



*BY Developers FOR Developers*

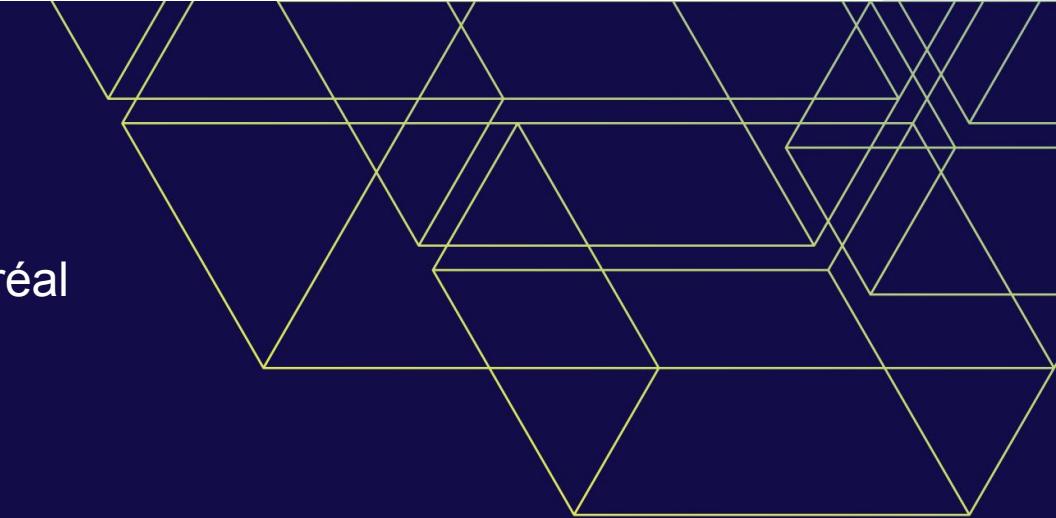
Storage Developer Conference  
September 22-23, 2020

# Tracing and Visualizing FS Internals with eBPF Superpowers

**Suchakra Sharma**  
ShiftLeft Inc.

**Hani Nemati**  
Microsoft





# Suchakra Sharma

PhD, École Polytechnique de Montréal

Staff Scientist  
**ShiftLeft Inc**

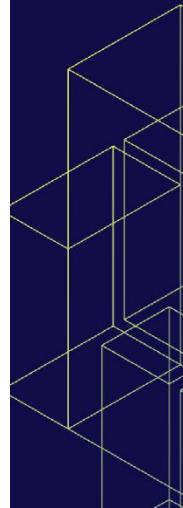
# Hani Nemati

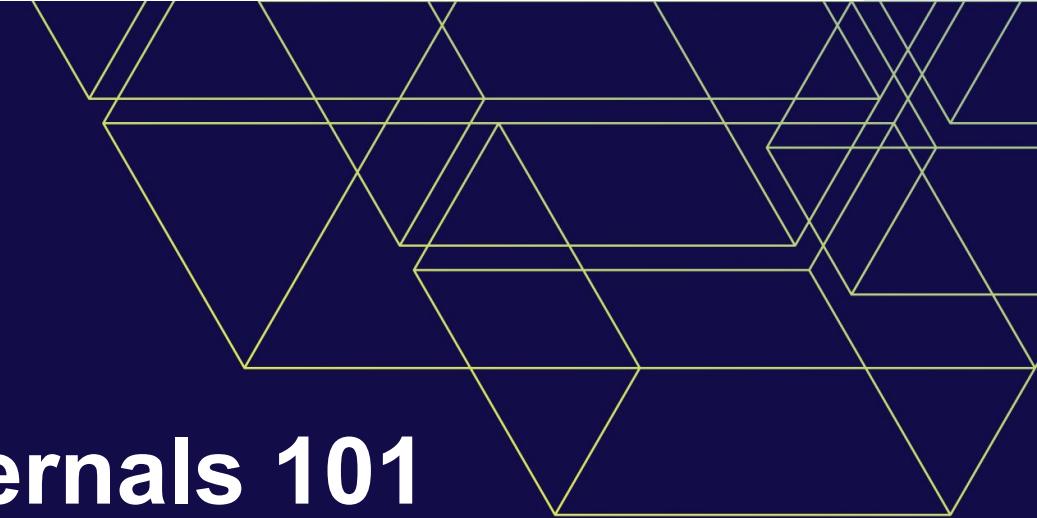
PhD, École Polytechnique de Montréal

Software Engineer (Kernel)  
**Microsoft**

# Agenda

- **FS Internals 101**
- **Journey of a `read()`**
  - Tracing Control-Flow via `ftrace`
- **FS and Performance**
  - Probing the system via `/proc`
  - BCC and `bpftrace` Tool Collection
- **Building tools with eBPF**
  - Case Study: Read-ahead
  - DIY new BPF tool - `readaheadstat`
  - Visualizing Performance

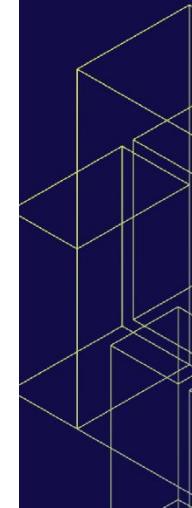




# FS Internals 101

# Why Filesystem?

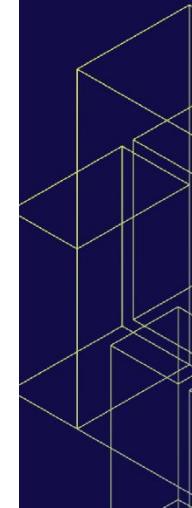
- **Keep track of things**
  - What many files are open or being read?
  - Which user/process is operating with files?
  - How many files can be created?
  - What data is most frequently accessed
- **Structure**
  - How are they structured?
  - What happens when we move them across?
- **Abstraction and Uniformity**
  - What if some of our files are on a network storage?
  - What if we copy a file from two different storage devices?



# FS Internals

APP

```
int fd = open("foo", R)
```



# FS Internals

APP

```
int fd = open("foo", R)
```

open()

sys\_open

Syscall Interface

vfs\_open

Virtual File System

ext4\_file\_open

Logical File System

Block Devices

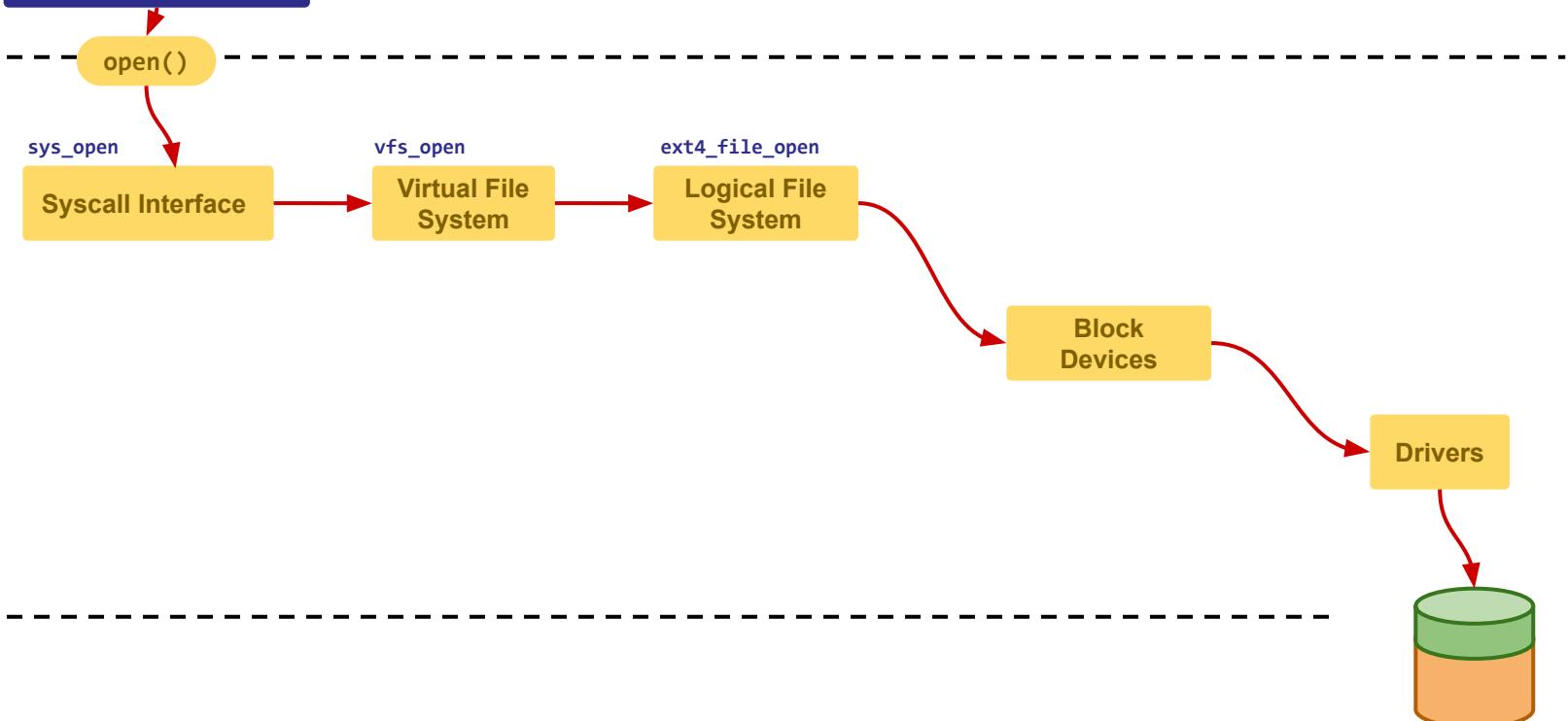
Drivers



# FS Internals

APP

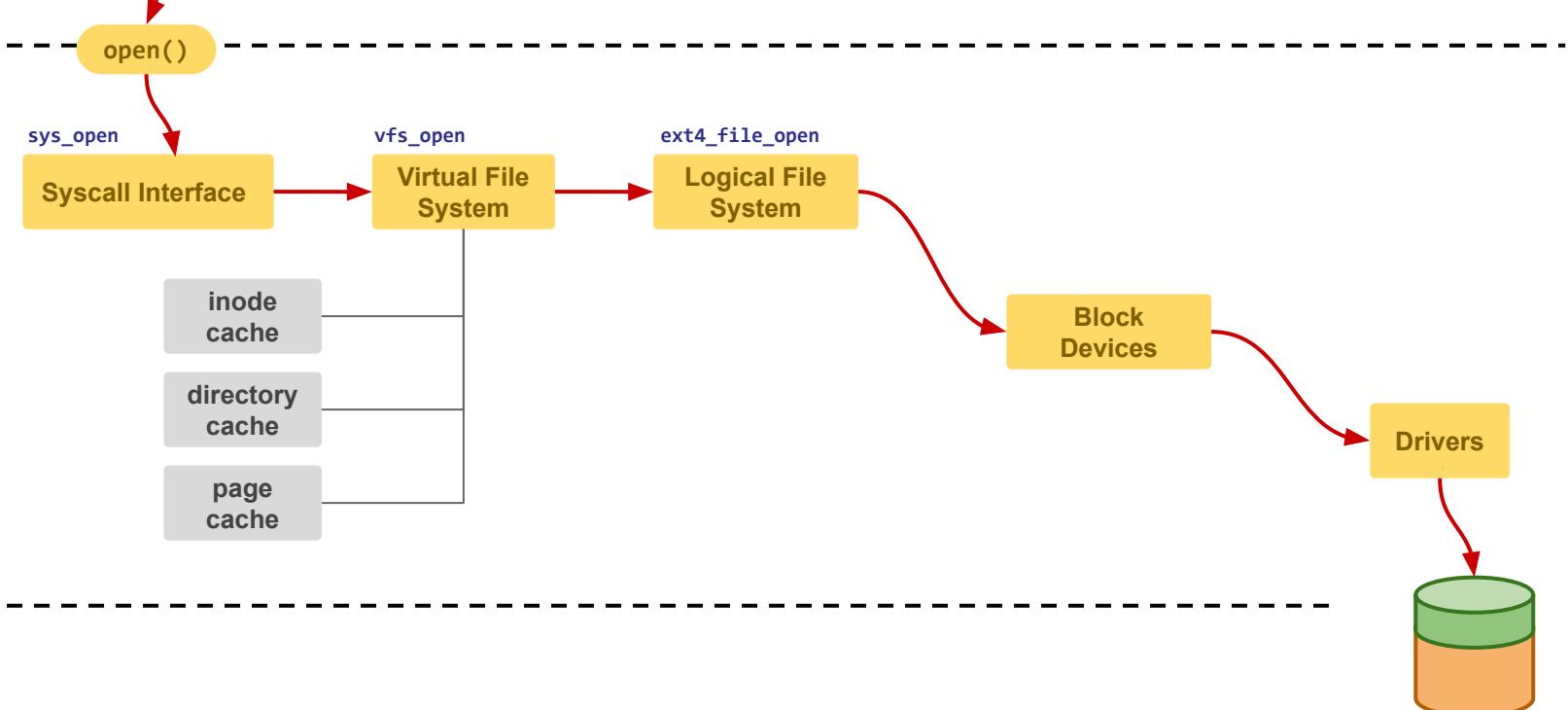
```
int fd = open("foo", R)
```

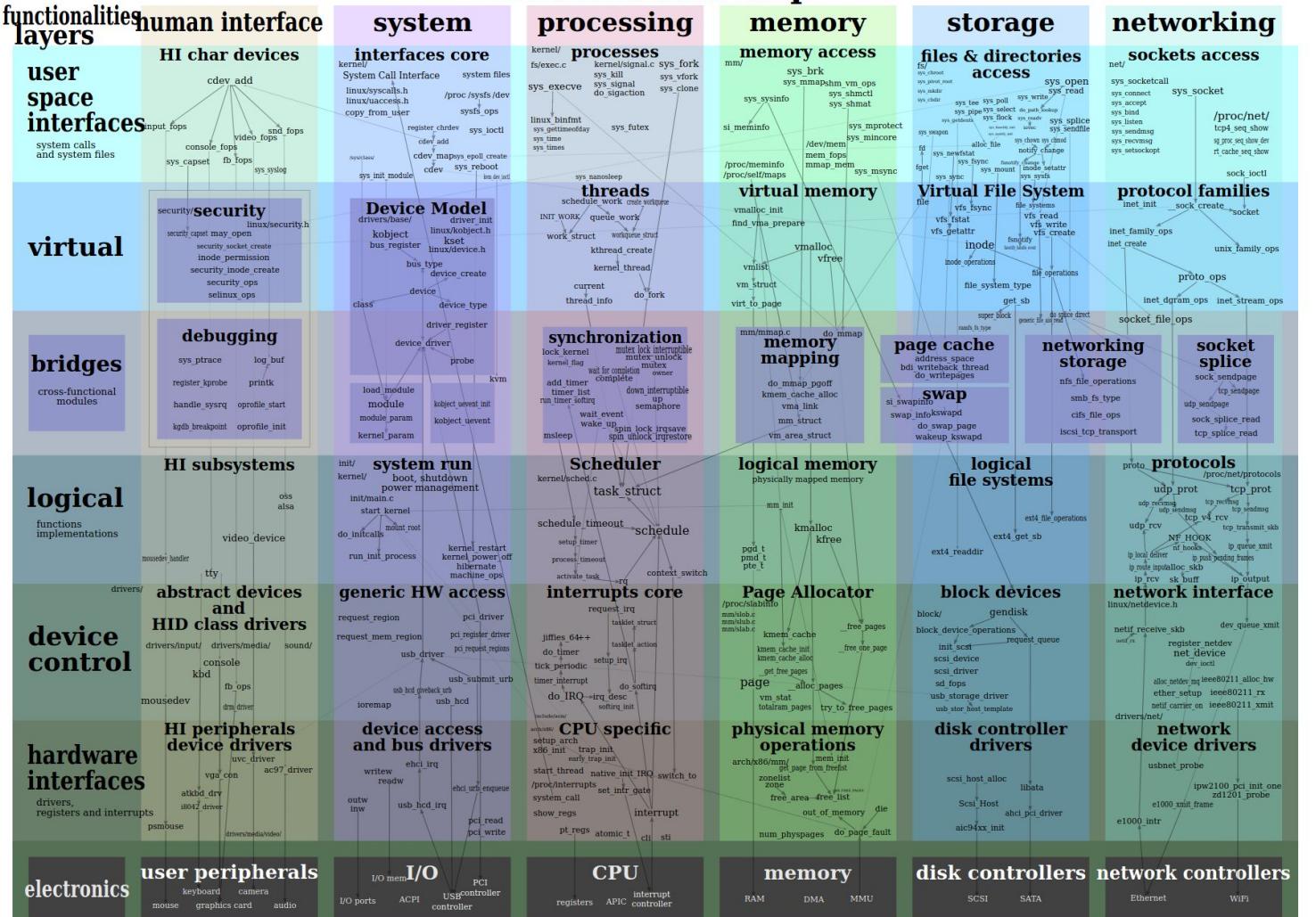


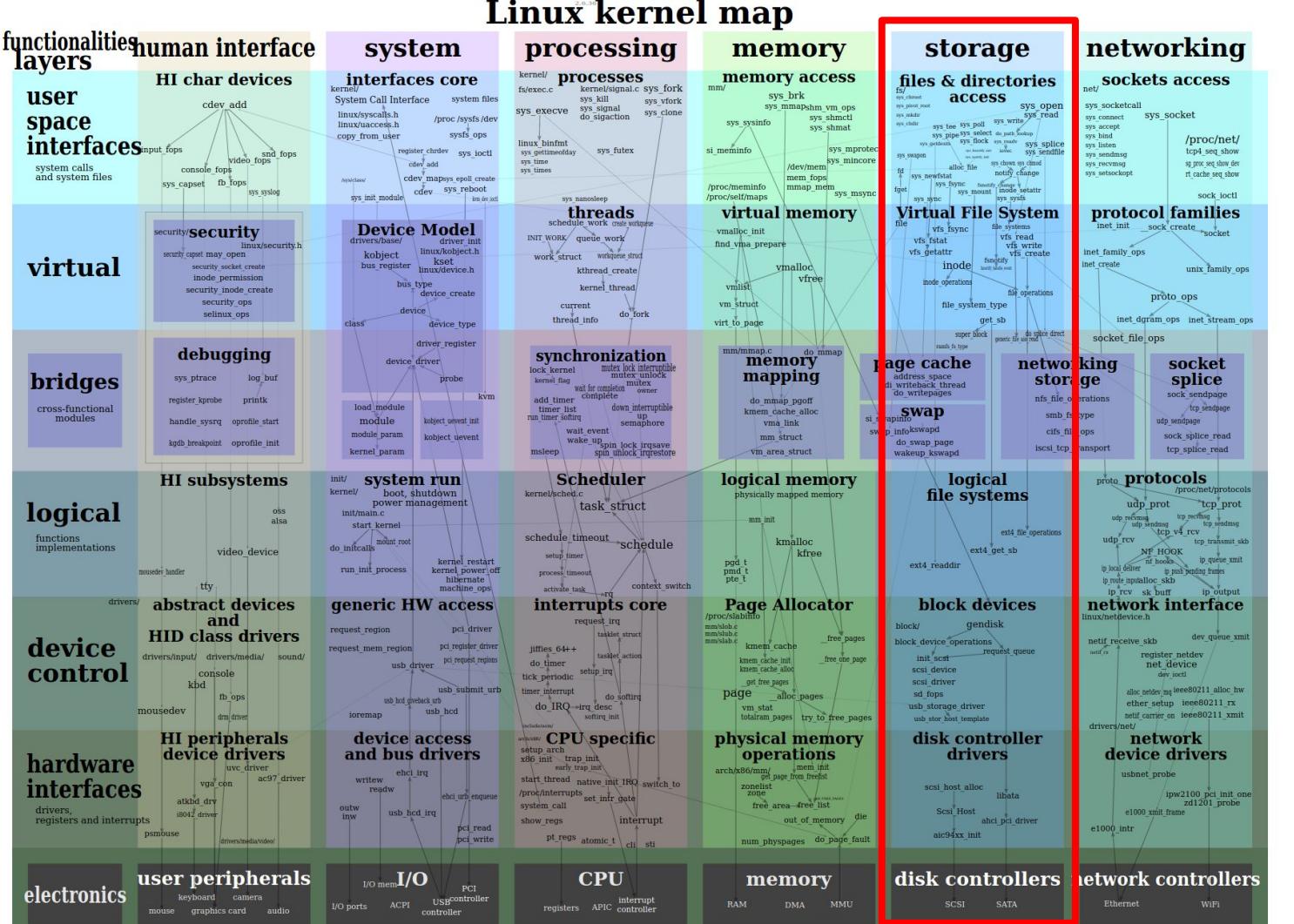
# FS Internals

APP

```
int fd = open("foo", R)
```

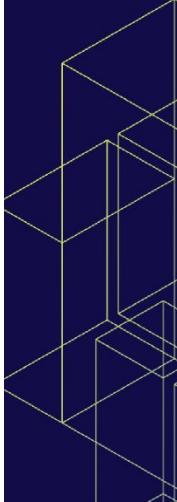






# Tracing a read() with ftrace

```
$ sudo trace-cmd record -p function -P 2535  
$ sudo trace-cmd report
```



# Tracing a read() with ftrace

```
$ sudo trace-cmd record -p function -P 2535  
$ sudo trace-cmd report
```

# Tracing a read() with ftrace

```
$ sudo trace-cmd record -p function -P 2535  
$ sudo trace-cmd report
```

```
do_syscall_64
    _x04_sys_read
    _sys_read
        _fdget_pos
            _fget_light
    vfs_readdir
        _rv_verify_area
            security_file_permission
                apparent_file_permission
                    common_file_perm
                        aa_file_perm
            rsnotify_parent
                rsnify
    _vfs_readdir
        new_sync_file
            ext4_file_read_iter
                generic_file_read_iter
                    generic_file_buffered_read
                        _csm_resched
                        _rcu_qs
                pagecache_get_page
                    find_get_entry
                        PageHigh
                mark_page Accessed

_cond_resched
    rcu_all_qs
        _cond_resched

rcu_all_qs
    pagecache_get_page
        find_get_entry
            PageHigh
    touch_atime
        atime_needs_update
            current_time
                Ktime_get_course_real_ts64
                timestamp_truncate
    sb_start
        _mnt_want_write
    current_time
        Ktime_get_course_real_ts64
        timestamp_truncate
    update_time
        generic_update_time
            __mark_inode_dirty
                ext4_inode_dirty
                    ext4_mark_start_sb
                    ext4_journal_check_start
                    _cond_resched

rcu_all_qs
    jbd2_journal_start
        _keev_cache_alloc
        _cond_resched

rcu_all_qs
    should_folio_lab
        memcg_kmem_put_cache
        start_this_handle
            _rcu_head_lock
            add_transaction_credits
        ext4_mark_inode_dirty
            _cond_resched

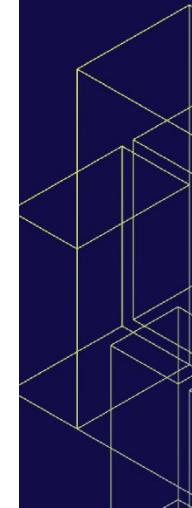
rcu_all_qs
    ext4_reserve_inode_write
        _ext4_get_inode_loc
            ext4_get_group_desc
            ext4_group_table
                __getblk_gfp
                    _find_get_block
```

- 1 read syscall
- 5 other syscalls
- 555 function calls

# 8MB in 2 seconds

# Tracing a read() with ftrace

```
do_syscall_64
  _x64_sys_read
    ksys_read
      _fget_pos
        fget_light
        vfs_read
          rw_verify_area
            security_file_permission
              apparm_file_permission
                common_file_perm
                  aa_file_perm
                    fnotify_parent
                      fnotify
        _vfs_read
          new_sys_read
            ext4_file_read_iter
              generic_file_read_iter
                generic_file_buffered_read
                  _cond_resched
                    rCU_all_qs
                    pagecache_get_page
                      find_get_entry
                        PageHuge
                      mark_page_accessed
        _cond_resched
        rCU_all_qs
        _cond_resched
        pagecache_get_page
          find_get_entry
            PageHuge
          touch_atime
            atime_needs_update
              current_time
                ktime_get_coarse_real_ts64
                timestamp_truncate
              __sb_start_write
                __mnt_want_write
                current_time
                  ktime_get_coarse_real_ts64
                  timestamp_truncate
                update_time
                  generic_update_time
                    __mark_inode_dirty
                      ext4_dirty_inode
                      __ext4_journal_start_sb
                        __ext4_journal_check_start
                          _cond_resched
        rCU_all_qs
          jbd2_journal_start
            kmem_cache_alloc
              _cond_resched
        rCU_all_qs
          should_failslab
            memcg_kmem_put_cache
            start_this_handle
              __raw_read_lock
                add_transaction_credits
              ext4_mark_inode_dirty
                _cond_resched
        rCU_all_qs
          ext4_reserve_inode_write
            __ext4_get_inode_loc
              __ext4_get_group_desc
                __ext4_mnt_table
                  __getblk_gfp
                    __find_get_block
```

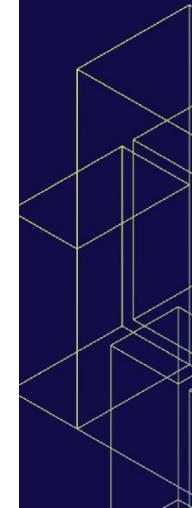


# Tracing a read() with ftrace

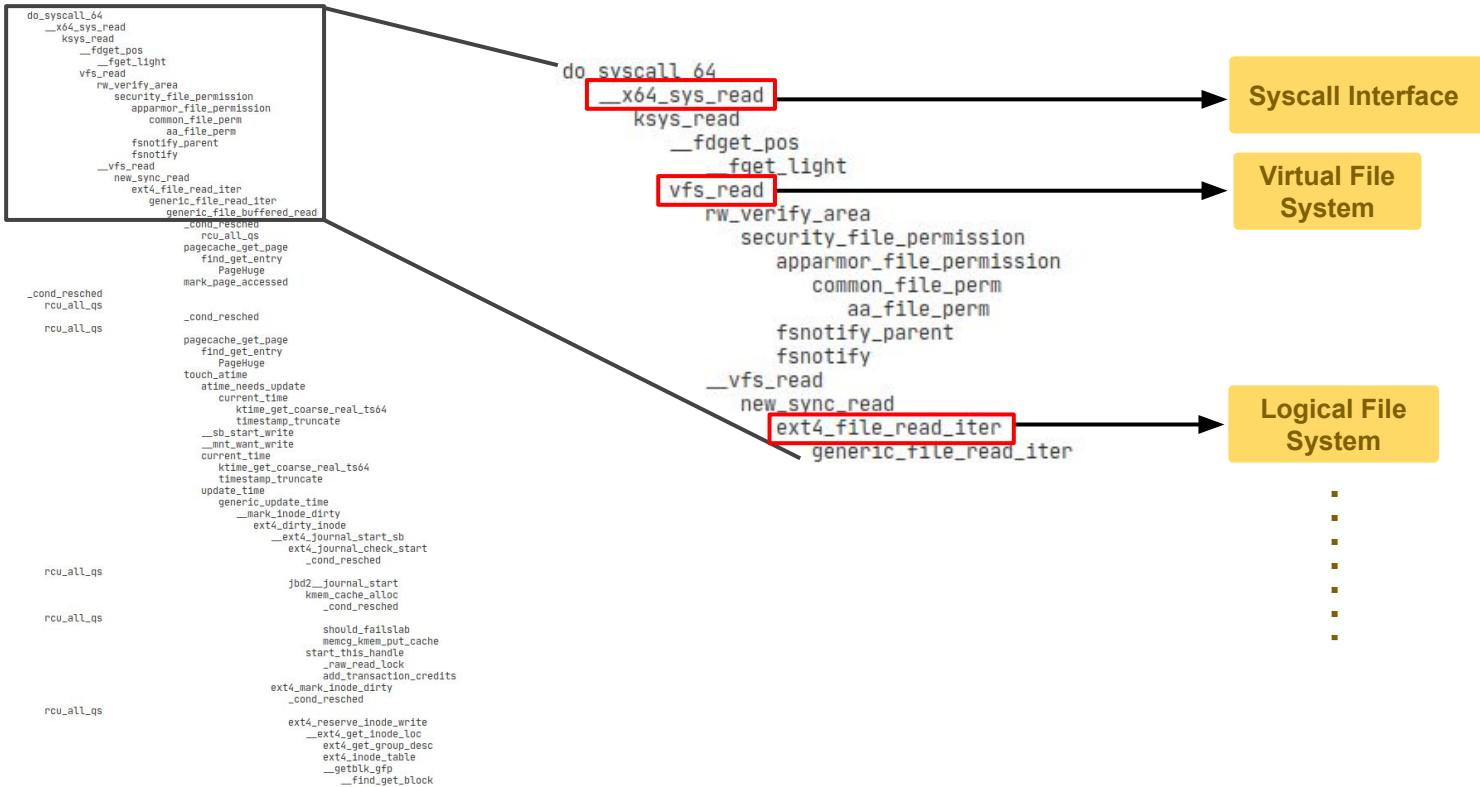
```

do_syscall_64
  _x64_sys_read
    ksys_read
      __fdget_pos
        fget_light
        vfs_read
          RW_verify_area
            security_file_permission
              apparmor_file_permission
                common_file_perm
                  aa_file_perm
                    fsnotify_parent
                      fsnotify
        __vfs_read
          new_sync_read
            ext4_file_read_iter
              generic_file_read_iter
                generic_file_buffered_read
                  -cond_resched
                    rCU_all_qs
                    pagecache_get_page
                      find_get_entry
                        PageHuge
                      mark_page_accessed
        _cond_resched
        rCU_all_qs
        _cond_resched
        pagecache_get_page
          find_get_entry
            PageHuge
          touch_atime
            atime_needs_update
              current_time
                ktime_get_coarse_real_ts64
                timestamp_truncate
              __sb_start_write
                __mnt_want_write
                current_time
                  ktime_get_coarse_real_ts64
                  timestamp_truncate
                  update_time
                    generic_update_time
                      __mark_inode_dirty
                        ext4_dirty_inode
                          __ext4_journal_start_sb
                            ext4_journal_check_start
                            _cond_resched
        rCU_all_qs
        jbd2_journal_start
          kmem_cache_alloc
            _cond_resched
              should_failslab
                memcg_kmem_put_cache
                start_this_handle
                  __raw_read_lock
                    add_transaction_credits
                    ext4_mark_inode_dirty
                    _cond_resched
                    ext4_reserve_inode_write
                      __ext4_get_inode_loc
                        ext4_get_group_desc
                        ext4_mnt_table
                          __getblk_gfp
                            __find_get_block
  rCU_all_qs

```

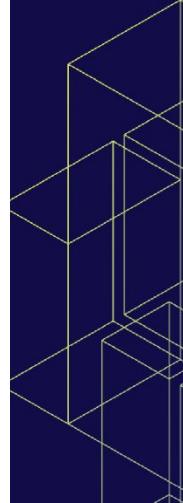


# Tracing a read() with ftrace



# Tracing a read() with ftrace

```
example-2535 [007] 1580764.896732: function:  
example-2535 [007] 1580764.896732: function:  
example-2535 [007] 1580764.896732: function:  
example-2535 [007] 1580764.896733: function:  
example-2535 [007] 1580764.896733: function:  
example-2535 [007] 1580764.896734: function:  
          do_syscall_64  
          __x64_sys_read  
          ksys_read  
          __fdget_pos  
          __fget_light  
          vfs_read  
  
          ...  
  
example-2535 [007] 1580764.896770: function:  
example-2535 [007] 1580764.896770: function:  
          __wake_up_common_lock  
          __raw_spin_lock_irqsave  
  
          ...  
  
example-2535 [007] 1580764.896794: function:  
          irq_exit  
  
          ...  
  
example-2535 [007] 1580764.896808: function:  
example-2535 [007] 1580764.896808: function:  
example-2535 [007] 1580764.896809: function:  
example-2535 [007] 1580764.896810: function:  
example-2535 [007] 1580764.896810: function:  
          __mnt_drop_write  
          __sb_end_write  
          fsnotify_parent  
          fsnotify  
          fpregs_assert_state_consistent
```



# Tracing a read() with ftrace

# Tracing a read() with ftrace

```
example-2535 [007] 1580764.896732 function: do_syscall_64
example-2535 [007] 1580764.896732 function: __x64_sys_read
example-2535 [007] 1580764.896732 function: ksys_read
example-2535 [007] 1580764.896733 function: __fdget_pos
example-2535 [007] 1580764.896733 function: __fget_light
example-2535 [007] 1580764.896734 function: vfs_read
example-2535 [007] 1580764.896770 function: __wake_up_common_lock
example-2535 [007] 1580764.896770 function: _raw_spin_lock_irqsave
example-2535 [007] 1580764.896794 function: irq_exit
example-2535 [007] 1580764.896808 function: __mnt_drop_write
example-2535 [007] 1580764.896808 function: __sb_end_write
example-2535 [007] 1580764.896809 function: fsnotify_parent
example-2535 [007] 1580764.896810 function: fsnotify
example-2535 [007] 1580764.896810 function: fpregs_assert_state_consistent
```

Syscall Enter

IRQ Entry

IRQ Exit

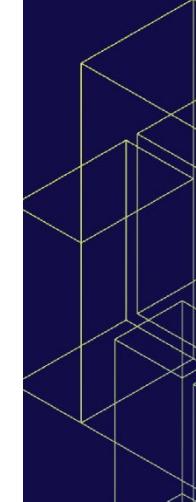
Syscall Exit

# **FS and Performance**

# What can /proc tell us?

```
# cat /proc/diskstats
```

Major	Minor	Device	rrqm_s	rrqm_k	rqesv_s	rqesv_k	rqtw_s	rqtw_k	avgq_s	avgq_k	avgqu_s	avgqu_k	avgquw_s	avgquw_k	avgquwz_s	avgquwz_k	avgquwzv_s	avgquwzv_k	avgquwzvz_s	avgquwzvz_k
259	0	nvme0n1	1823879	1141588	58777929	328845	3512978	2244953	151315050	7397322	0	2759684	7859940	88819	0	543261592	26313	273795	107459	
259	1	nvme0n1p1	1443	0	27072	2273	2	0	2	34	0	392	2310	4	0	1034184	1	0	0	
259	2	nvme0n1p2	573	398	59794	77	367	115	151352	415	0	776	525	91	0	645136	33	0	0	
259	3	nvme0n1p3	1820929	1141190	58667503	326372	3238845	2244838	151163696	7287293	0	2747016	7639943	88724	0	541582272	26277	0	0	



<sup>1</sup><https://www.kernel.org/doc/Documentation/iostats.txt>

# What can /proc tell us?

```
# cat /proc/diskstats
```

```
259      0 nvme0n1 1823879 1141588 58777929 328845 3512978 2244953 151315050
7397322 0 2759684 7859940 88819 0 543261592 2631 273795 107459
259      1 nvme0n1p1 1443 0 27072 2273 2 0 2 34 0 392 2310 4 0 1034184 1 0 0
259      2 nvme0n1p2 573 398 59794 77 367 115 151352 415 0 776 525 91 0
645136 33 0 0
259      3 nvme0n1p3 1820929 1141190 58667503 326372 3238845 2244838
151163696 7287293 0 2747016 7639943 88724 0 541582272 26277 0 0
```

Milliseconds spent in  
all IOs

Milliseconds spent in  
all disk reads

<sup>1</sup><https://www.kernel.org/doc/Documentation/iostats.txt>

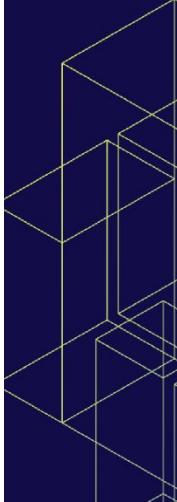
# What can /proc tell us?

```
# cat /proc/1475/io
```

rchar:	723810543
wchar:	1609216447277
syscr:	27156292
syscw:	142054696
<b>read_bytes:</b>	<b>30699520</b>
write_bytes:	6516736
cancelled_write_bytes:	5160960

Number of bytes  
requested to be read  
from block device by this  
process



<sup>1</sup><https://www.kernel.org/doc/Documentation/filesystems/proc.txt>

# Building tools over /proc and /sys

```
# iostat -d nvme0n1
```

```
Linux 5.7.1-050701-generic (isengard)      09/11/2020      _x86_64_      (8 CPU)
```

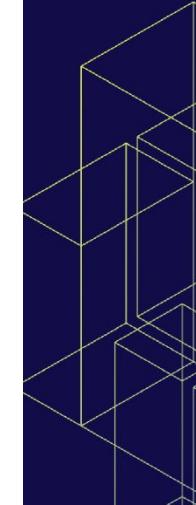
Device	tps	kB_read/s	kB_wrtn/s	kB_read	kB_wrtn
nvme0n1	7.61	42.62	109.59	30521808	78476133



Total kilobytes written  
per second

```
# iotop
```

*(Exercise left for the reader)*



# More Tracing Techniques

## Targeted Analysis

```
example-2535 [007] 1580764.896732: function:  
example-2535 [007] 1580764.896732: function:  
example-2535 [007] 1580764.896732: function:  
example-2535 [007] 1580764.896733: function:  
example-2535 [007] 1580764.896733: function:  
example-2535 [007] 1580764.896734: function:  
do_syscall_64  
    __x64_sys_read  
        ksys_read  
            __fdget_pos  
                __fget_light  
                    vfs_read
```

Syscall Enter

```
example-2535 [007] 1580764.896770: function:  
example-2535 [007] 1580764.896770: function:  
__wake_up_common_lock  
    __raw_spin_lock_irqsave
```

IRQ Entry

```
example-2535 [007] 1580764.896794: function:  
irq_exit
```

IRQ Exit

```
example-2535 [007] 1580764.896808: function:  
example-2535 [007] 1580764.896808: function:  
example-2535 [007] 1580764.896809: function:  
example-2535 [007] 1580764.896810: function:  
example-2535 [007] 1580764.896810: function:  
__mnt_drop_write  
    __sb_end_write  
        fsnotify_parent  
            fsnotify  
                tpregs_assert_state_consistent
```

Syscall Exit

<sup>1</sup><https://github.com/iovisor/bpftrace/tree/master/tools>

# More Tracing Techniques

## Targeted Analysis



```
example-2535 [007] 1580764.896732: function:  
example-2535 [007] 1580764.896732: function:  
example-2535 [007] 1580764.896732: function:  
example-2535 [007] 1580764.896733: function:  
example-2535 [007] 1580764.896733: function:  
example-2535 [007] 1580764.896734: function:
```

```
do_syscall_64  
__x64_sys_read  
ksys_read  
__fget_pos  
__fget_light  
vfs_read
```

**Syscall Enter**

```
example-2535 [007] 1580764.896770: function:  
example-2535 [007] 1580764.896770: function:
```

```
_wake_up_common_lock  
_raw_spin_lock_irqsave
```

**IRQ Entry**

```
example-2535 [007] 1580764.896794: function:
```

irq\_exit

**IRQ Exit**

```
example-2535 [007] 1580764.896800: function:  
example-2535 [007] 1580764.896800: function:  
example-2535 [007] 1580764.896809: function:  
example-2535 [007] 1580764.896810: function:  
example-2535 [007] 1580764.896810: function:
```

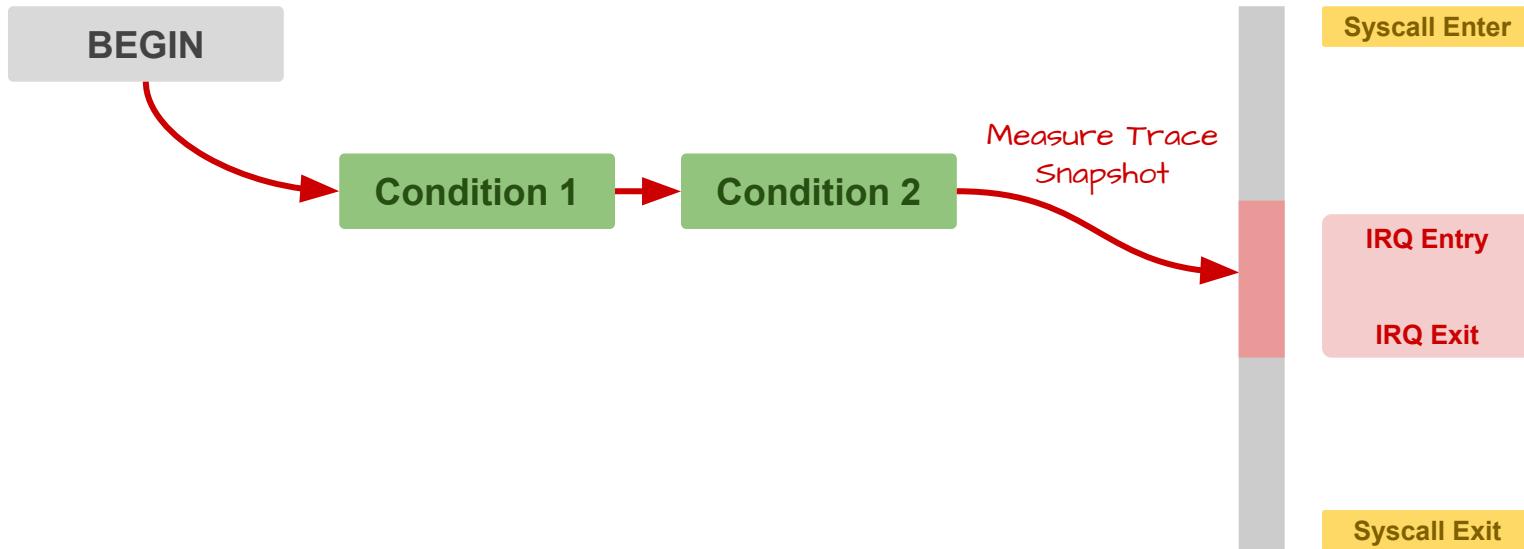
```
_mnt_drop_write  
_sb_end_write  
fsnotify_parent  
fsnotify  
tfregs_assert_state_consistent
```

**Syscall Exit**

<sup>1</sup><https://github.com/iovisor/bpftrace/tree/master/tools>

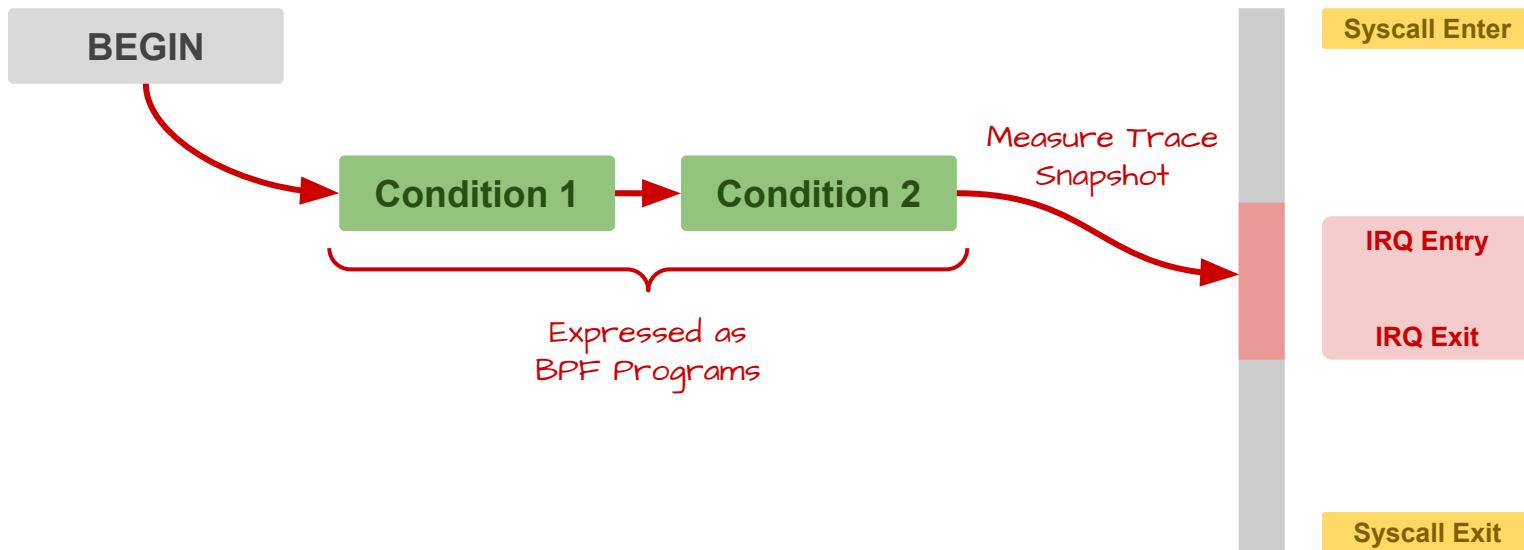
# More Tracing Techniques

## Live and Programmatic Analysis



# More Tracing Techniques

## Live and Programmatic Analysis



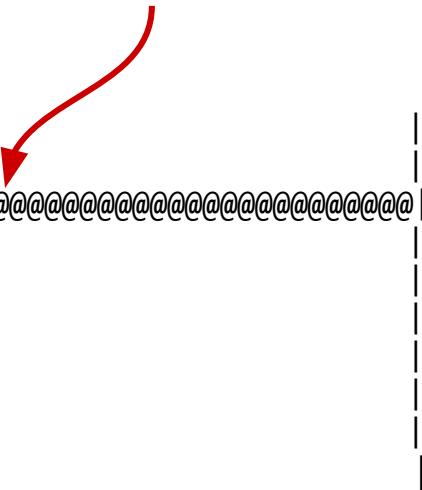
# Tracing with BPF

# bpftrace Tools

```
# ./biolatency.bt
Attaching 3 probes...
Tracing block device I/O... Hit Ctrl-C to end.
^C
```

Most block operations took between 1024 to 2048 micro-seconds

@usecs:		
[256, 512)	2	
[512, 1K)	10	@
[1K, 2K)	426	@aaaaaaaaaaaaaaaaaaaaaa
[2K, 4K)	230	@aaaaaaaaaaaaaaaaaaaaaa
[4K, 8K)	9	@
[8K, 16K)	128	@aaaaaaaaaaaaaaaaaaaaaa
[16K, 32K)	68	@aaaaaaaaaaaa
[32K, 64K)	0	
[64K, 128K)	0	
[128K, 256K)	10	@



<sup>1</sup><https://github.com/iovisor/bpftrace/tree/master/tools>

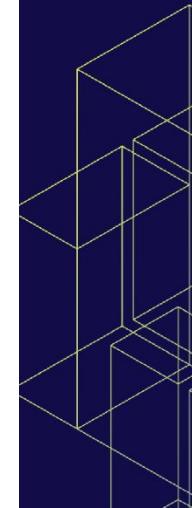
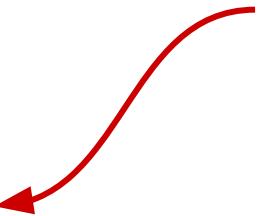
# Tracing with BPF

## bpftrace Tools

```
# ./biosnoop.bt  
Attaching 4 probes...
```

TIME(ms)	COMM	PID	LAT(ms)
611	bash	4179	10
611	cksum	4179	0
627	cksum	4179	15
<b>641</b>	<b>cksum</b>	<b>4179</b>	<b>13</b>
644	cksum	4179	3
658	cksum	4179	13
673	cksum	4179	14
686	cksum	4179	13

Each block event per process with block operation latency



<sup>1</sup><https://github.com/iovisor/bpftrace/tree/master/tools>

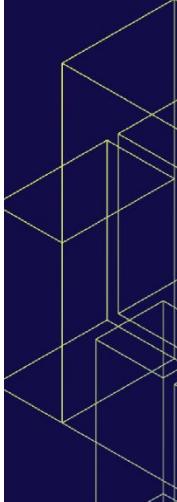
# Tracing with BPF

## bpftrace Tools

- biostack.bt
- bitesize.bt
- dcsnoop.bt
- mdflush.bt

## BCC Tools

- biostat.py
- bitesize.py
- bitehist.py
- biotop.py



<sup>1</sup><https://github.com/iovisor/bpftrace/tree/master/tools>

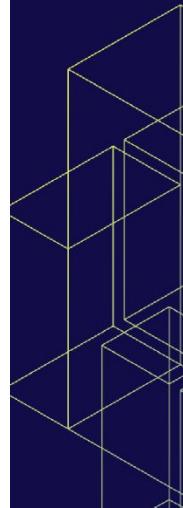


# Building tools with



# eBPF for Observability

- **Probe the OS, Apps, Network**
  - Single programmatic way of tracing and monitoring each and every aspect of the software infrastructure
  - Hook to any userspace or kernel function\* or predefined tracepoints, then extract and visualize data
  - Enable and disable probes, live filtering, low runtime overhead
- **Workflow**
  - Create your observation tool in userspace
  - Compile it and send it to kernel and attach to hooks
  - It runs and returns data back via buffer or maps



# eBPF for Observability



# eBPF for Observability

TOOL

eBPF Program

Control

USERSPACE

KERNEL

TARGET APP



read()

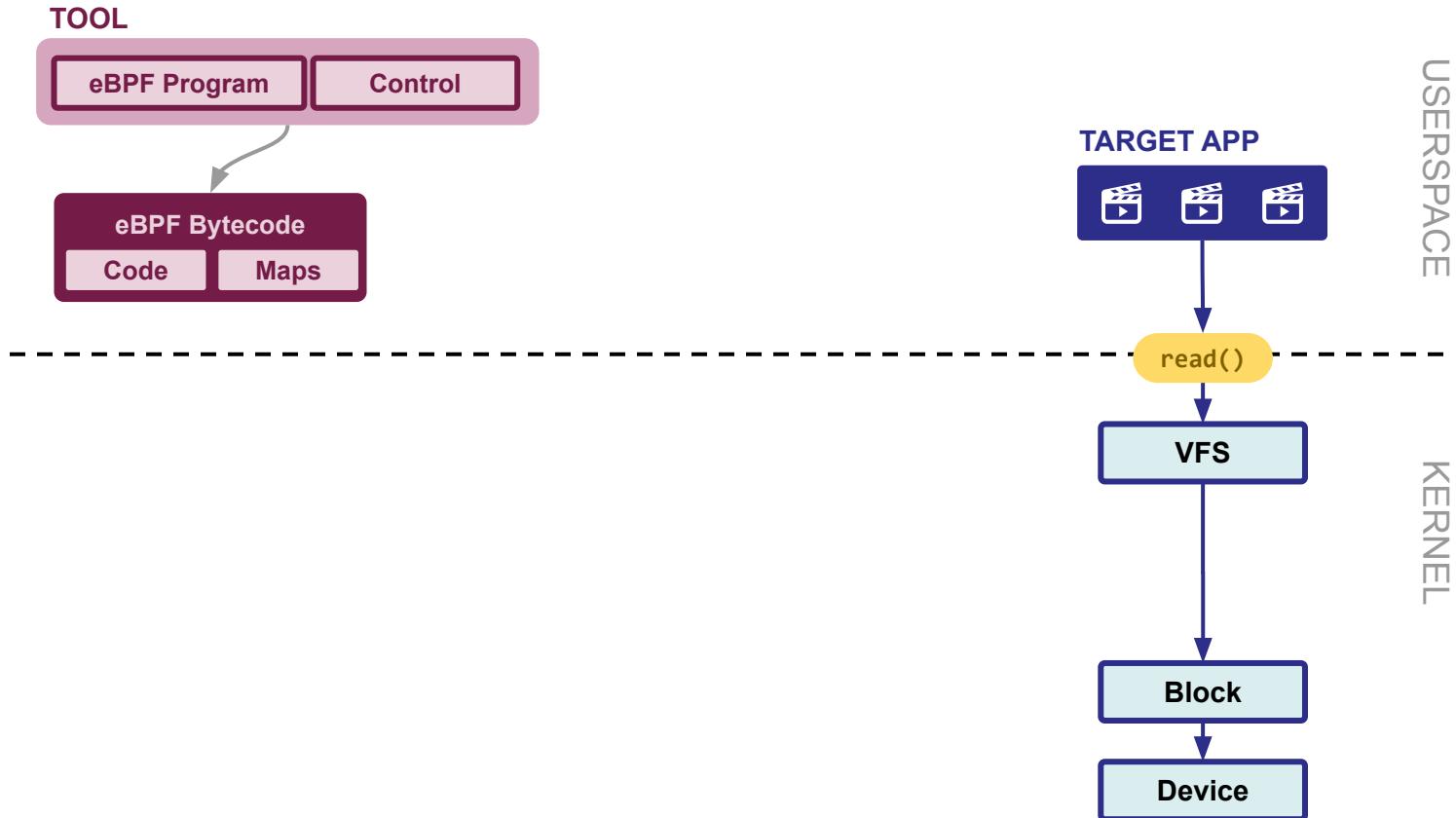


VFS

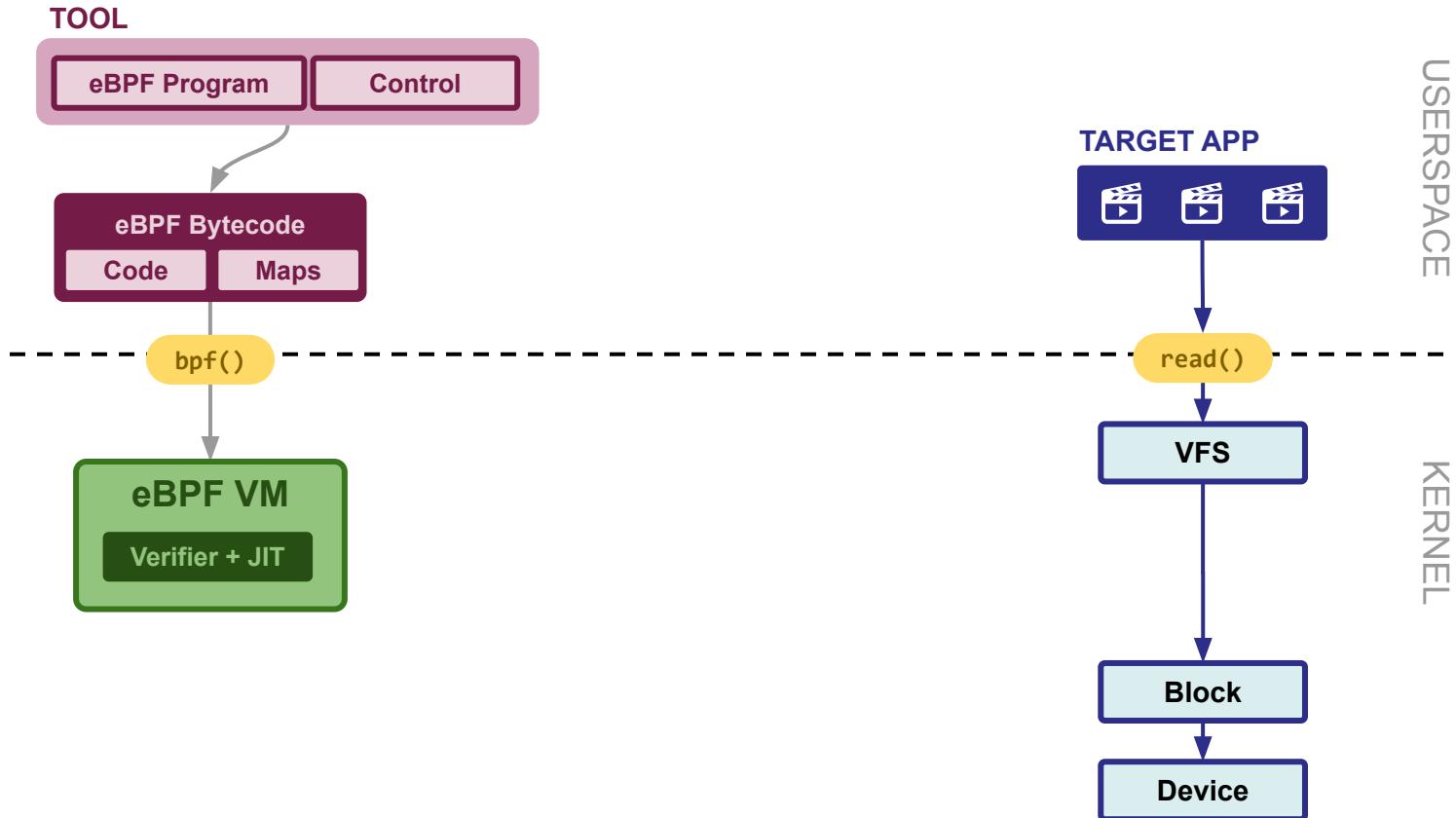
Block

Device

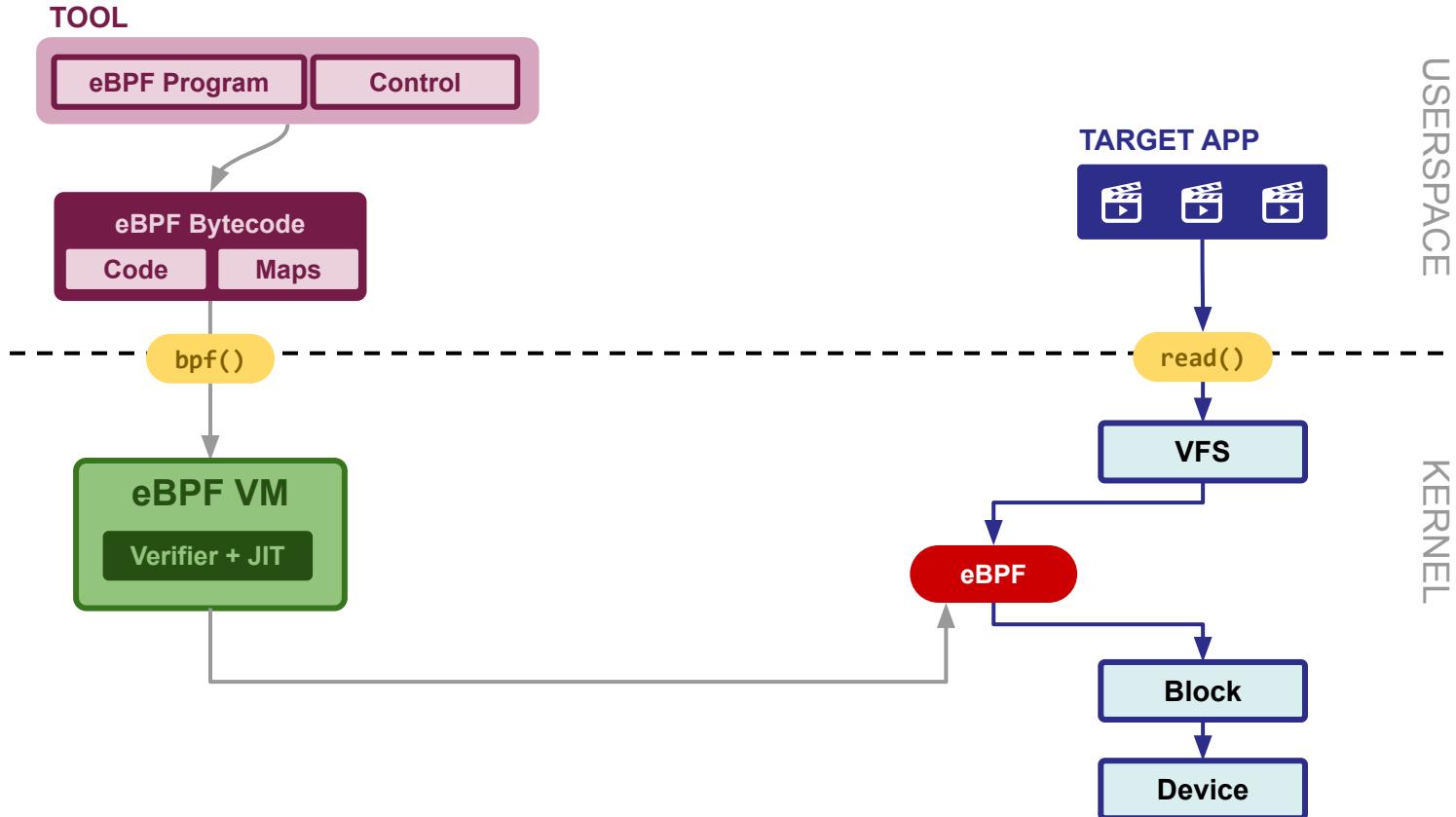
# eBPF for Observability



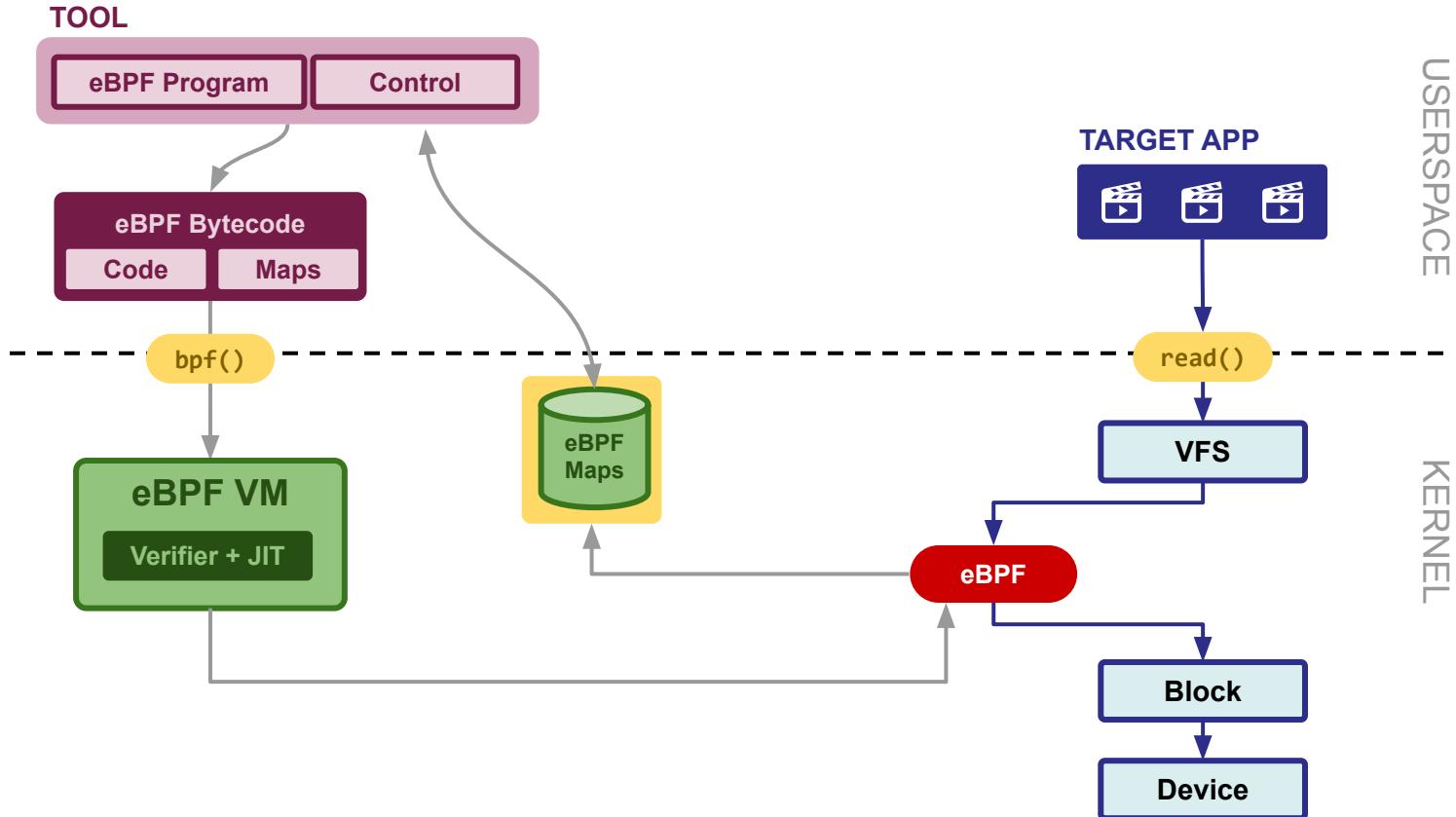
# eBPF for Observability



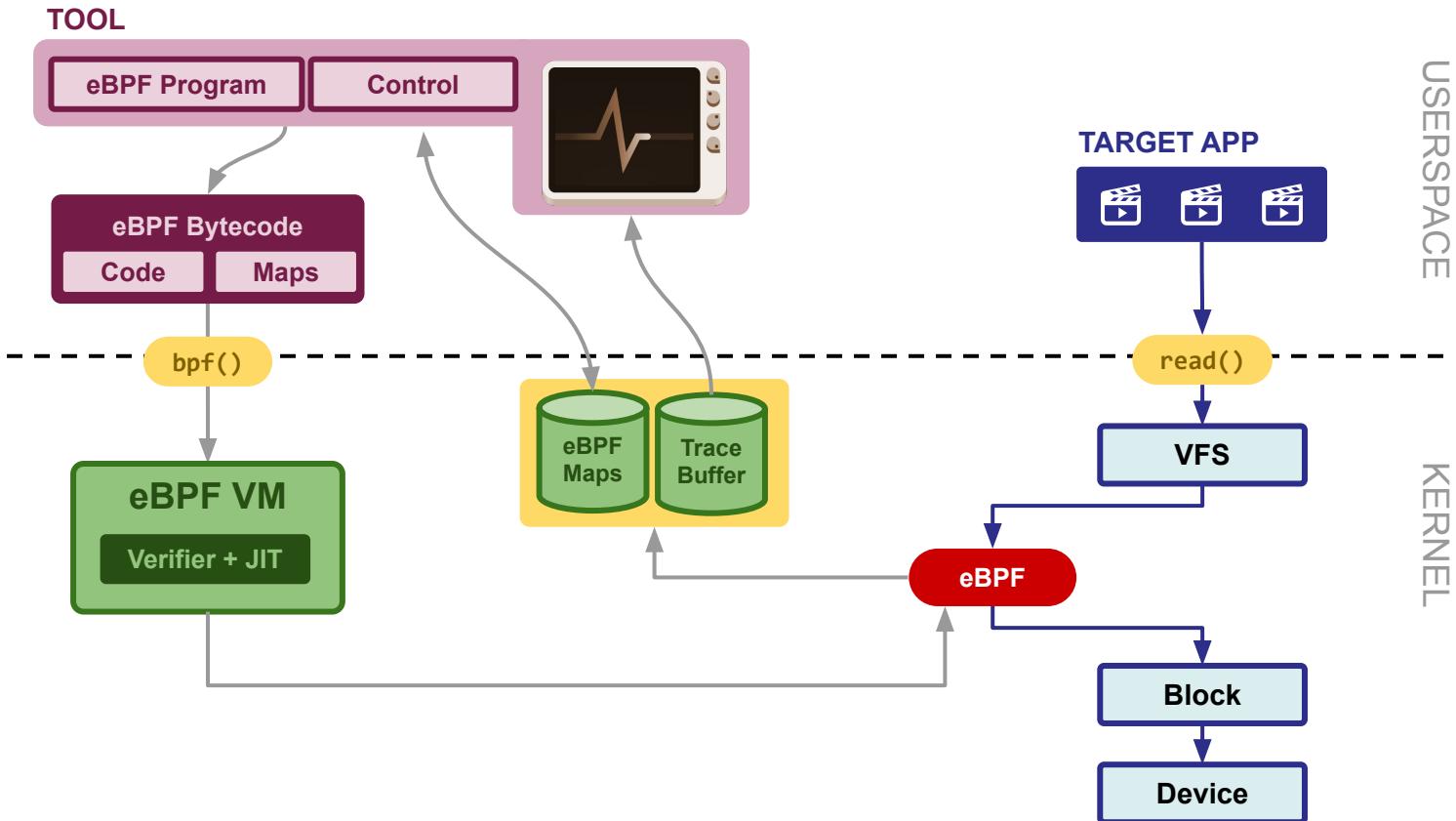
# eBPF for Observability



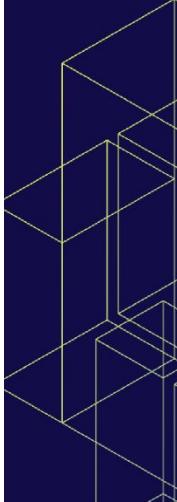
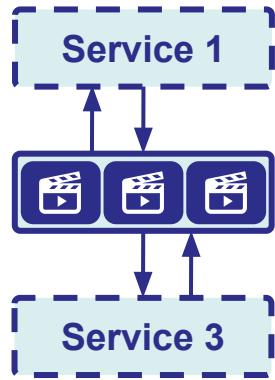
# eBPF for Observability



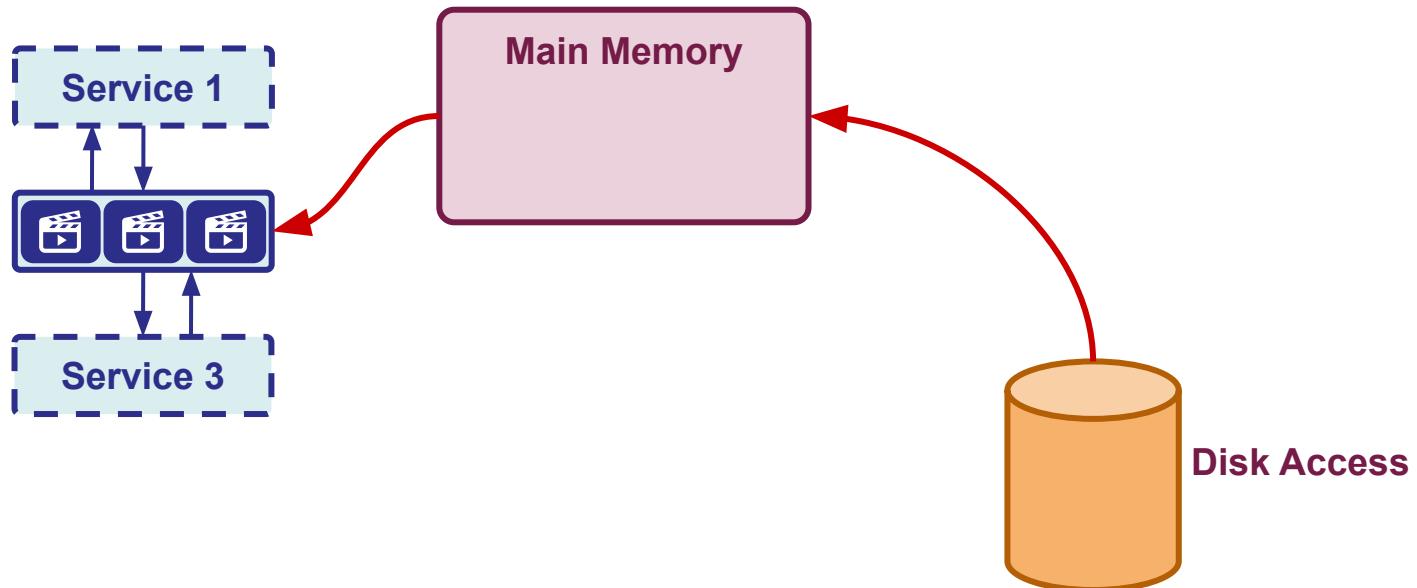
# eBPF for Observability



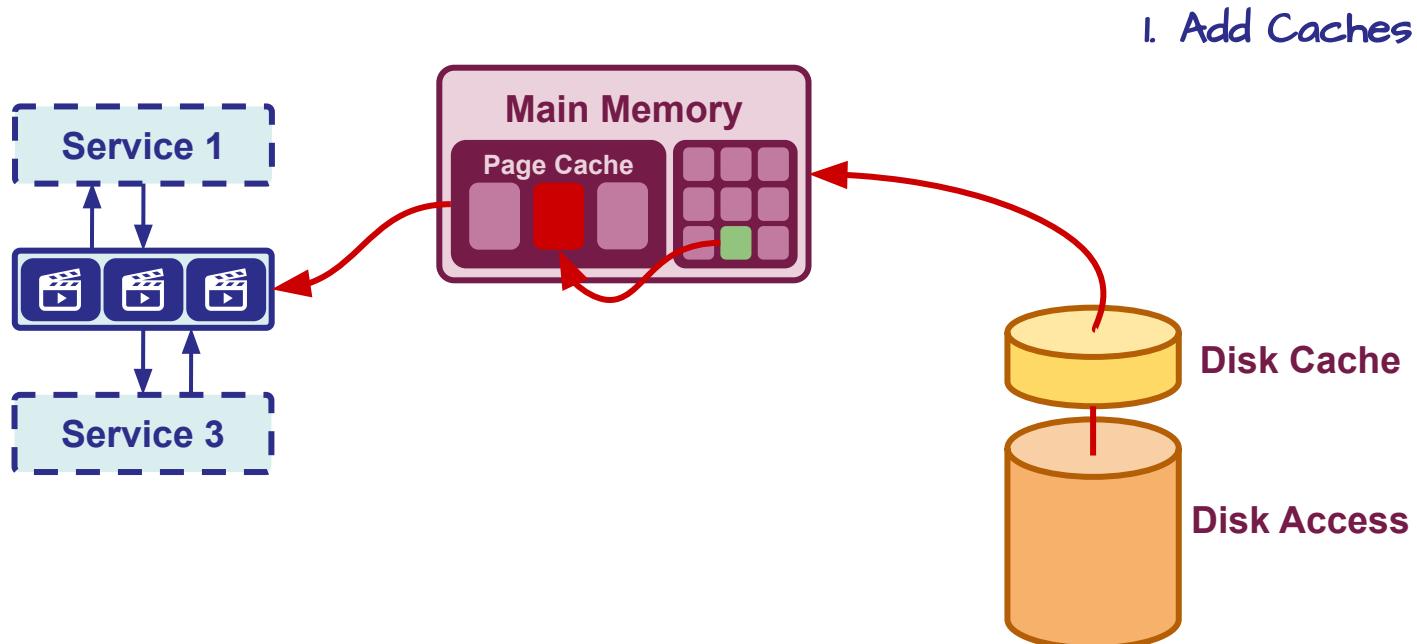
# Case Study: Read-ahead



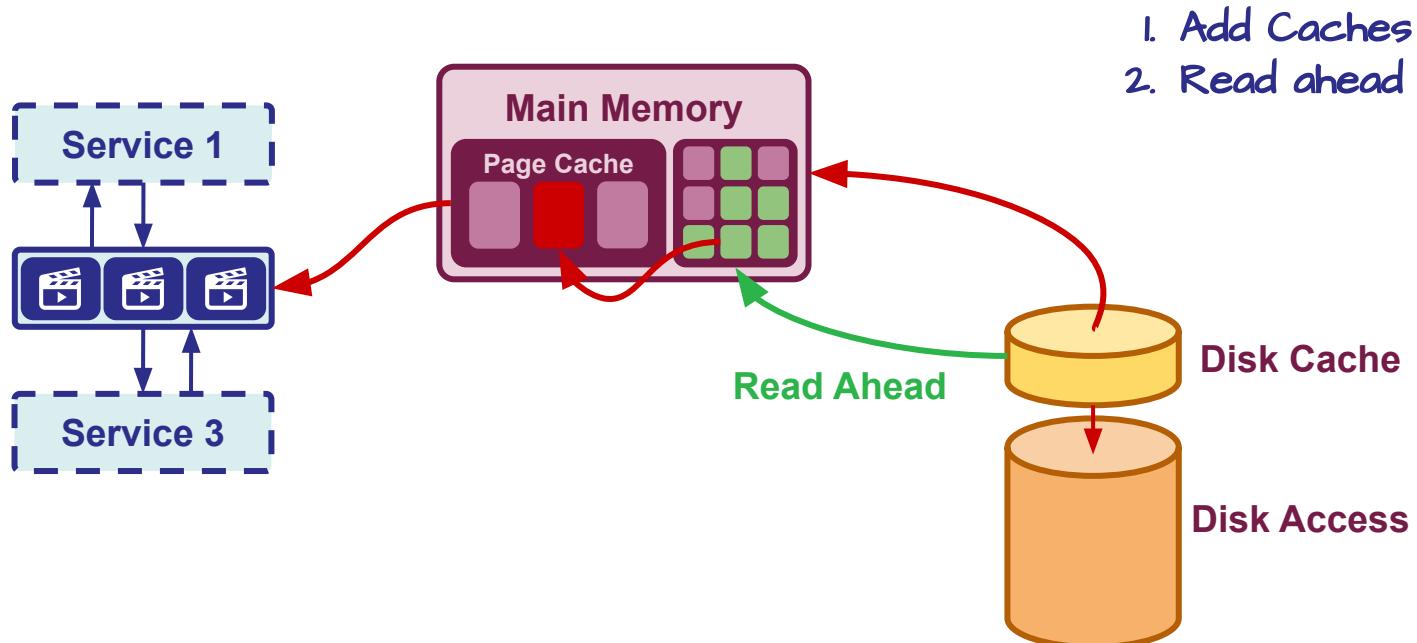
# Case Study: Read-ahead



# Case Study: Read-ahead

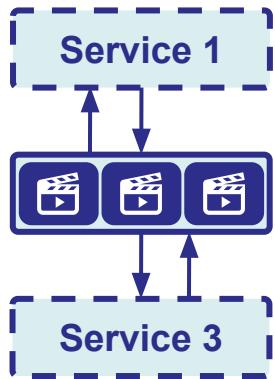


# Case Study: Read-ahead

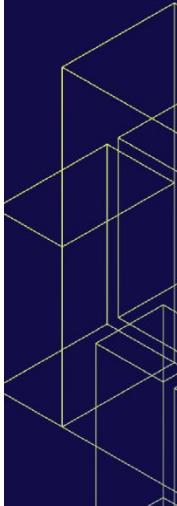


# Case Study: Read-ahead

## How?

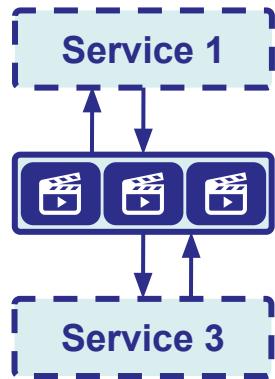


- Apps can “advise” the OS to use efficient read-ahead mechanism for sequential reads (like the one here)
  - `madvise (addr, len, MADV_SEQUENTIAL)`
- OS enables an aggressive read-ahead of pages in memory since it now know faults will be less due to seq read.

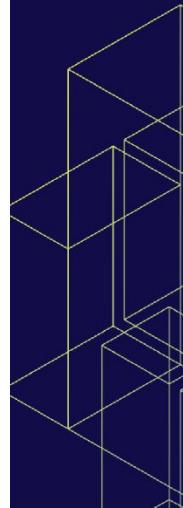


# Case Study: Read-ahead

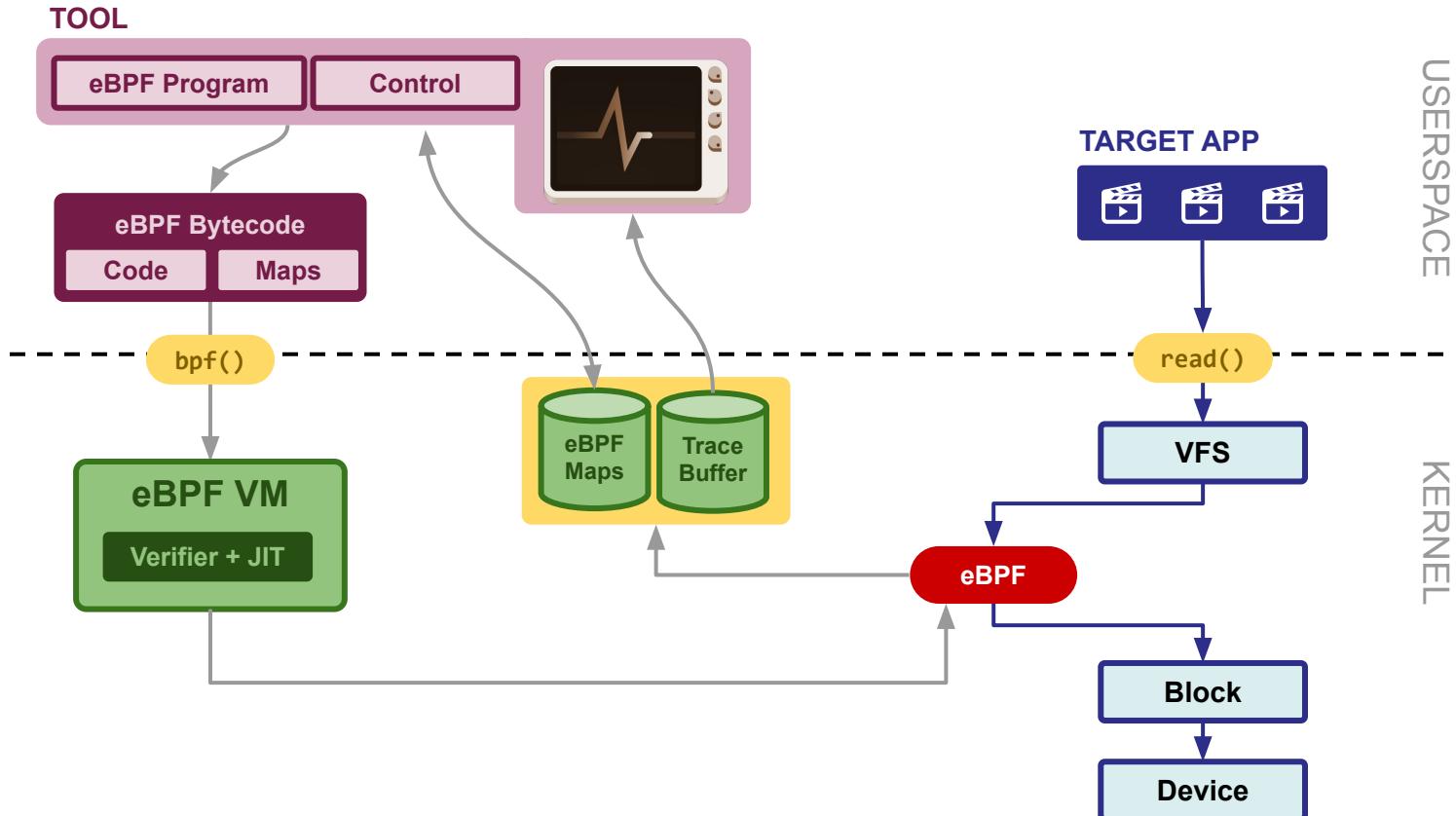
## Problems



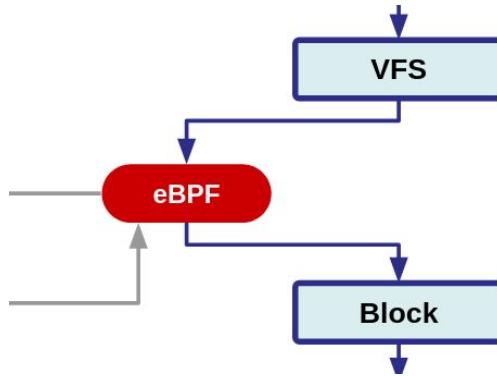
- Do apps really know how to use it?
- Modern apps are unaware of underlying infrastructure
- Lack of insight on how the different hardwares/OS optimize this really
- Has been a lack of **dynamic and deep-observability** tools to “see” the problems in real time



# eBPF for Observability

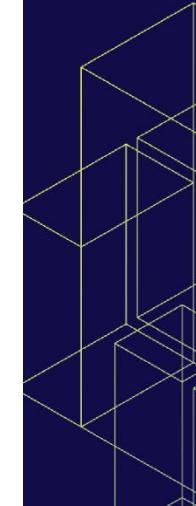
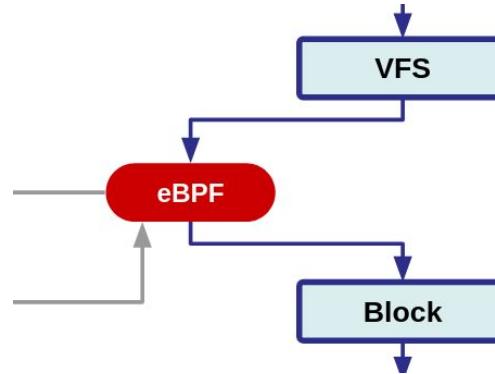


# Building readaheadstat



# Building readaheadstat

- Where and what are the probes?
  - Dynamic:*
    - Kprobes/Kretprobes (Kernel)**
    - USDT/Uprobes (Userspace)
  - Static:* Tracepoints



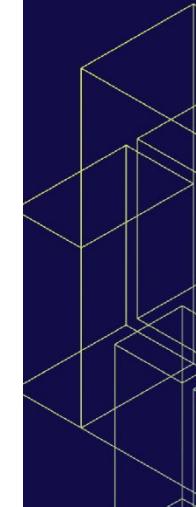
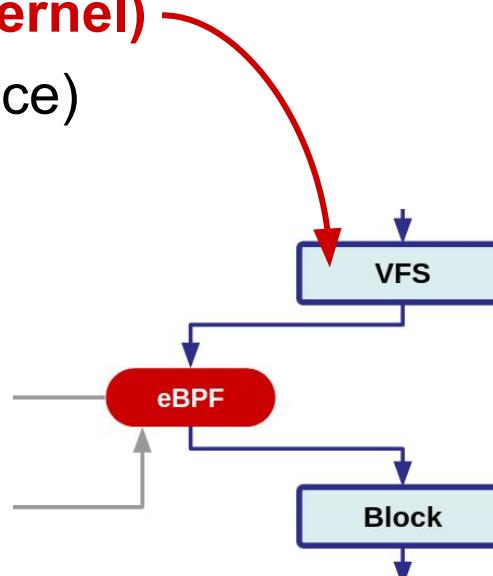
# Building readaheadstat

- Where and what are the probes?

- Dynamic:*

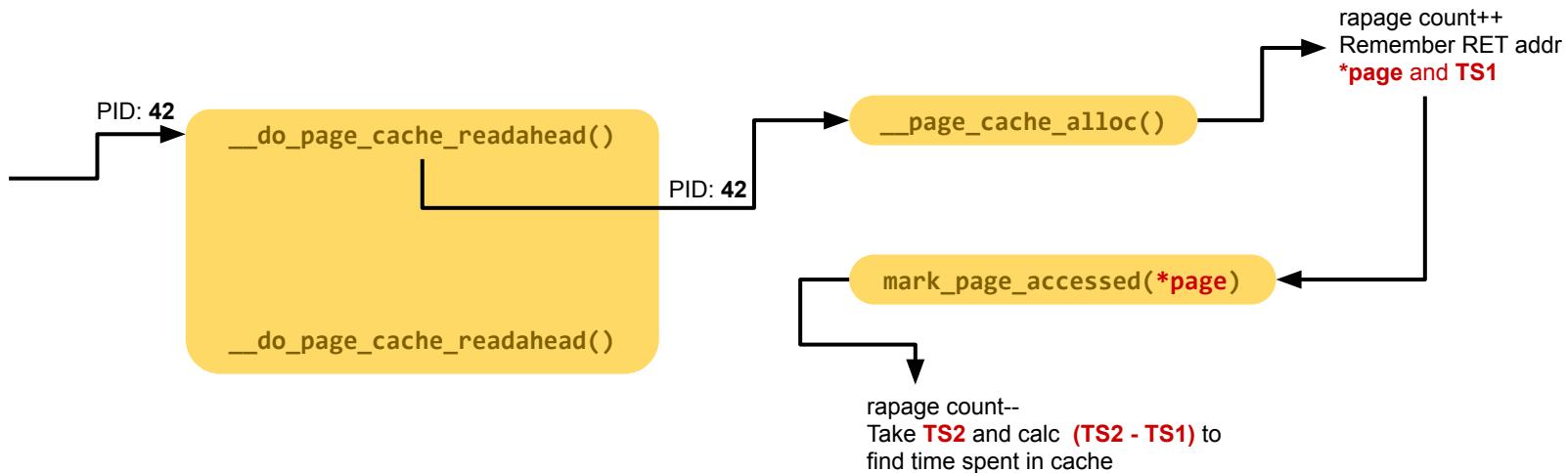
- Kprobes/Kretprobes (Kernel)**
- USDT/Uprobes (Userspace)

- Static:* Tracepoints



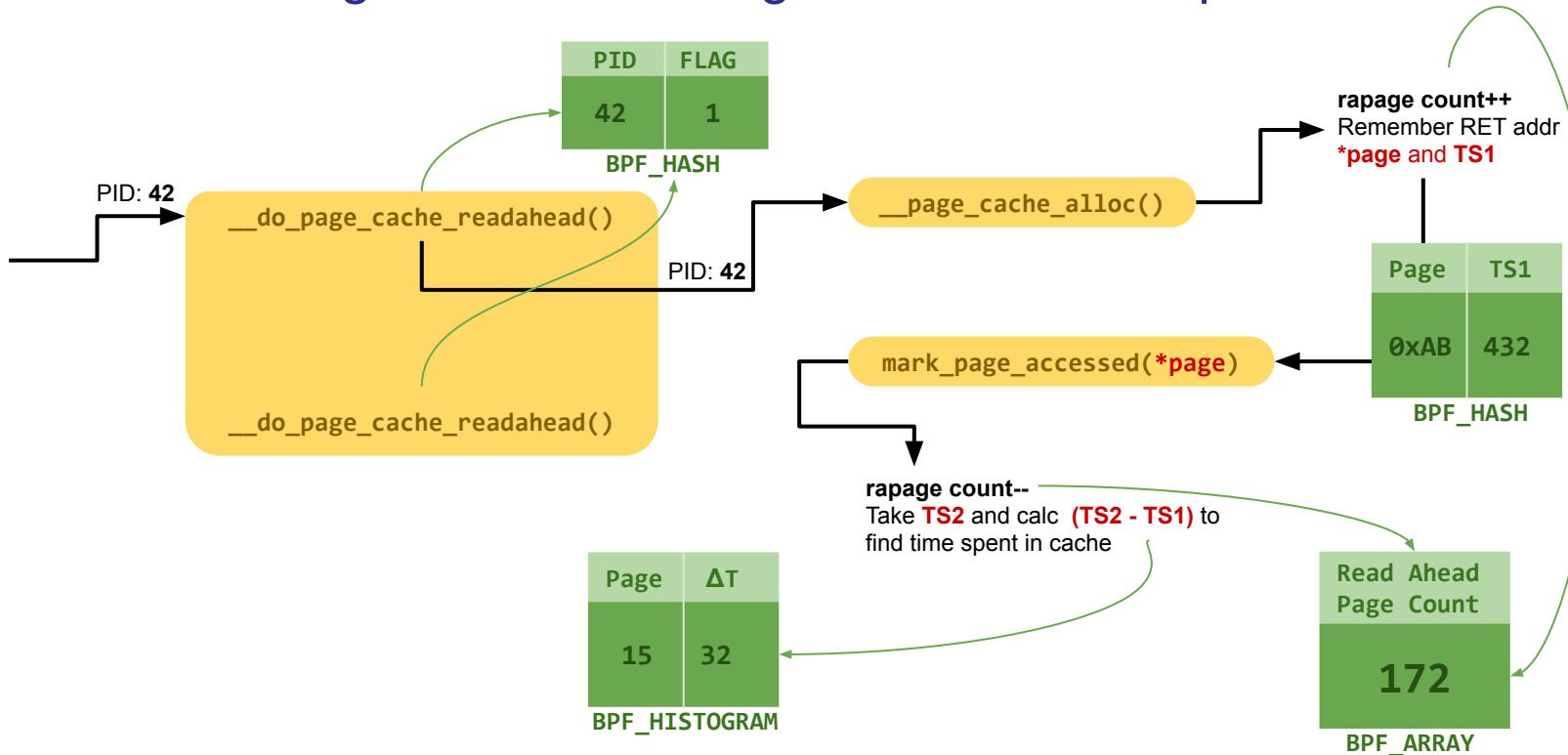
# Building readaheadstat

- Selecting Kprobes/Kretprobes



# Building readaheadstat

- Maintaining state and storing data in hash-maps





# readaheadstat Demo

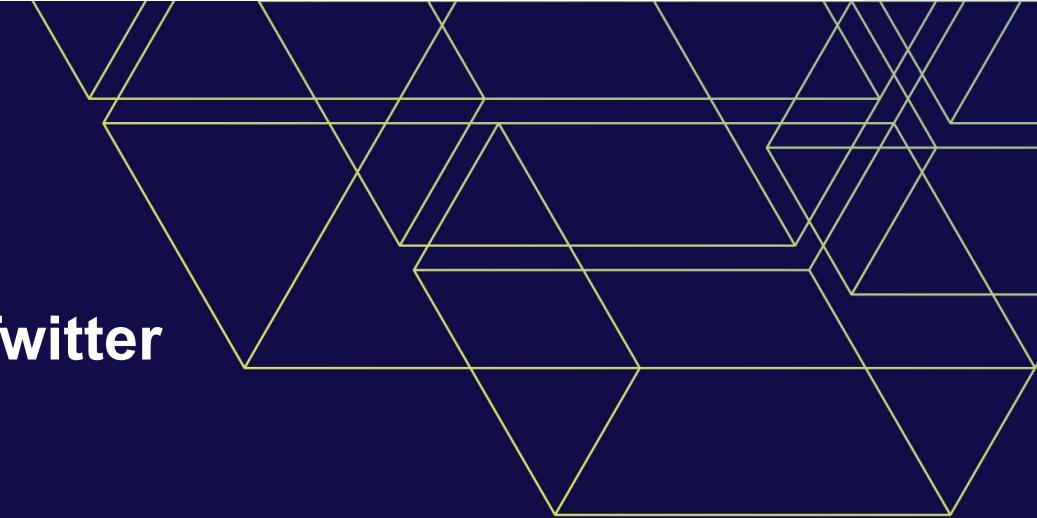


Download and Try

<https://github.com/tuxology/readaheadstat>

# Building readaheadstat





**You can reach out to us on Twitter**

[@tuxology](https://twitter.com/tuxology)

[@hani\\_nemati](https://twitter.com/hani_nemati)

**Or e-mail us**

[suchakra@shiftleft.io](mailto:suchakra@shiftleft.io)

[hanemati@microsoft.com](mailto:hanemati@microsoft.com)