# BPF for chaos and tracing in Kubernetes

*Wenbo Zhang*

# **About the Presenter**

- Wenbo Zhang
  - A PingCAP Development Engineer, focusing on performance analysis and diagnosis of Linux kernel
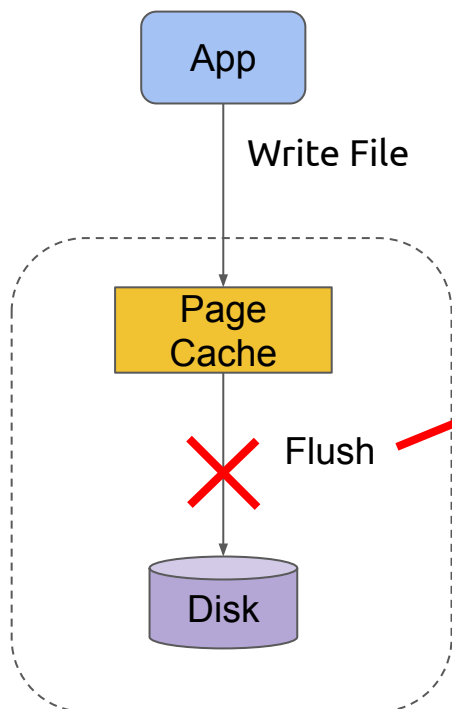
# Agenda

- Why we need Chaos
- Kernel Chaos with BPF
- BPF with native support for containers
- Tracing with BPF

# Why we need Chaos

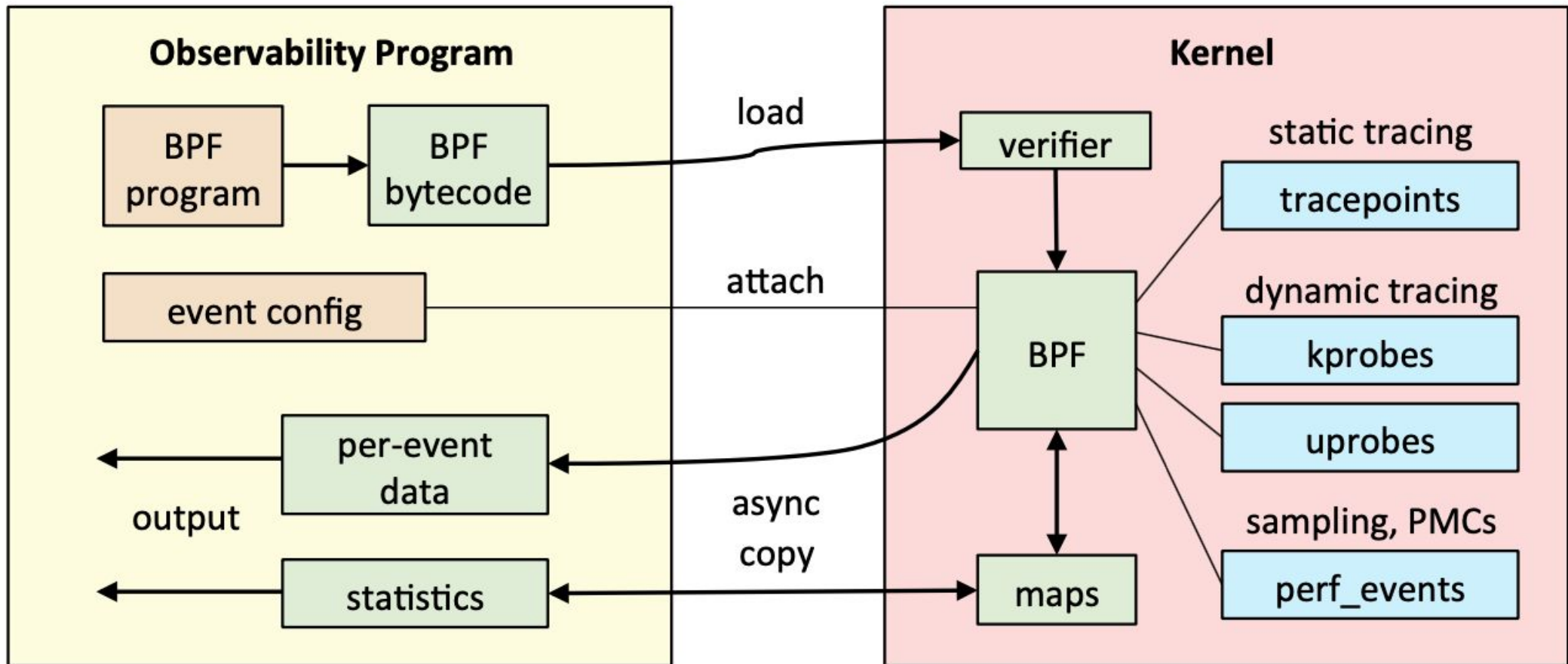- Error happens, any types, anytime, anywhere, any device



```
[17988717.953807] 0 pages in swap cache
[17988717.953808] Swap cache stats: add 0, delete 0, find 0/0
[17988717.953809] Free swap  = 0kB
[17988717.953810] Total swap = 0kB
[17988717.953811] SLUB: Unable to allocate memory on node -1 (gfp=0x20)
[17988717.953813]   cache: kmalloc-8192, object size: 8192, buffer size: 8192,
[17988717.953815]   node 0: slabs: 78, objs: 312, free: 21
[17988717.953816]   node 1: slabs: 37, objs: 148, free: 0
```

Linux kernel bugs happen in real world!!!

https://pingcap.com/blog/try-to-fix-two-linux-kernel-bugs-while-testing-tidb-operator-in-k8s/

# Kernel Chaos With BPF

- Enhanced BPF Tracing Internals

# Kernel Chaos With BPF

- ## Disadvantages of fault injection framework
  - Just making kmalloc() fail universally is unlikely to be helpful
  - The parameters control mechanism is somewhat awkward to use and is not as flexible as one might like

- ## BPF override return
  - Fault injection for specific paths
  - Support precise filtering
  - Never crash the kernel
    - ALLOW_ERROR_INJECTION
    - override function should only change integer error values

# Examples

- Inject congestion wait on a special task's fsync path

# Examples

- Inject alloc inode failure on a special open file path

# Examples

● Inject alloc inode failure on a special symlink file path

# Examples

- Accuracy depends on your familiarity with kernel code
    - eg (from bcc):

*struct disk_part_tbl *tbl = d->part_tbl;*

*struct hd_struct **parts = (void *)tbl + sizeof(struct disk_part_tbl);*

*struct hd_struct **partp = parts + bio->bi_partno;*

*struct hd_struct *p = *partp;*

*dev_t disk = p->__dev.devt;*

**disk == MKDEV(254,16) && bio->bi_iter.bi_sector == 128**

# Where to start

- Injection of syscalls is a good entry point
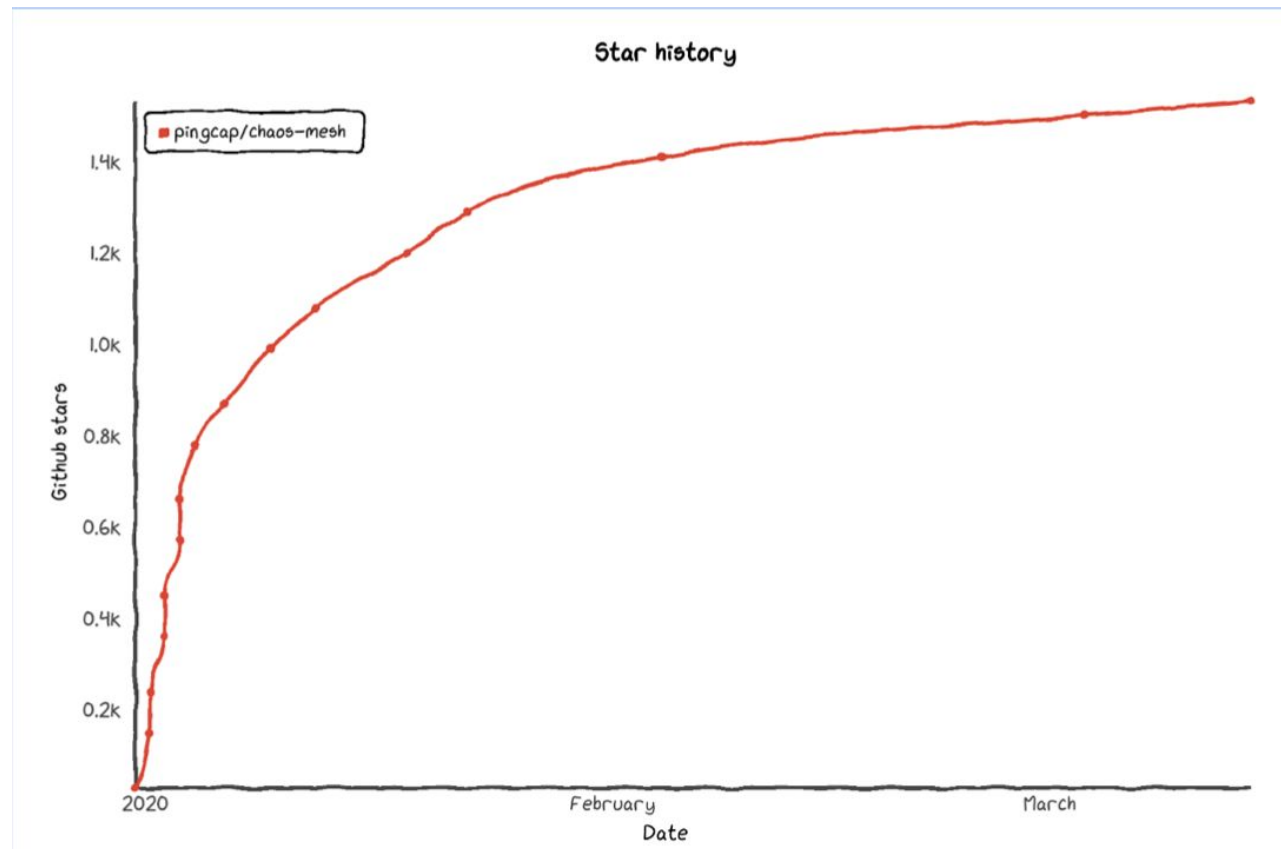
# Where to start

- You can try kernel chaos on

Chaos Mesh®

# Chaos Mesh®

- ## Community and Ecosystem

  - 1.9k stars on Github
  - 41 contributors
    - PingCAP
    - Red Hat
    - DailyMotion
    - Nvidia
    - Xpeng Motors
    - China Merchants Bank
    - Meituan Dianping
  - 400+ commits
  - Monthly meeting



Star history

# BPF With Native Support For Containers

- ## Challenges
  - PIDs in host don't match those seen in containers
  - Kernel currently doesn't have a container ID

- ## Solutions
  - If in process context, we can read nsproxy struct in the kernel
    - *(struct task_struct*)task->nsproxy->pid_ns_for_children->ns.inum*
    - *(struct task_struct*)task->nsproxy->uts_ns->name.nodename*

```
struct nsproxy {
    atomic_t count;
    struct uts_namespace *uts_ns;
    struct ipc_namespace *ipc_ns;
    struct mnt_namespace *mnt_ns;
    struct pid_namespace *pid_ns_for_children;
    struct net           *net_ns;
    struct time_namespace *time_ns;
    struct time_namespace *time_ns_for_children;
    struct cgroup_namespace *cgroup_ns;
};
```

```
→ ns pwd
/proc/1/ns
→ ns ll
total 0
lrwxrwxrwx 1 root root 0 Jul  7 22:35 cgroup -> 'cgroup:[4026531835]'
lrwxrwxrwx 1 root root 0 Jul  7 22:35 ipc -> 'ipc:[4026531839]'
lrwxrwxrwx 1 root root 0 Jul  3 03:57 mnt -> 'mnt:[4026531840]'
lrwxrwxrwx 1 root root 0 Jul  7 22:35 net -> 'net:[4026532008]'
lrwxrwxrwx 1 root root 0 Jul  7 08:22 pid -> 'pid:[4026531836]'
lrwxrwxrwx 1 root root 0 Jul  7 22:35 pid_for_children -> 'pid:[4026531836]'
lrwxrwxrwx 1 root root 0 Jul  7 22:35 time -> 'time:[4026531834]'
lrwxrwxrwx 1 root root 0 Jul  7 22:35 time_for_children -> 'time:[4026531834]'
lrwxrwxrwx 1 root root 0 Jul  7 22:35 user -> 'user:[4026531837]'
lrwxrwxrwx 1 root root 0 Jul  7 22:35 uts -> 'uts:[4026531838]'
```
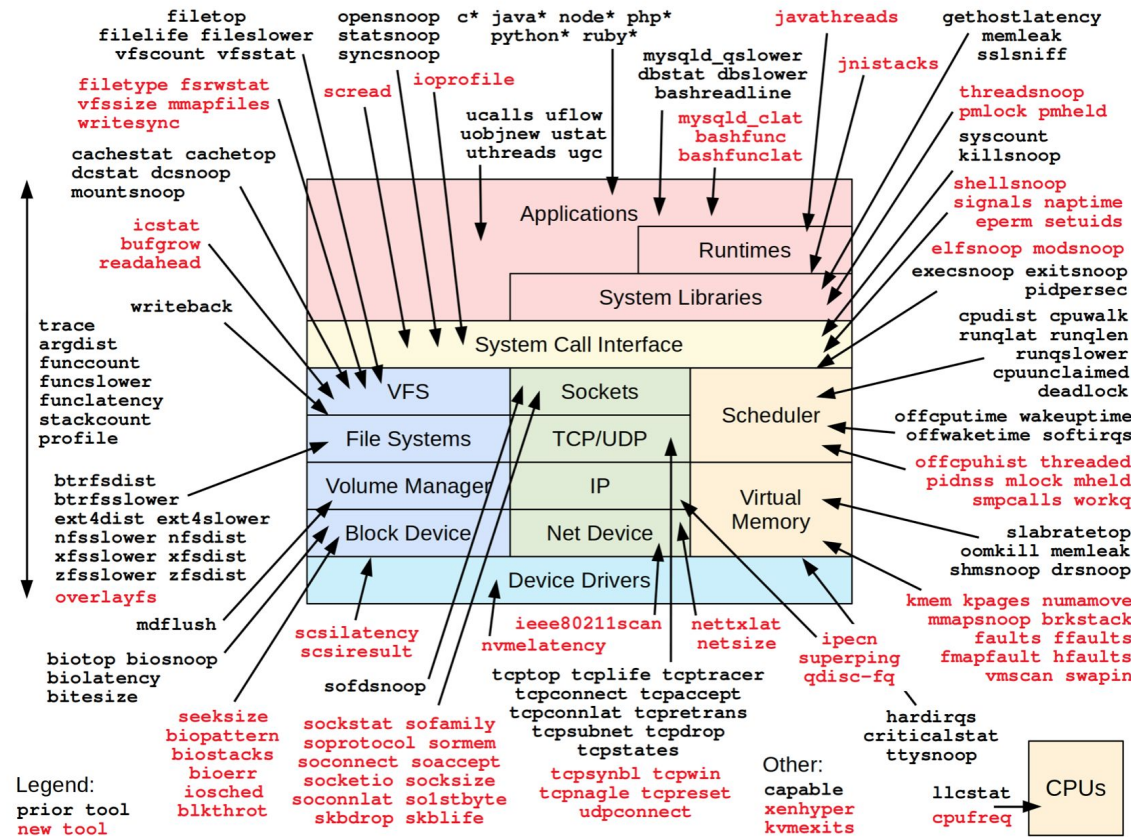
● Example

```
kprobe:finish_task_switch
{
    $prev = (struct task_struct *)arg0;
    $curr = (struct task_struct *)curtask;
    $prev_pidns = $prev->nsproxy->pid_ns_for_children->ns.inum;
    $curr_pidns = $curr->nsproxy->pid_ns_for_children->ns.inum;
    if ($prev_pidns != $curr_pidns) {
        @[$prev_pidns, $prev->nsproxy->uts_ns->name.nodename] = count();
    }
}
```

CLOUD NATIVE + OPEN SOURCE
Virtual Summit China 2020

- ## BPF Tools
  - ○ Over 150+ that you can run to find performance wins and troubleshoot software
  - ○ With kubectl-trace's help bpftrace progs (already a pr to support bcc progs) can be scheduled in Kubernetes cluster

# K8s tracing with BPF

- ## BCC drawbacks
  - Clang/LLVM combo is a big library, resulting in big fat binaries that need to be distributed with your application
  - Clang/LLVM combo is resource-heavy, so when you are compiling BPF code at start up, you'll use a significant amount of resources, potentially tipping over a carefully balanced production workload. And vice versa, on a busy host, compiling a small BPF program might take minutes in some cases.
  - You are making a big bet that the target system will have kernel headers present, which most of the time is not a problem, but sometimes can cause a lot of headaches. This is also an especially annoying requirement for kernel developers, because they often have to build and deploy custom one-off kernels as part of their development process. And without a custom-built kernel header package, no BCC-based application will work on such kernels, stripping developers of a useful set of tools for debugging and monitoring.
  - BPF program testing and development iteration is quite painful as well, as you are going to get even most trivial compilation errors only in runtime, once you recompile and restart your user-space control application. This certainly increases friction and is not helping to iterate fast.

# Tracing with BPF

- Trace node with [kubectl-trace](#)

- Trace pod with kubectl-trace

- Container-Specific Tools
  - runqlat --pidnss -m
  - pidnss
  - blkthrot
  - overlayfs

# Thank you!