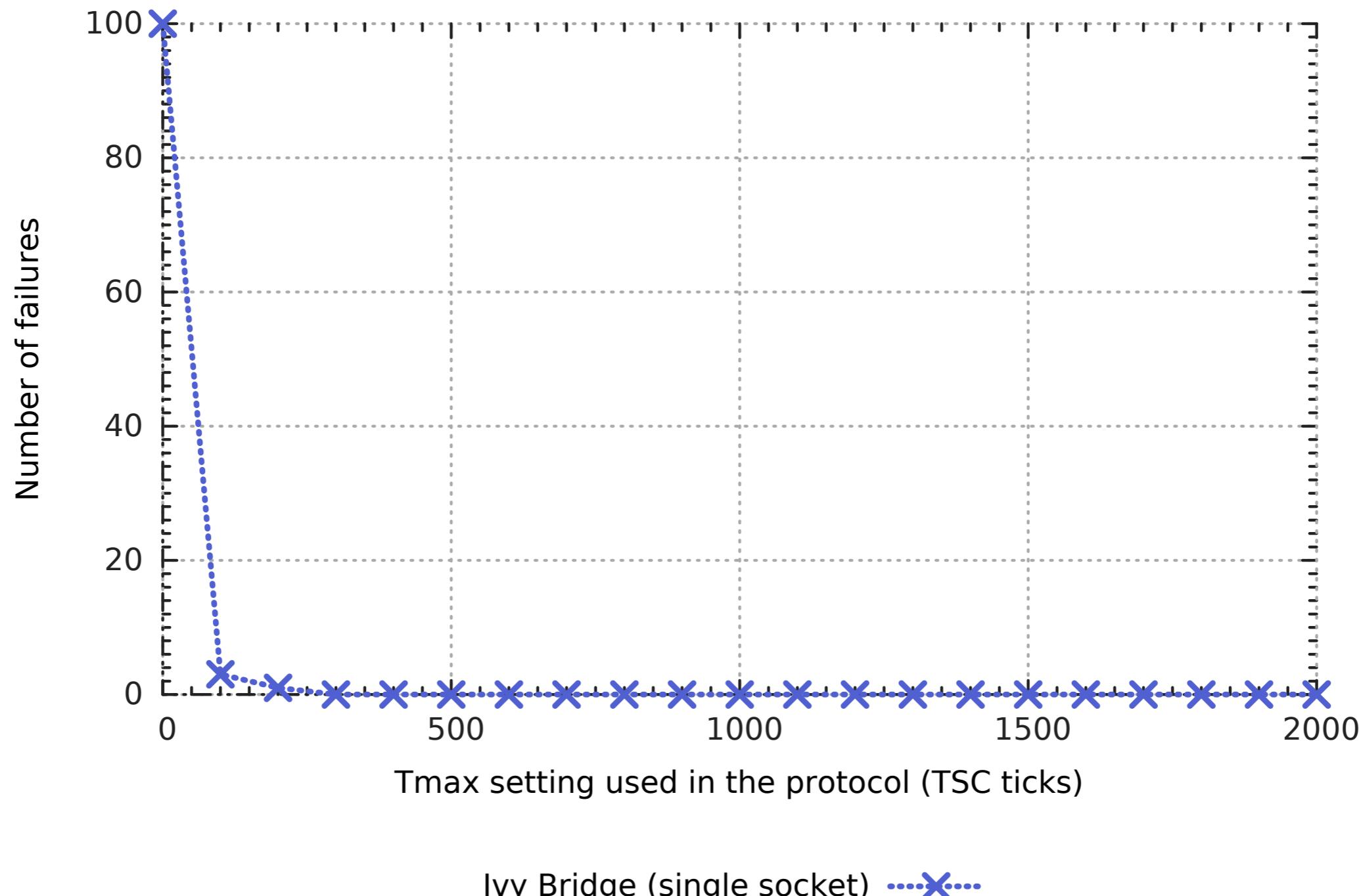
# Validating $T_{max}$

- A single thread toggling a call site
- Stress test with multiple threads attempting to execute the call site.
- ***Many different configurations for***
  - ▶ ***Straddle position***
  - ▶ ***# of calling threads***
- About 100 such test cases

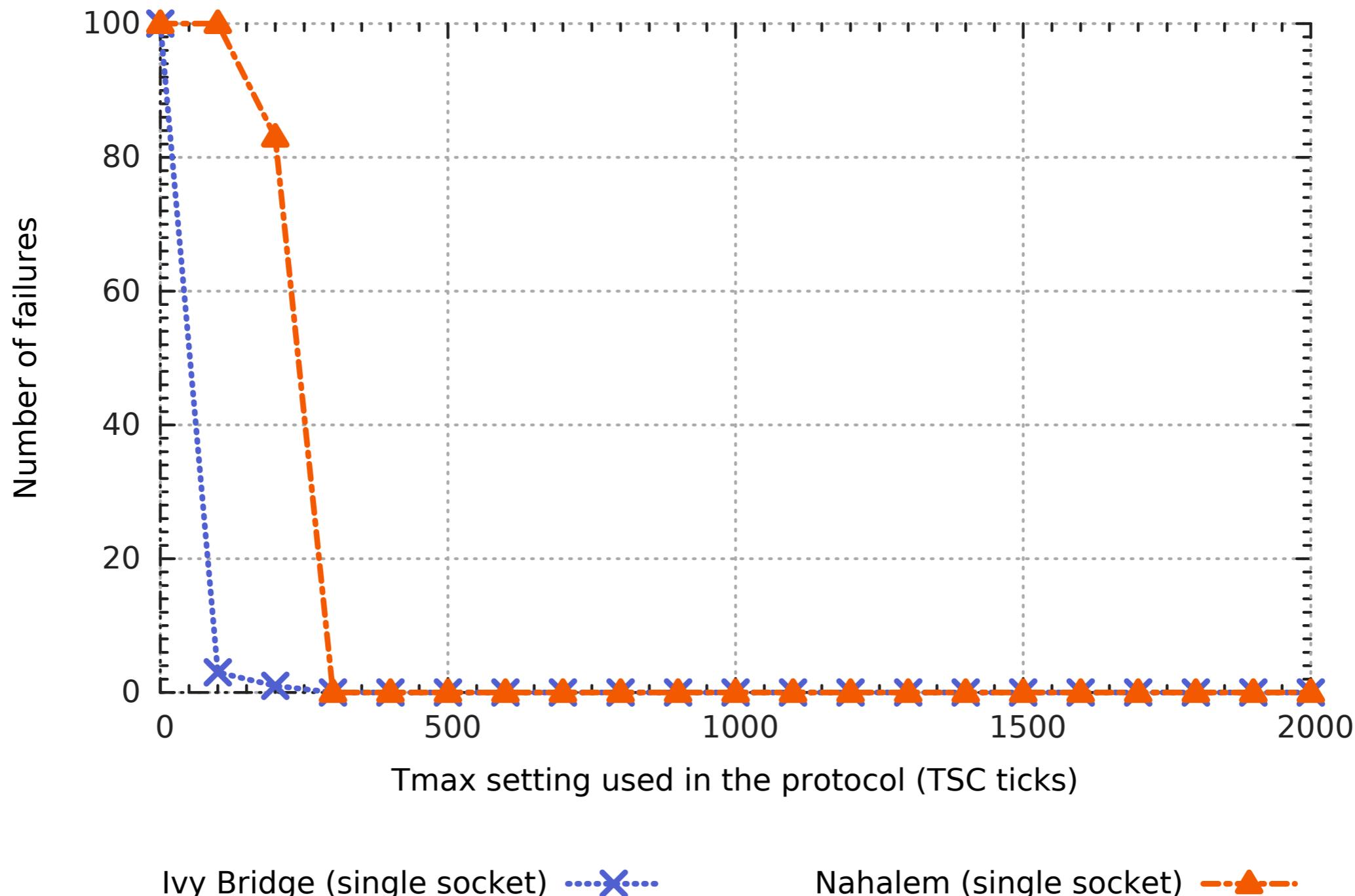
# Validating $T_{max}$

- A single thread toggling a call site
- Stress test with multiple threads attempting to execute the call site.
- Many different configurations for
  - Straddle position
  - # of calling threads
- ***About 100 such test cases***

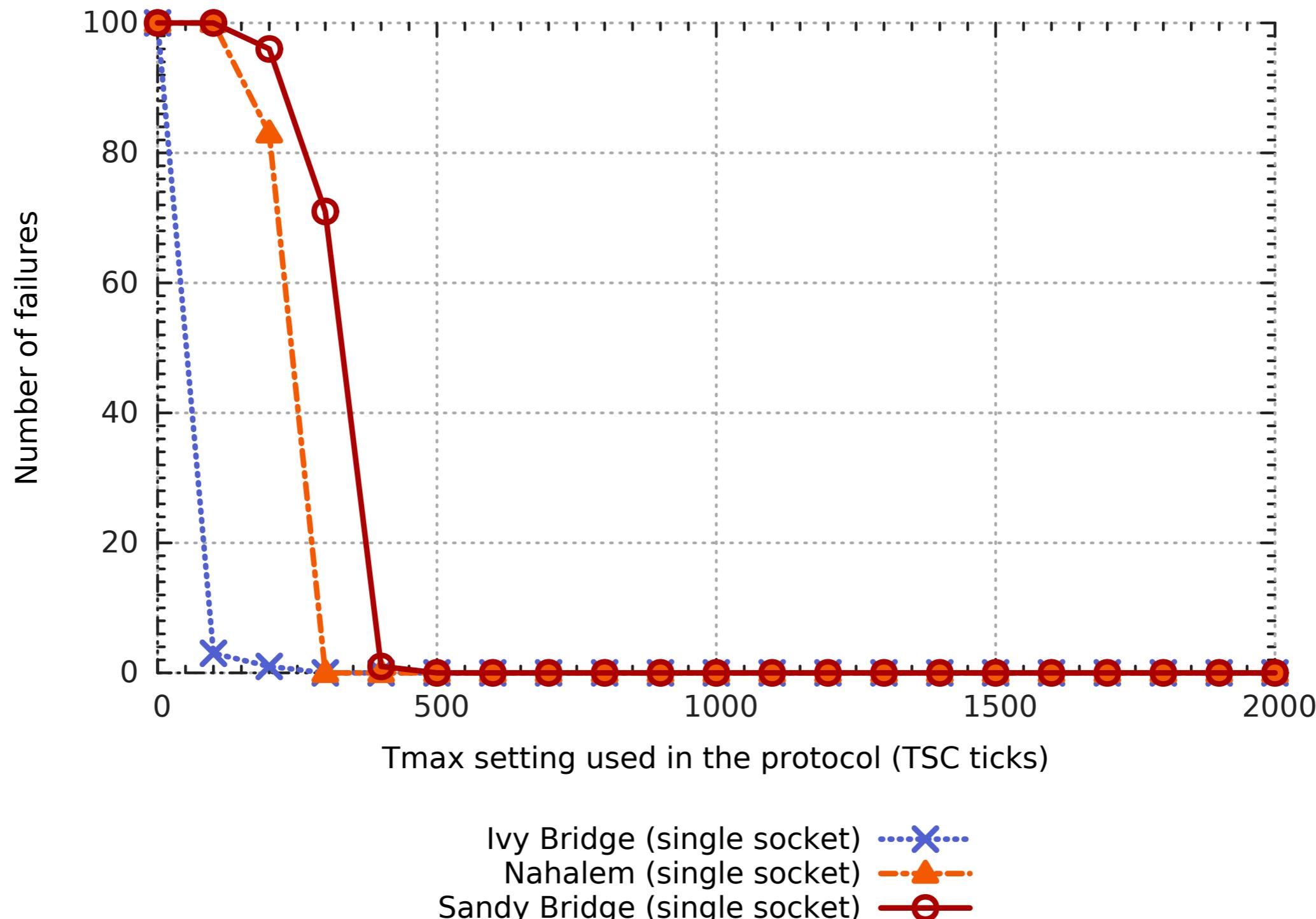
# Validating $T_{max}$



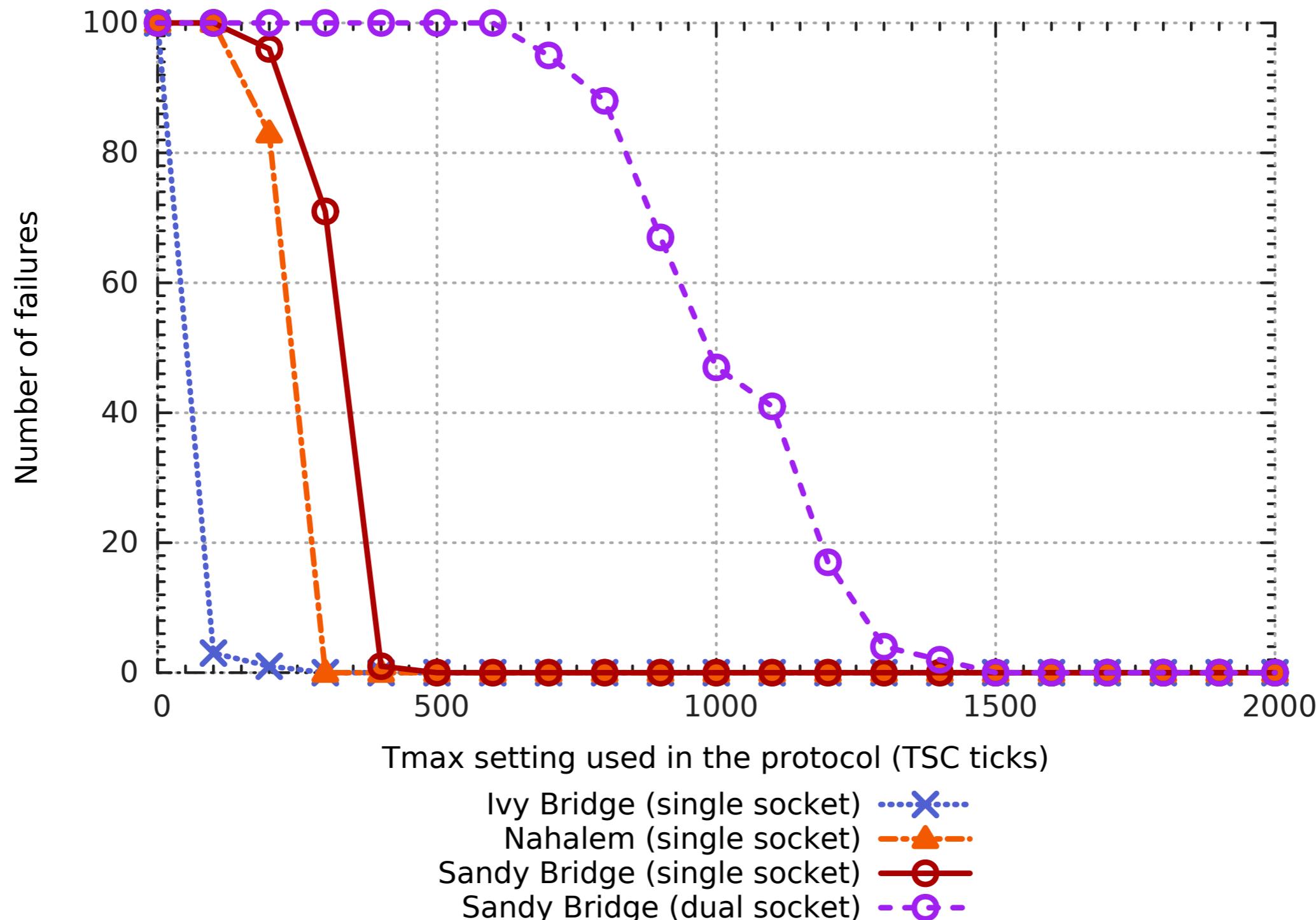
# Validating $T_{max}$



# Validating $T_{max}$



# Validating T<sub>max</sub>



# Safe *wordpatch*

```
bool patch(address, value)
    if (not is_straddler(address))
        write(address, value);
        return true;
    else if (int3_lock(address))
        wait();
        write_back(address, value);
        wait();
        write_front(address, value);
        return true;
    else return false;
```

# Safe *wordpatch*

```
// Write an interrupt byte into the address,  
void int3_lock(address)  
    return write_front(address, INT3);
```

```
bool patch(address, value)  
    if (not is_straddler(address))  
        write(address, value);  
        return true;  
    else if (int3_lock(address))  
        wait();  
        write_back(address, value);  
        wait();  
        write_front(address, value);  
        return true;  
    else return false;
```

# Safe *wordpatch*

```
// Write an interrupt byte into the address,
void int3_lock(address)
    return write_front(address, INT3);

bool patch(address, value)
    if (not is_straddler(address))
        write(address, value);
        return true;
    else if (int3_lock(address))
        wait();
        write_back(address, value);
        wait();
        write_front(address, value);
        return true;
    else return false;
```

***No global thread quiescence***

# Call Toggling (*callpatch*)

If only CALLs are patched

**no Tmax/ Wait necessary!**

Consider each straddling position **separately**  
and **patch only one half**

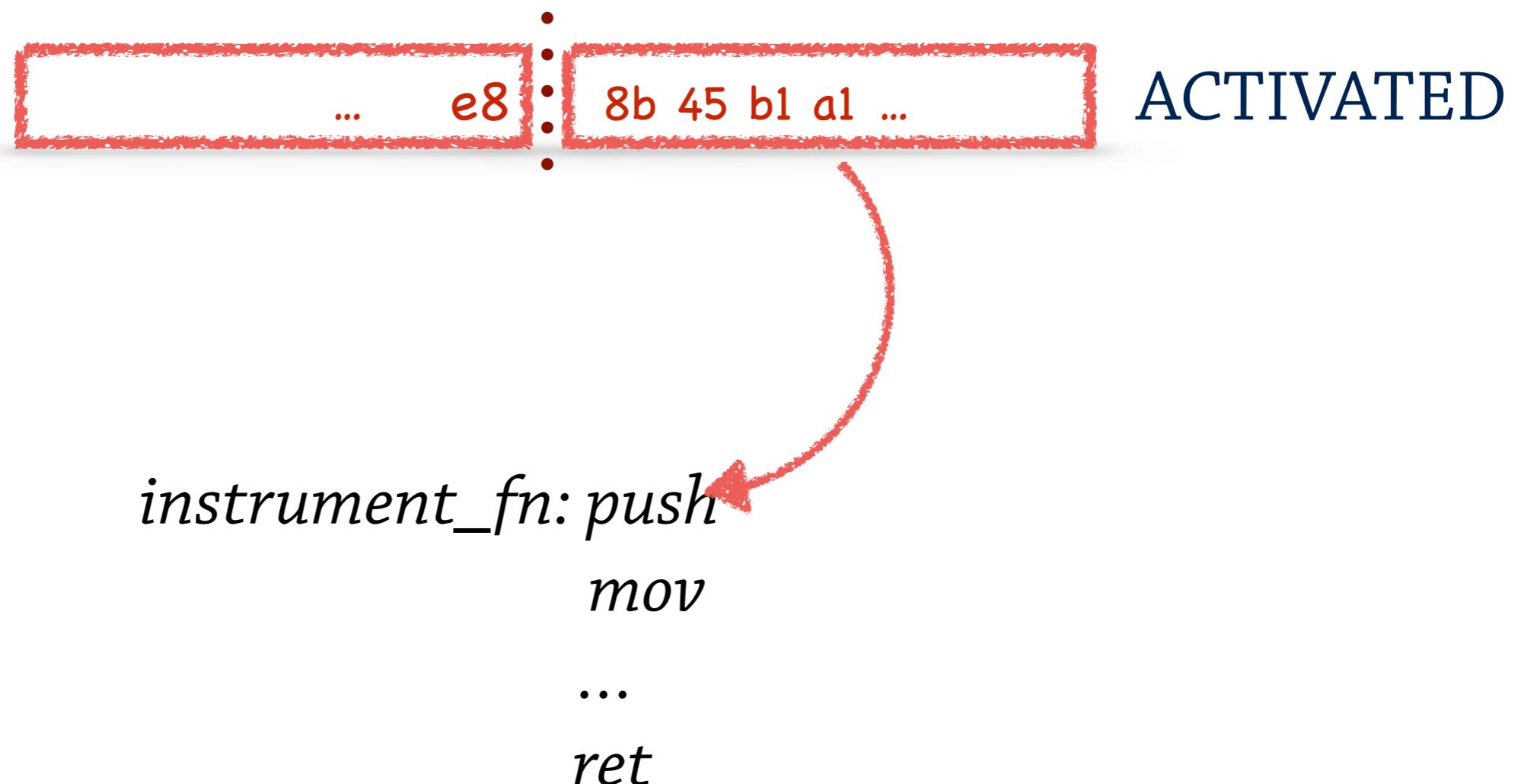
# CALL Anatomy

Op-code | Operands

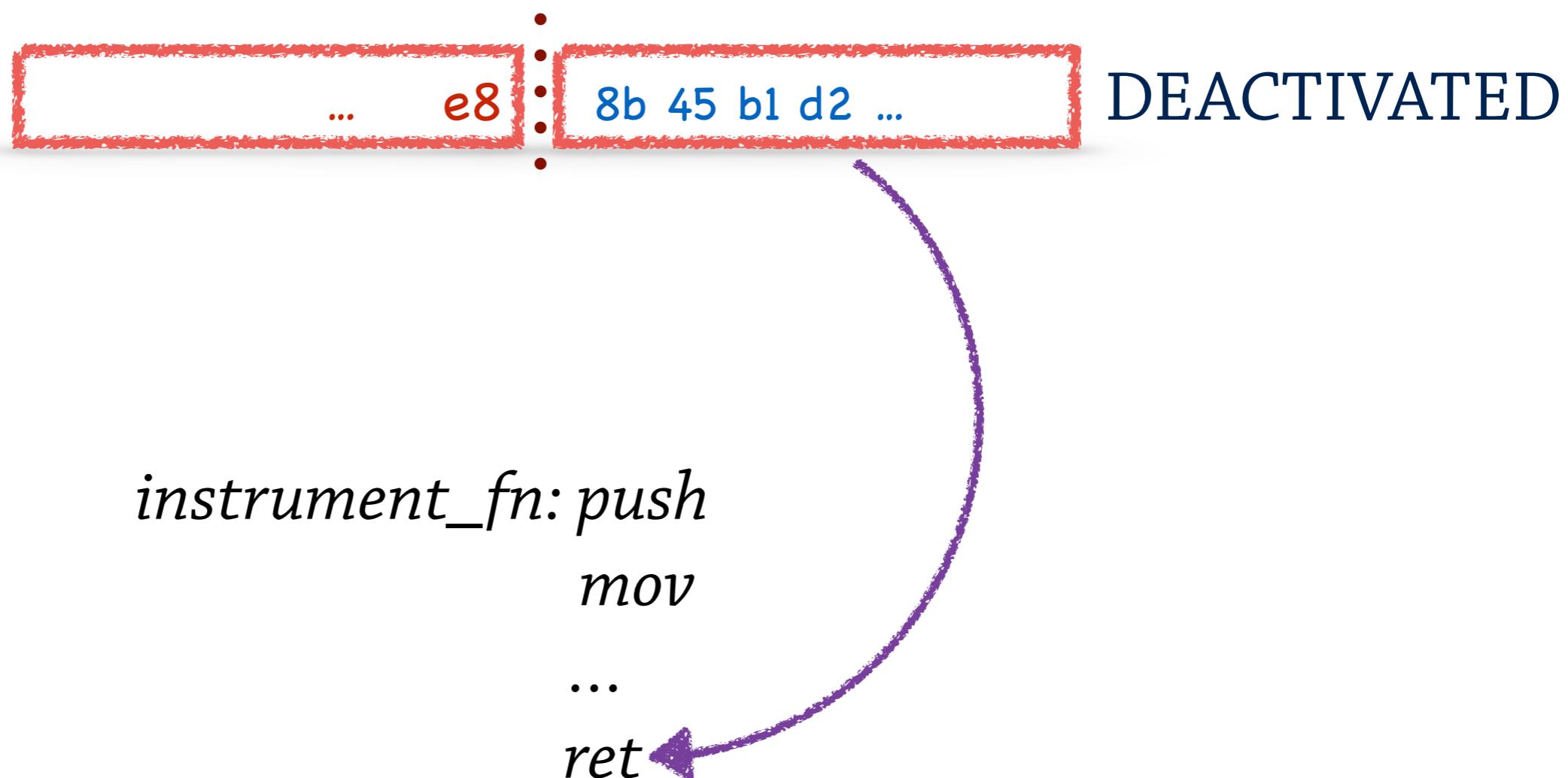
# CALL Anatomy



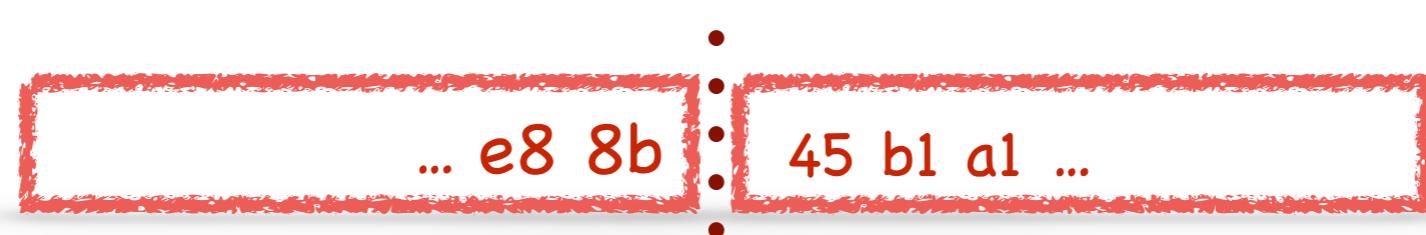
# [1 | 4] split



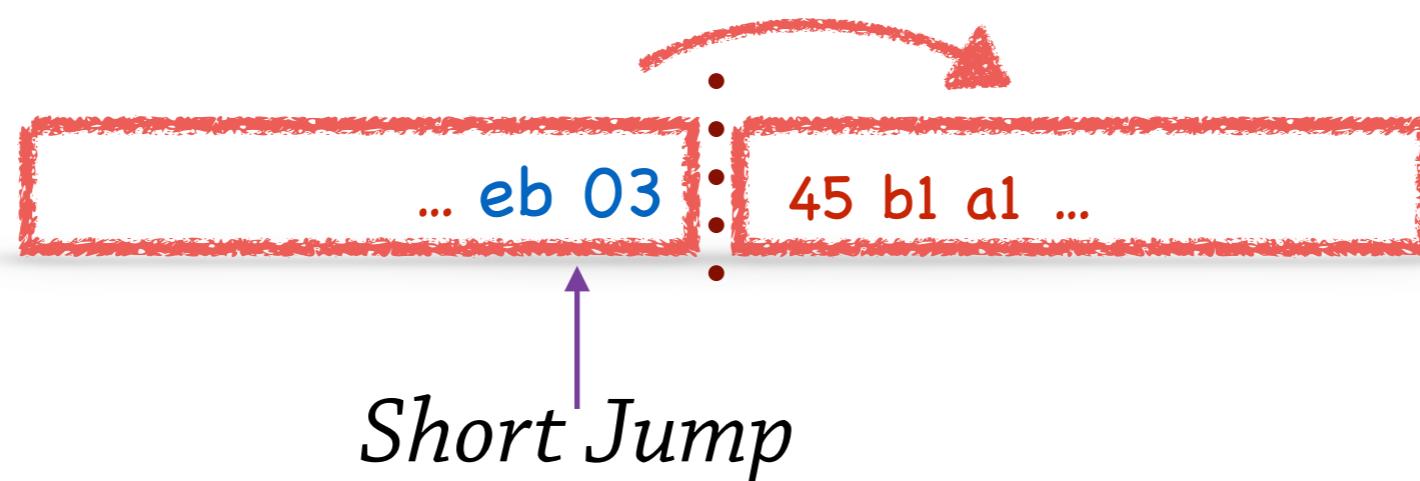
# [1 | 4] split



# [2 | 3] split



# [2 | 3] split

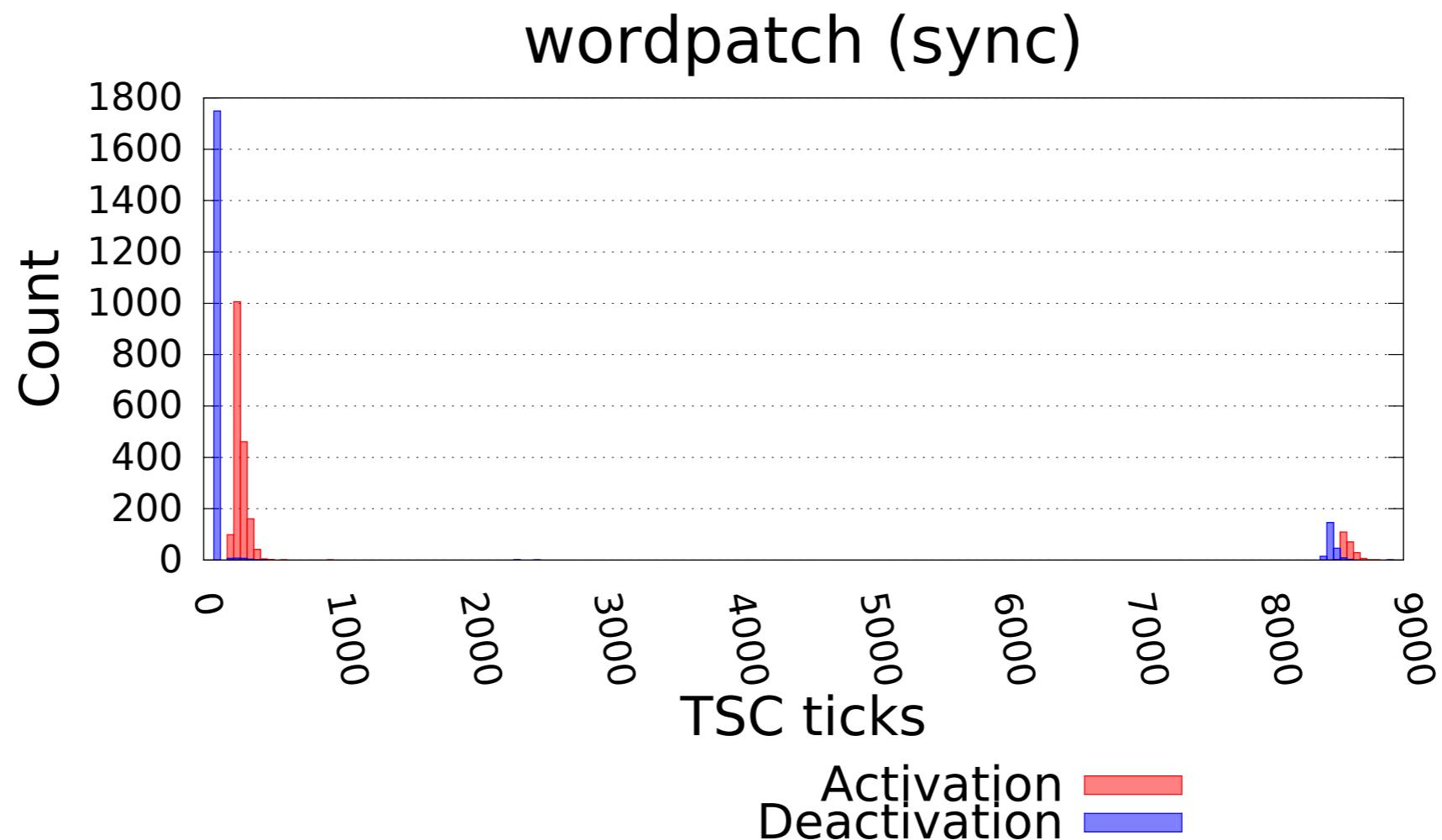


# [3 | 2] & [4 | 1] splits

***Reduces*** to [2 | 3] split  
short jump scheme

# Probing Costs

## ***Probe (De)Activation***

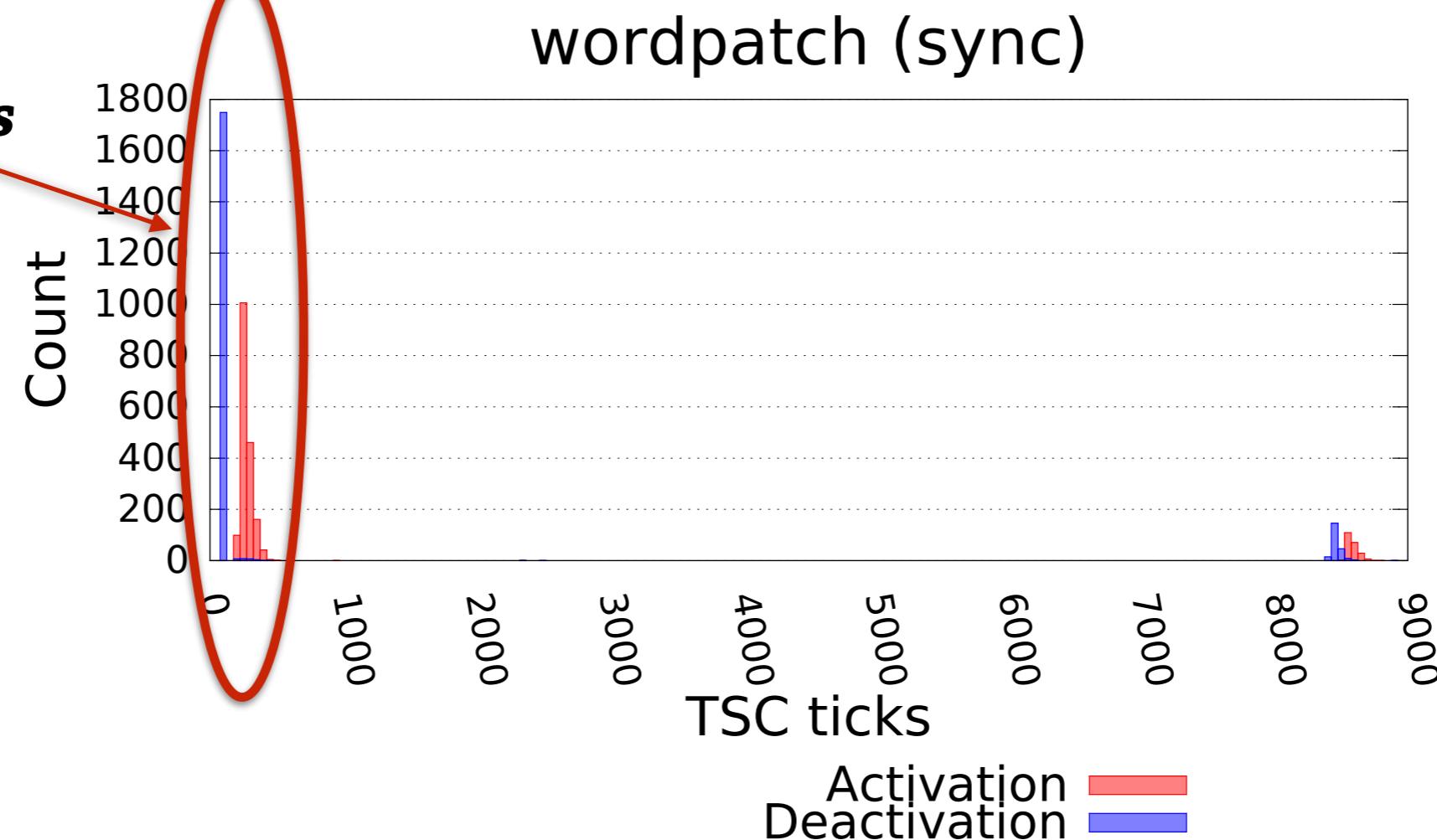


*On a synthetic application with 20000 probes*

# Probing Costs

## ***Probe (De)Activation***

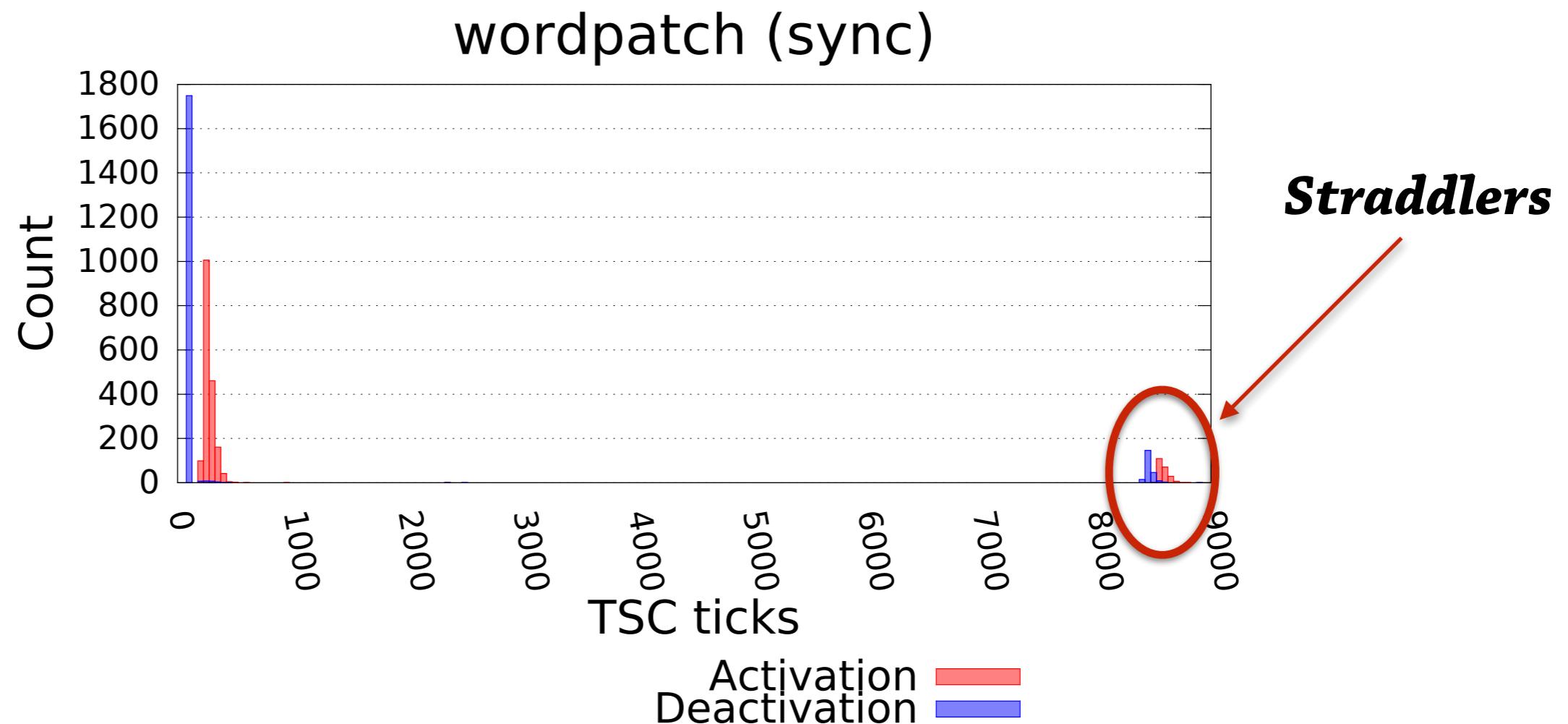
***Non straddlers***



*On a synthetic application with 20000 probes*

# Probing Costs

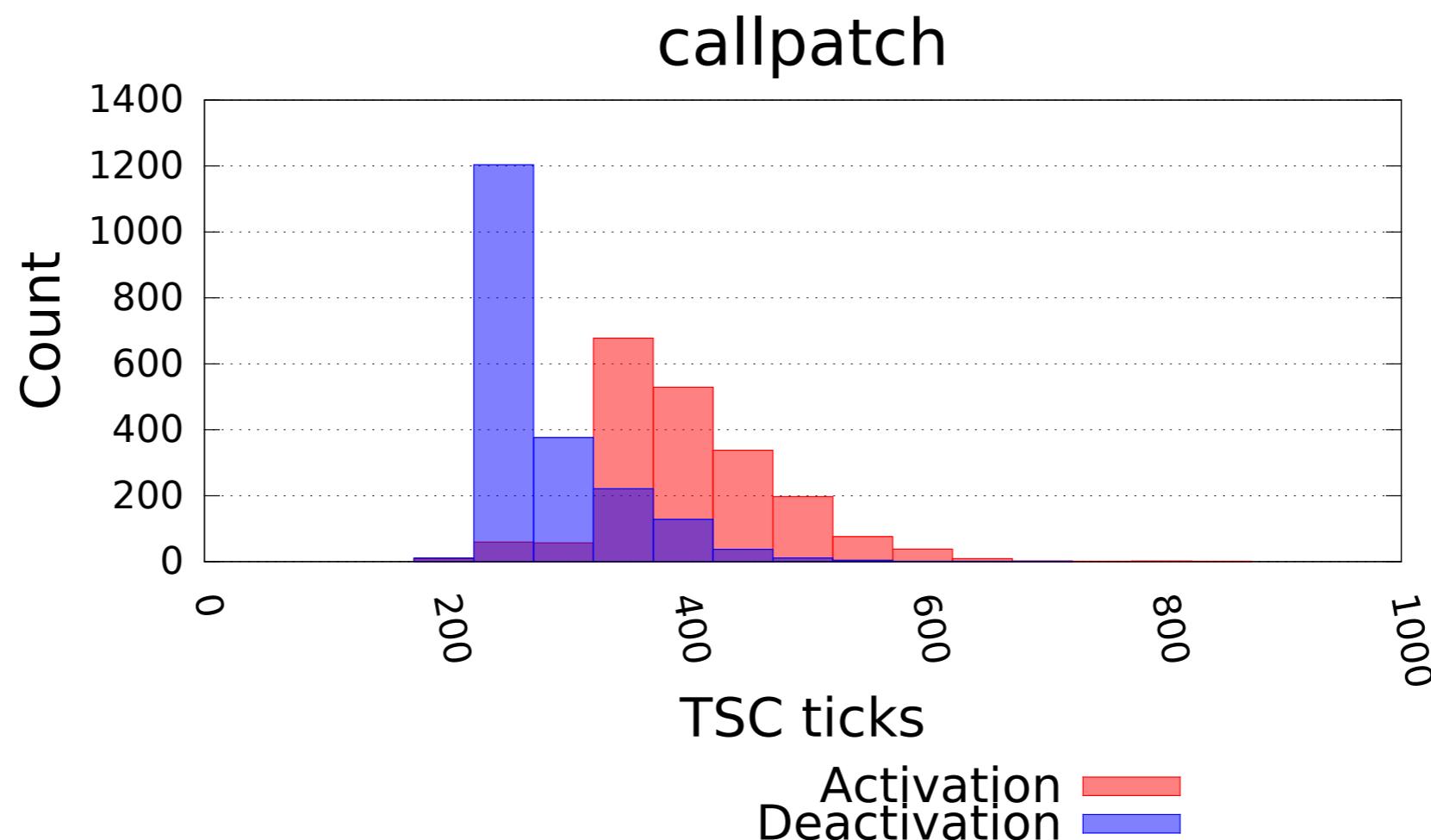
## ***Probe (De)Activation***



*On a synthetic application with 20000 probes*

# Probing Costs

## *Probe (De)Activation*



# Comparison with Other Options

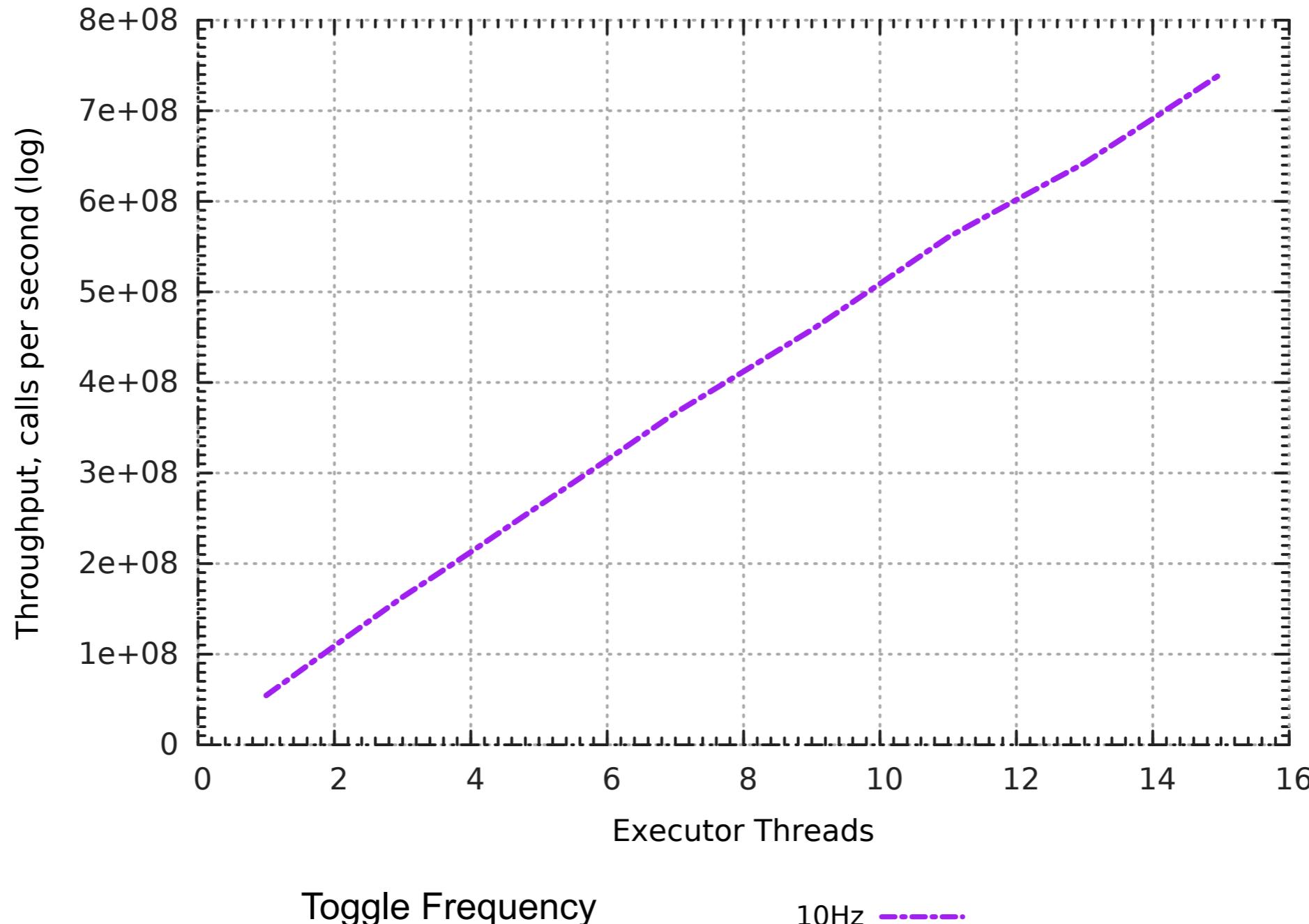
<b>Method</b>	<b>Activation</b>	<b>Deactivation</b>
<b>wordpatch (non straddlers)</b>	300	120
<b>wordpatch(straddlers)</b>	8850	8435
<b>callpatch</b>	432	315
<b>JVM VolatileCallSite</b>	995	995
<b>Dyninst</b>	1,929,244	995,447

# Comparison with Other Options

<b>Method</b>	<b>Activation</b>	<b>Deactivation</b>
<b>wordpatch (non straddlers)</b>	300	120
<b>wordpatch(straddlers)</b>	8850	8435
<b>callpatch</b>	432	315
<b>JVM VolatileCallSite</b>	995	995
<b>Dyninst</b>	1,929,244	995,447

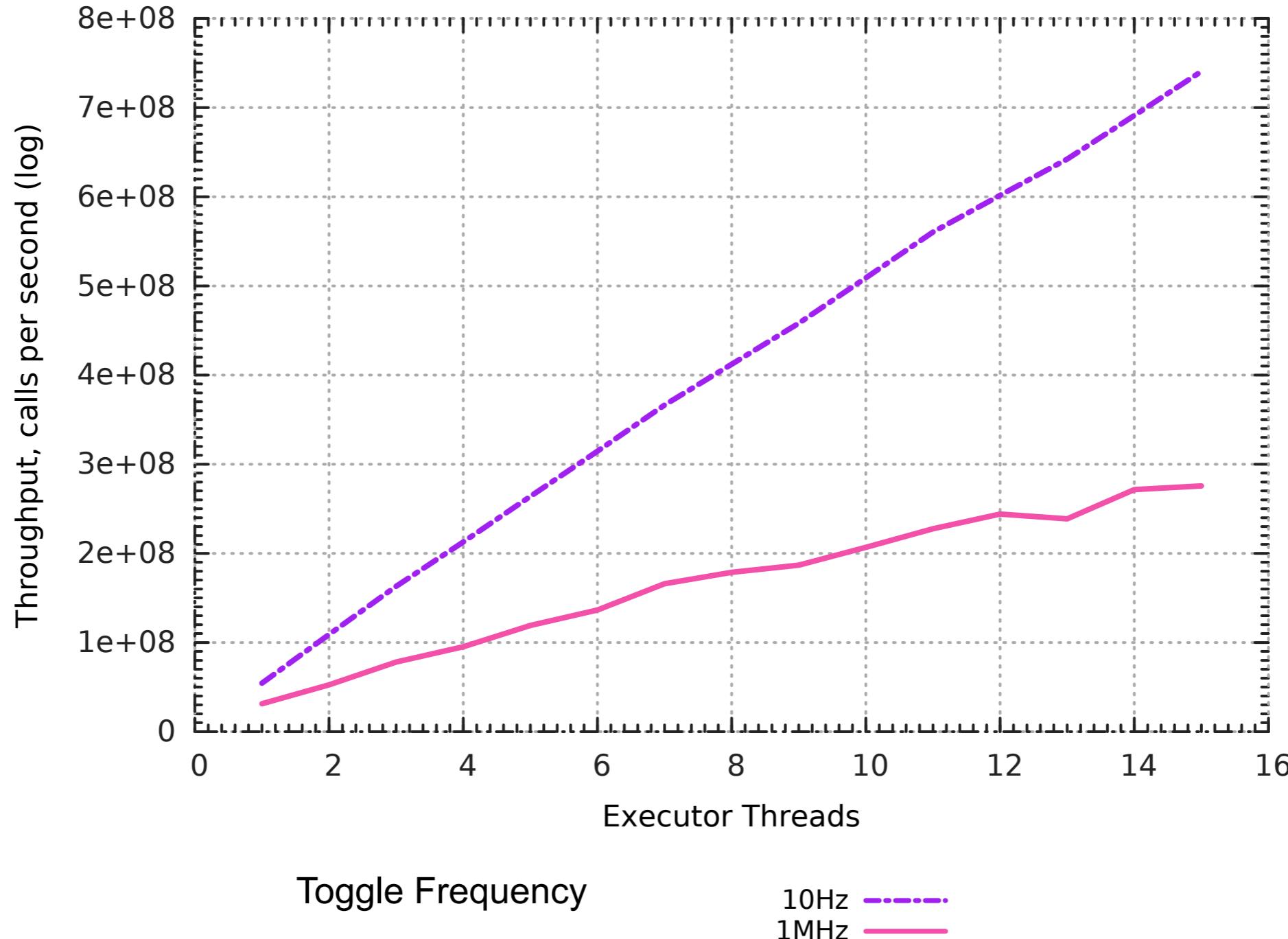
# Scalability

***Callpatch***



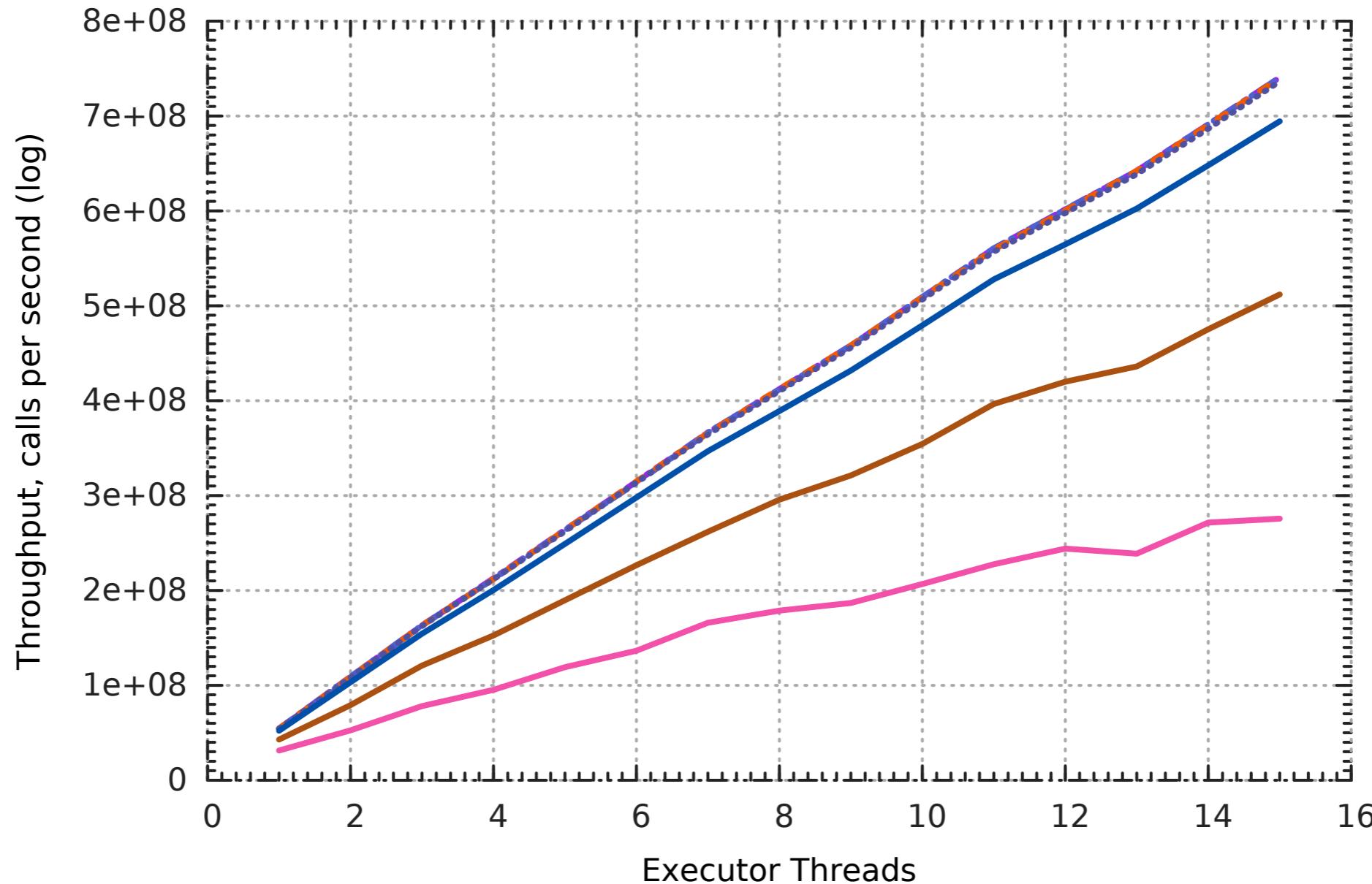
# Scalability

***Callpatch***



# Scalability

***Callpatch***

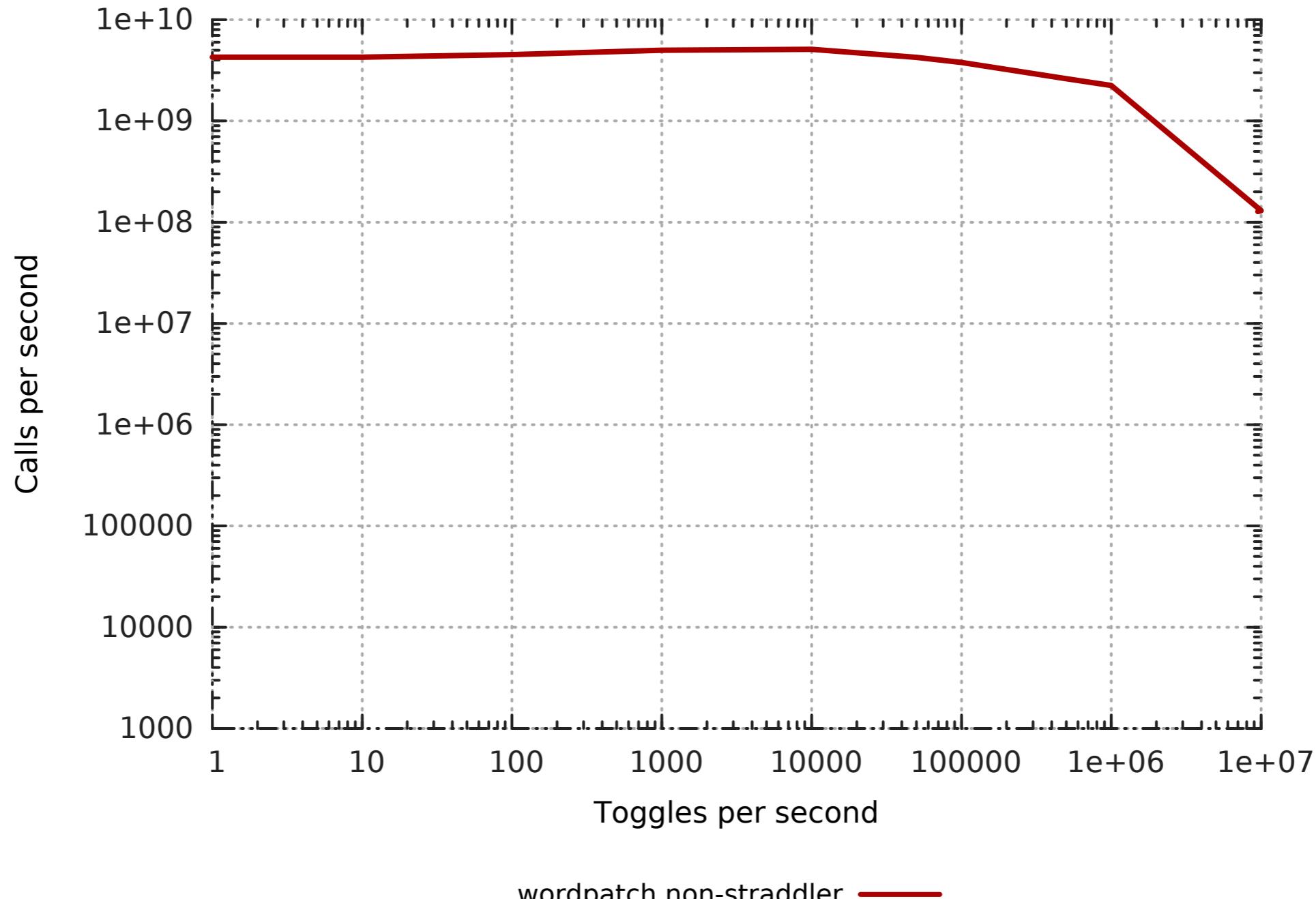


Toggle Frequency

10Hz	100kHz
100Hz	500kHz
1kHz	1MHz
10kHz	

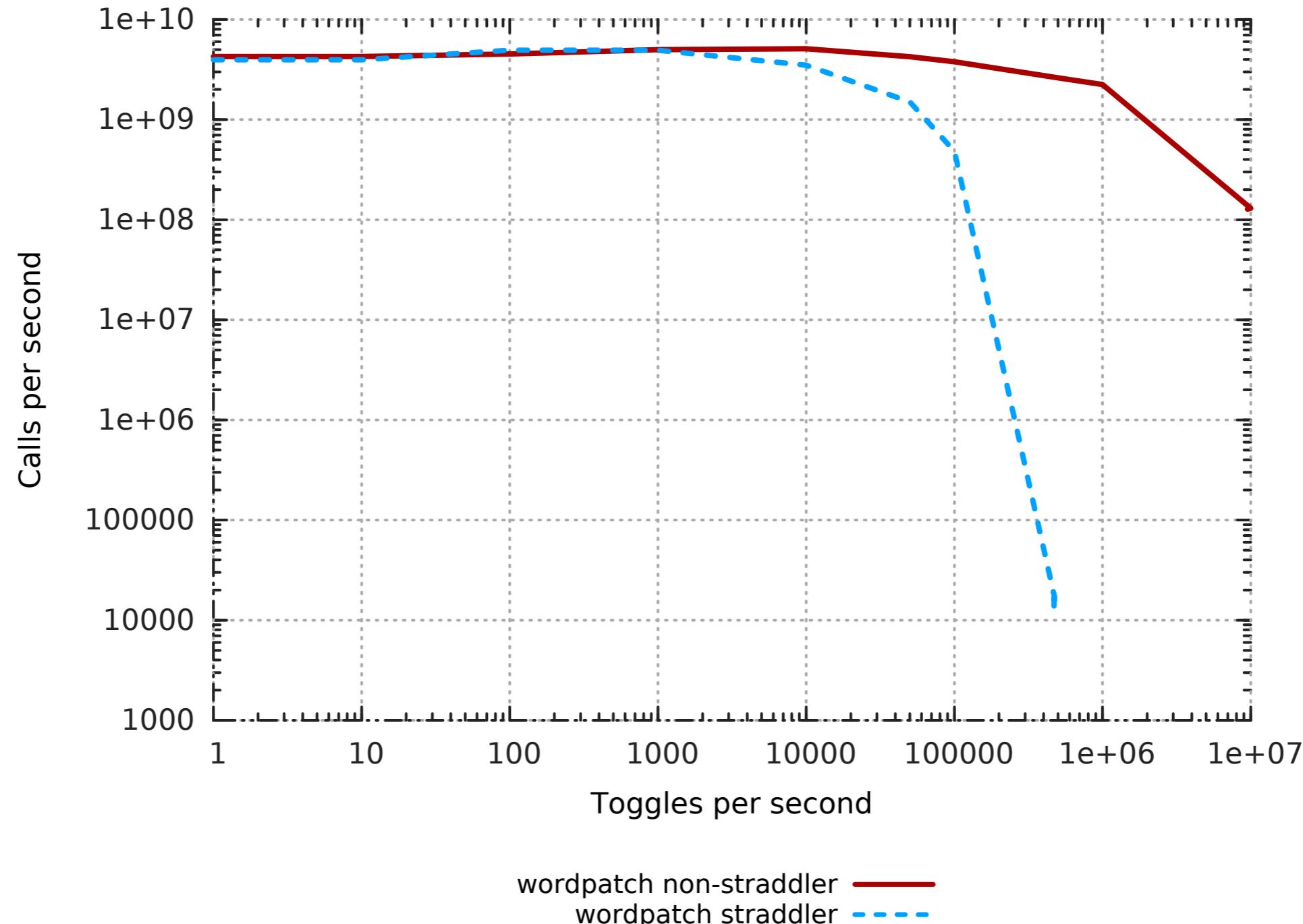
# Comparison

*Wordpatch*



# Comparison

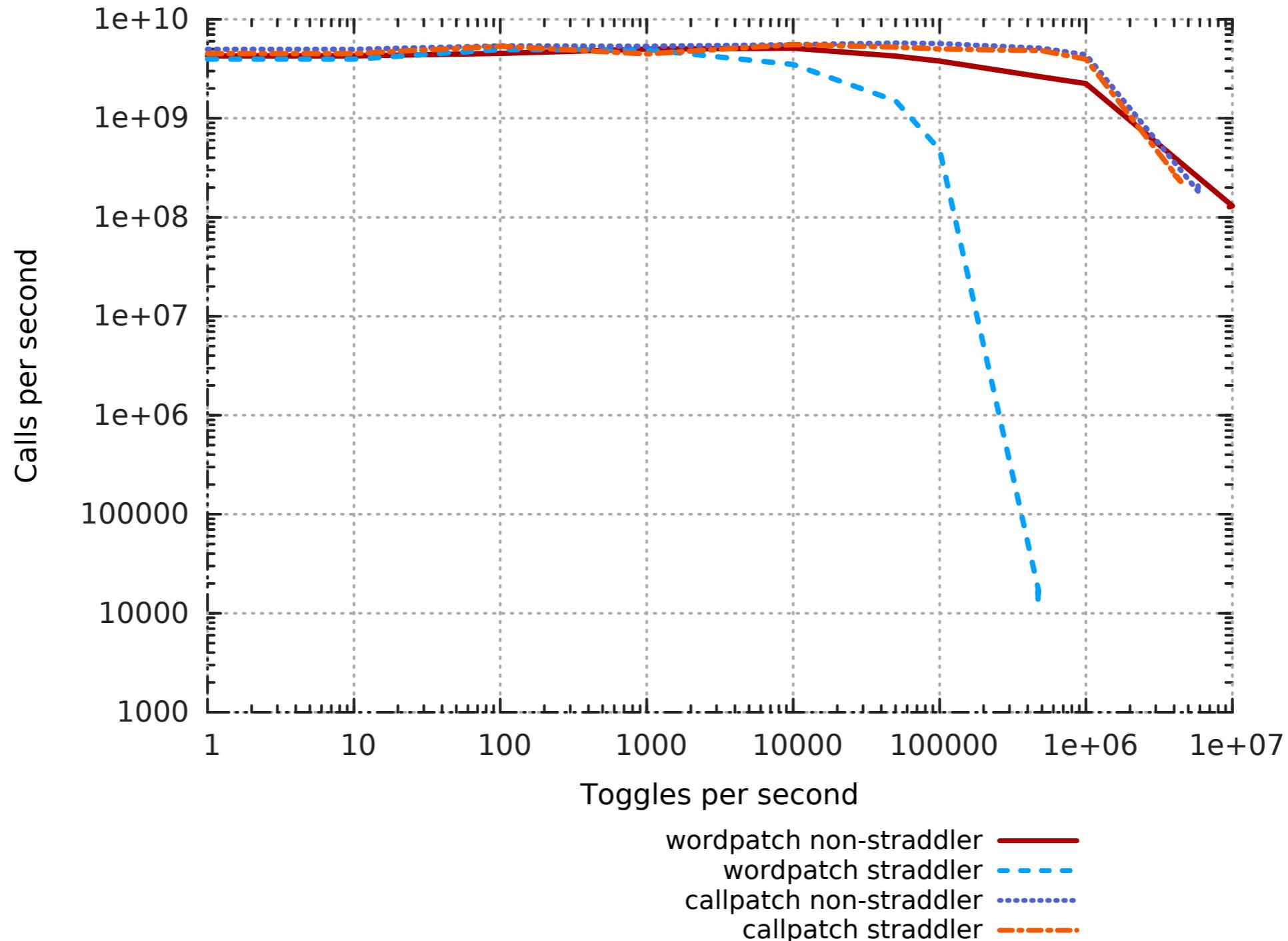
## *Wordpatch*



*Throughput is largely unaffected through **100Khz** frequency.*

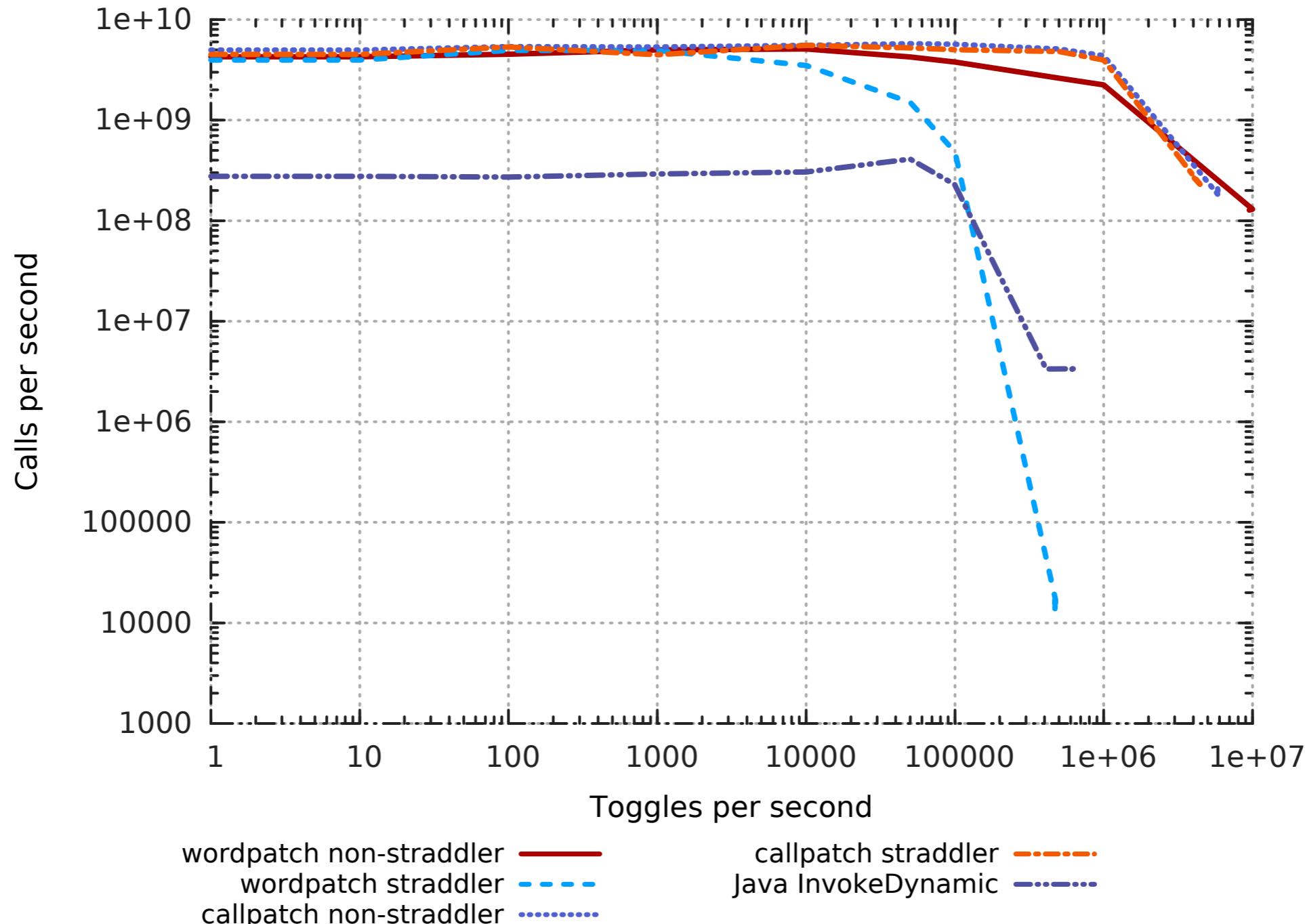
# Comparison

## *Wordpatch*



# Comparison

## *Wordpatch*



# Lightweight Profiler

- Measures function call counts, average execution times
- Does sampling with self deactivated probes
- Monitor thread wakes up periodically and reactivate inactive probes.

# Application Benchmarks

Benchmark	#probes	#toggles/sec	#samples/sec	slowdown%	toggle%
<b>h264ref</b>	638	24,015	115,526	6	0.1
<b>bzip</b>	162	2231	10,638	0.6	0.01
<b>sjeng</b>	142	14,616	73,060	3	0.06
<b>perl</b>	1408	42,276	211,346	11	0.2
<b>nbody</b>	252	8052	34,888	2	0.02
<b>hull</b>	164	1185	5636	0.2	0.005
<b>black scholes</b>	8	730	3650	1	0.002

# Application Benchmarks

Benchmark	#probes	#toggles/sec	#samples/sec	slowdown%	toggle%
<b>h264ref</b>	638	24,015	115,526	6	0.1
<b>bzip</b>	162	2231	10,638	0.6	0.01
<b>sjeng</b>	142	14,616	73,060	3	0.06
<b>perl</b>	1408	42,276	211,346	11	0.2
<b>nbody</b>	252	8052	34,888	2	0.02
<b>hull</b>	164	1185	5636	0.2	0.005
<b>black scholes</b>	8	730	3650	1	0.002

# Wrapping up

- ✓ Safe cross modification protocols on x86 for general instruction patching and call instruction patching
- ✓ No stop the world requirements
- ✓ A rapid probe toggling implementation
- ✓ Useful for profilers, JIT compilers, live software updaters or in other DBIs for x86 etc.



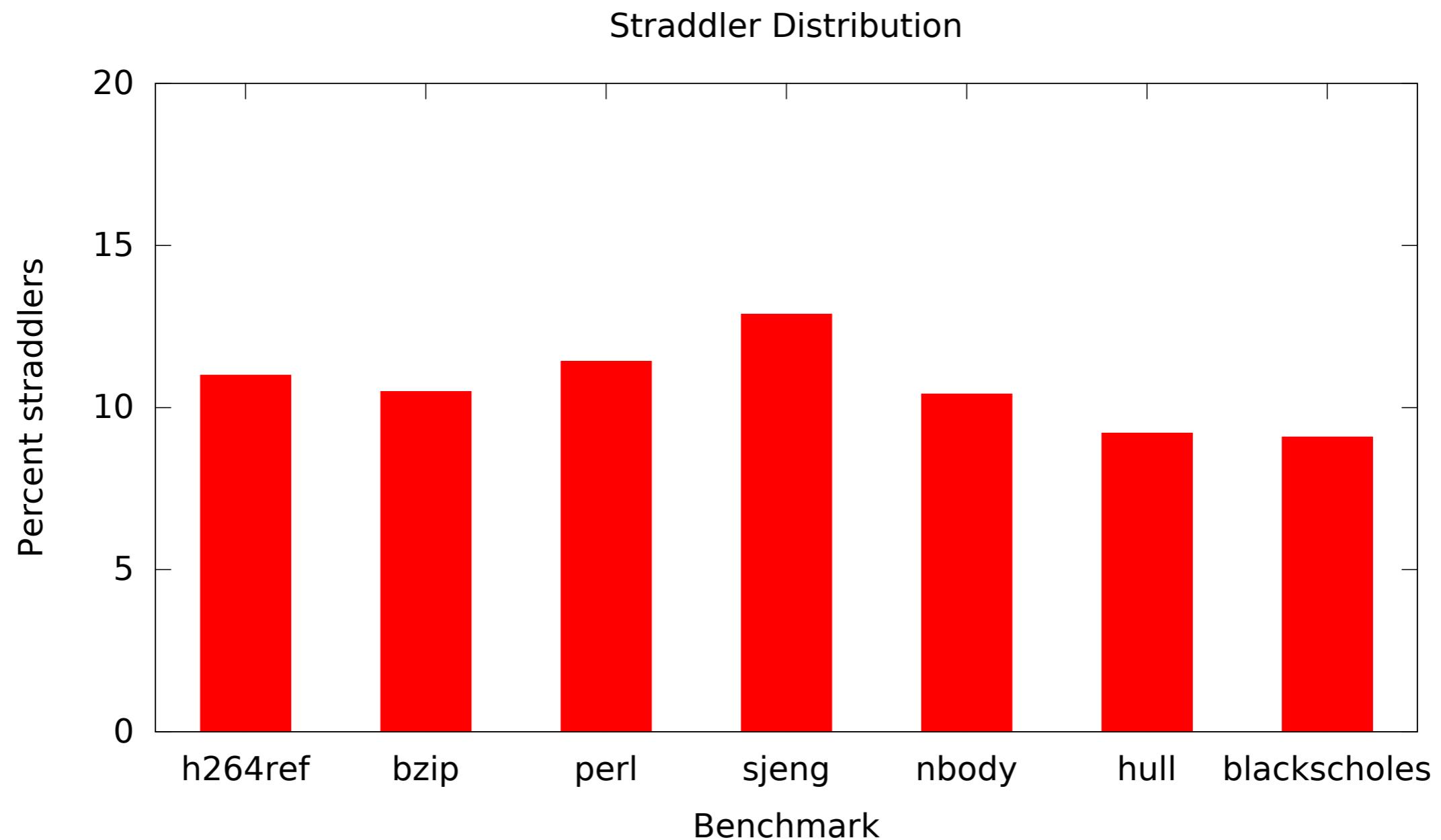
*Buddhika Chamith*  
[budkahaw@indiana.edu](mailto:budkahaw@indiana.edu)  
*github : @chamibuddhika*

# Safe wordpatch

```
// Write an interrupt byte into the address,
void int3_lock(address)
    return write_front(address, INT3);

bool patch(address, value)
    if (not is_straddler(address))
        write(address, value);
        return true;
    else if (int3_lock(address))
        wait();
        write_back(address, value);
        wait();
        write_front(address, value);
        return true;
    else return false;
```

# Straddler Distribution

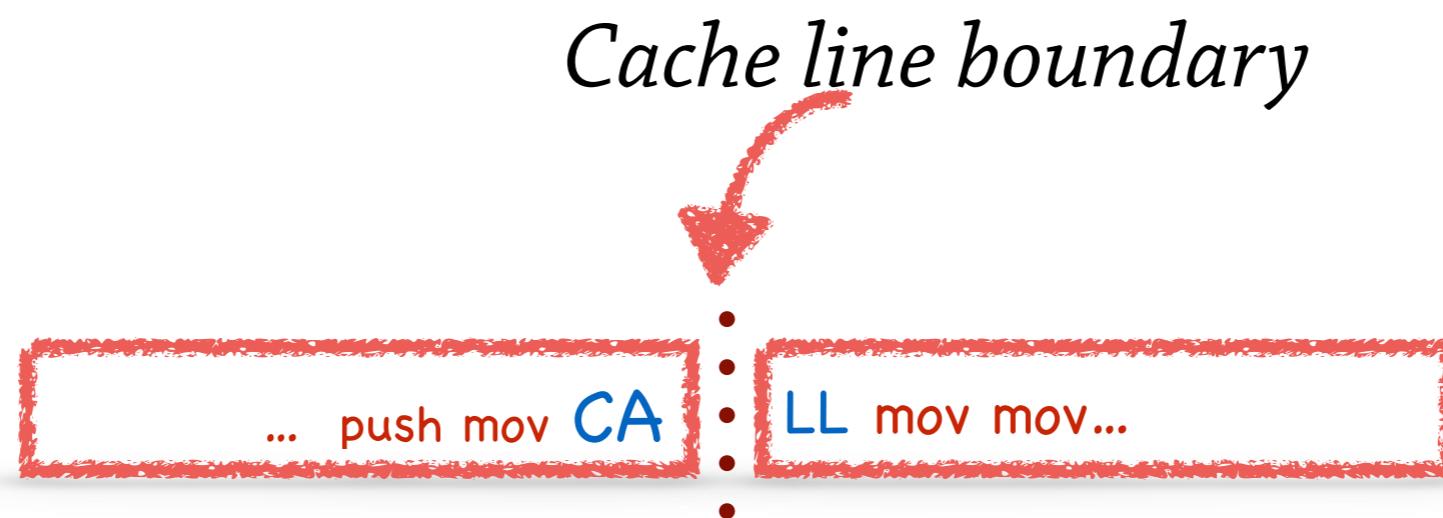


# Intel Cross Modification

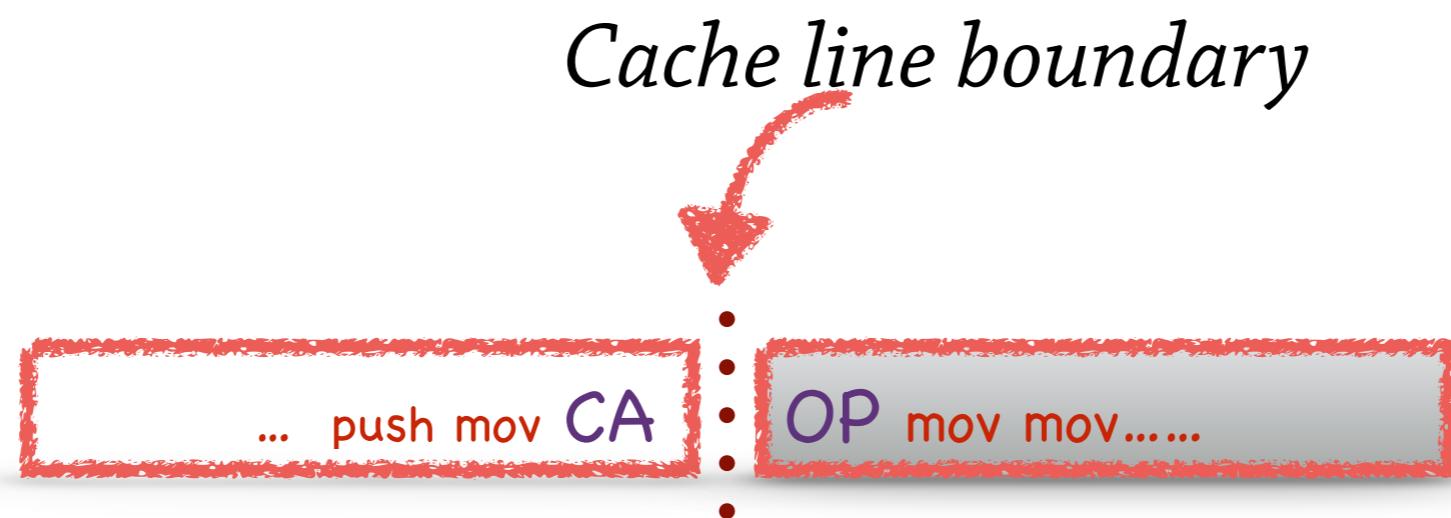
```
// -- Action of Modifying Processor --
memory_flag = 0;
Store modified code (as data) into code segment;
memory_flag = 1;

// -- Action of Executing Processor --
WHILE (memory_flag != 1)
Wait for code to update;
ELIHW;
Execute serializing instruction; // e.g: CPUID
Begin executing modified code;
```

# Woes of cmpxchg



# Woes of cmpxchg



**INCORRECT**

# Latent Costs

## ***Potential sources***

- 5 byte NOOP : ~1-2 cycles
- Callpatch : Inactive state jumps
  - ▶ Short : ~2 cycles
  - ▶ Degenerate RET call : ~4 cycles

# Probing Costs

Method	Activation	Deactivation	Invocation
wordpatch :			
non straddlers	300	120	35
straddlers	8850	8435	35
callpatch	432	315	35
JVM VolatileCallSite	995	995	55
Dyninst	1,929,244	995,447	320