

# oomd

**Daniel Xu**

Production Engineer

# Out of memory daemon

- User-space
- Dependency free
- Deterministic, faster, and more flexible than kernel OOM killer
- Open source (GPL2)
  - <https://github.com/facebookincubator/oomd>
  - <https://facebookmicrosites.github.io/oomd/>

# Agenda

- Motivation
- Mechanism
- Results
- Questions

# Motivation

# Kernel OOM killer

- Memory overcommit
- Configuration not very intuitive (what's with all the numbers?)
  - */proc/[pid]/oom\_adj*
  - */proc/[pid]/oom\_score*
  - */proc/[pid]/oom\_score\_adj*
  - */proc/sys/vm/oom\_kill\_allocating\_task*
  - */proc/sys/vm/panic\_on\_oom*
- Slow to act; often it's already too late by the time the kernel reacts
  - Tries to protect kernel health; user-space could be livelocked but the kernel could still be happily churning pages in and out

# Kernel OOM killer (cont)

- No context on logical composition of system
  - What should be killed together, what shouldn't be, etc.
- No way to customize kill action (modulo OOM eventfd; still slow though)
  - For some processes, a SIGTERM/SIGKILL is fine. Other might want a song and dance
    - Eg. a persistent process that manages containers
- Non-deterministic (Or at least really hard to get deterministic)

# OOM problems @ FB

- Continuous build and test platform
- Container & service platform (Tupperware)
- White box (commodity) top-of-rack switches (FBOSS)
- Other multi-tenant platforms
- *vm.panic\_on\_oom*

# Fbtax2

- See "Resource Control @ FB" talk
  - Work-conserving full-OS resource isolation across applications
  - Uses cgroup2, PSI, io.latency, btrfs, and oomd to isolate workloads on the same machine
    - <https://facebookmicrosites.github.io/cgroup2>
    - <http://git.cmpxchg.org/cgit.cgi/linux-psi.git/>
    - <https://facebookmicrosites.github.io/btrfs>



Mechanism

# How does it work?

- PSI
  - New kernel feature (by FB's Johannes Weiner) that quantifies lost wall clock time due to resource shortages
- Core of oomd is the plugin system
  - All detection and action behavior is customizable
  - Default OomDetector and OomKiller plugins are sensible and work well across a variety of workloads
- Not just memory, also monitors IO pressure and swap depletion
  - There to remediate when kernel resource isolation isn't enough

```
1 {
2   "target": "system.slice",
3   "kill_list": [
4     {"chef.service": { "max_usage": "1000" } },
5     {"sshd.service": { "max_usage": "inf" } },
6   ],
7   "oomdetector": "default",
8   "oomkiller": "default"
9 }
```

# oomd2

- Oomd as it currently exists works well
- More users, more use cases, new oomd design
- Still iterating and playing around with details
  - In general, we're onto something very useful here

```
1 {
2   "rulesets": [
3     <FIRE IF WORKLOAD.SLICE SLOWS BY > 5%>,
4     <FIRE IF SYSTEM.SLICE SLOWS BY > 40%>
5   ],
6   "actions": [
7     <KILL A MEMORY HOG IN SYSTEM.SLICE>
8   ]
9 }
10
```

 pseudocode

```
1  {
2    "rulesets": [
3      [
4        "name": "my first ruleset",
5        "detectors": [
6          [
7            "group1",
8            [
9              "pressure_rising_beyond",
10             "cgroup=workload.slice",
11             "resource=memory",
12             "threshold=5"
13           ],
14           [
15             "pressure_rising_beyond",
16             "cgroup=system.slice",
17             "resource=memory",
18             "threshold=40"
19           ]
20         ]
21       ],
22       "actions": [
23         [
24           "kill_by_memory_size_or_growth",
25           "cgroup=system.slice"
26         ]
27       ]
28     ]
29   ]
30 }
```

oomd2 config (real)

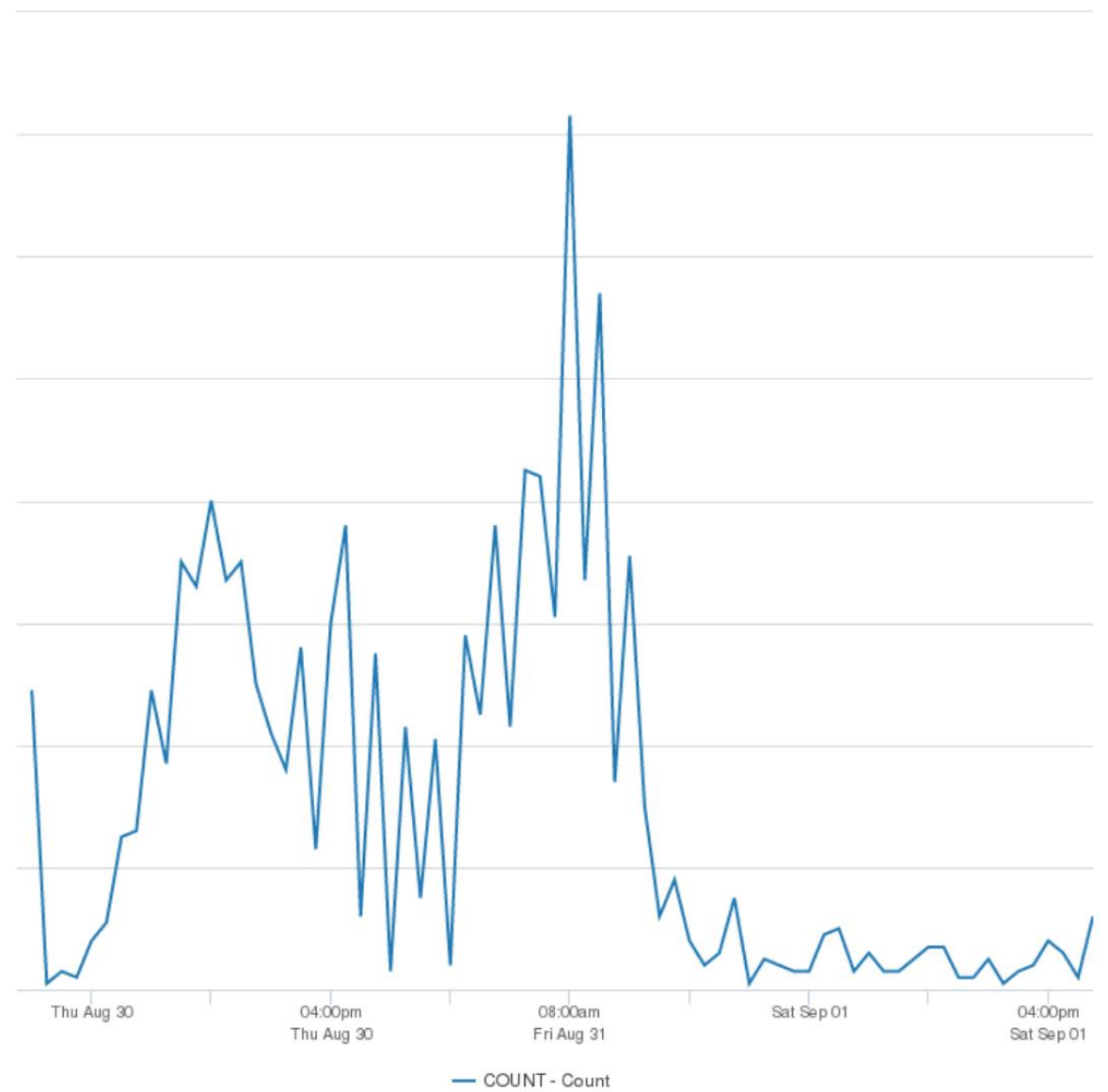
# Results

# Memory Usage



An oomd kill on a build host (kills memory hogging job)





panic\_on\_oom rate before and after oomd rollout

Questions