

September 14, 2022

# cgroup rstat's advanced adoption

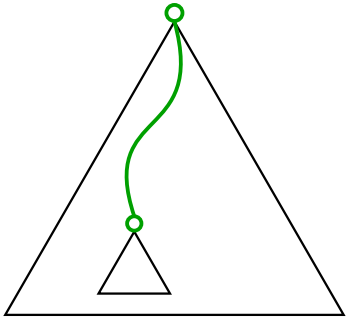
Michal Koutný, [mkoutny@suse.com](mailto:mkoutny@suse.com)



# cgroup recursive stats framework

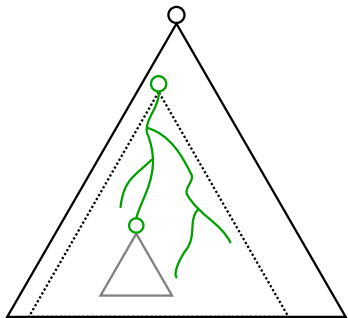
- ▶ in this talk
  - ▶ give overview to the unfamiliar (to join discussion)
  - ▶ background for analysis
  - ▶ results of my experiments
  - ▶ ideas, next steps
- ▶ quick intro
  - ▶ `$cgroup_dir/*.stat`

# Update scope



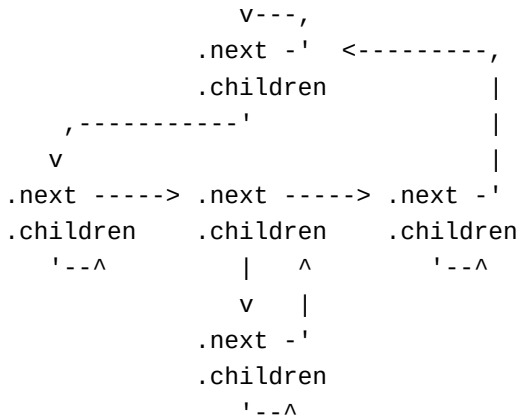
►  $O(\text{depth})$

# Flush scope



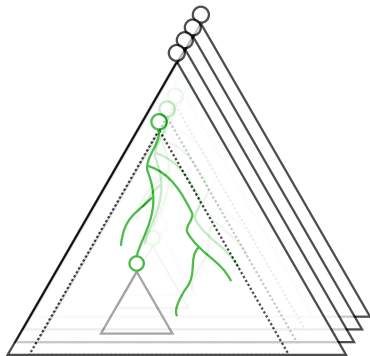
►  $O(\text{subtree\_size} \cdot \text{nr\_controllers})$

## Update tree structure



- ▶ struct cgroup\_rstat\_cpu
- ▶ topological sort overall
- ▶ .next list is LIFO processed

# Full flush scope



▶ per-cpu replicas of the update tree

# Trade-offs and benefits

- ▶ updates are per-cpu
- ▶ flushing subtree of interest
- ▶ flushing subtree with updates

# Common updaters

- ▶ core: CPU time accounting (total and components)
- ▶ io controller: BIO dispatched
- ▶ memory controller: state, events, LRU stats
- ▶ BPF kfunc<sup>1</sup>

---

<sup>1</sup>bpf-next: a319185be9f5ad13c2a296d448ac52ffe45d194c



# Common flushers

- ▶ imprecise
  - ▶ `cpu.stat`
  - ▶ `io.stat`
  - ▶ `memory.stat`<sup>2</sup>, `memory.zswap.current`
  - ▶ memcg: periodic flush (0.5 / s)
  - ▶ memcg: `workingset_refault`
- ▶ precise
  - ▶ memcg: `shrink_node` (reclaim)
  - ▶ memcg: `mem_cgroup_wb_stats` (writeback throttling)
  - ▶ BPF kfunc<sup>3</sup>
  - ▶ also cgroup removal

---

<sup>2</sup>Also memcg on v1

<sup>3</sup>bpf-next: a319185be9f5ad13c2a296d448ac52ffe45d194c

# Scalability

- ▶ locks
  - ▶ **cgroup\_rstat\_cpu\_lock**, per-cpu spinlock
    - ▶ update traversal
    - ▶ per-cpu tree flushing
  - ▶ **cgroup\_rstat\_lock**
    - ▶ all flushing (possible release with **cond\_resched**)
  - ▶ **stats\_flush\_lock** (memcg flush)
- ▶ individual controller flush callbacks
  - ▶ and time accounting in core
  - ▶ and BPF callbacks
- ▶ global accesses
  - ▶ **cgroup->self.parent** (update, flush)
  - ▶ **cgroup->rstat\_css\_list** (flush)
  - ▶ **cgroup\_subsys\_state->rstat\_css\_node** (flush)

# Memcg specialization

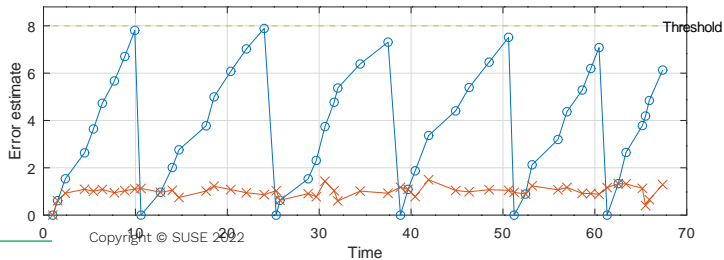
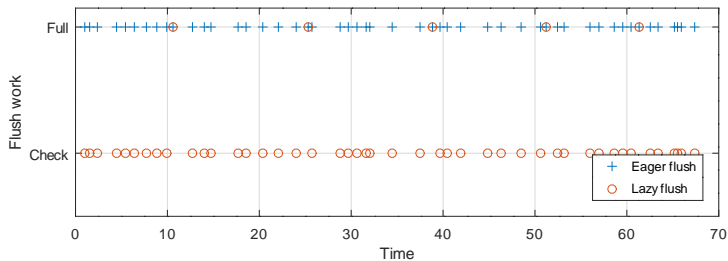
- ▶ expensive flush
  - ▶ `NR_VM_NODE_STAT_ITEMS` = 41
  - ▶ `NR_VM_EVENTS_ITEMS` = 107
  - ▶ `MEMCG_NR_STAT` = `NR_VM_NODE_STAT_ITEMS` + 7
- ▶ lazy flushing based on accumulated error estimate<sup>4</sup>
- ▶ error in  $O(\text{nr\_cpus} \cdot \text{MEMCG\_CHARGE\_BATCH})$  or  $O(\text{flush\_period} \cdot \text{cr}_{\text{max}})$

---

<sup>4</sup>error before rstat  $O(\text{nr\_cpus} \cdot \text{MEMCG\_CHARGE\_BATCH} \cdot \text{subtree\_size})$  but no flushing work was necessary



# Memcg specialization – theory



## Memcg specialization – theory

- ▶ average flush time in  $O(\text{nr\_items}^2)$ <sup>5</sup> (strawman?)
  - ▶ full flush in  $O(\text{nr\_items})$
  - ▶ frequency of reaching error estimate threshold in  $O(\text{nr\_items})$

---

<sup>5</sup>20220311160051.GA24796@blackbody.suse.cz

## Measurement setup

- ▶ 6.0.0-rc1 with a “masking” patch, 48 CPUs

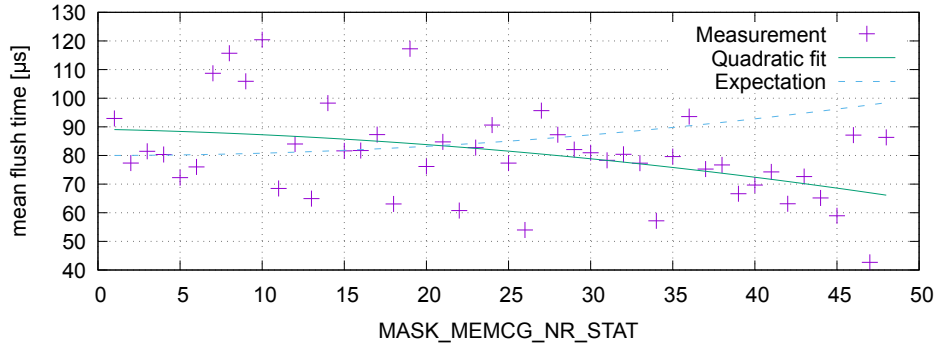
```
@@ -681,7 +685,8 @@ void __mod_memcg_state(  
-     memcg_rstat_updated(memcg, val);  
+     if (idx < MASK_MEMCG_NR_STAT)  
+         memcg_rstat_updated(memcg, val);
```

```
@@ -5404,7 +5411,7 @@ static void mem_cgroup_css_rstat_flush(  
-     for (i = 0; i < MEMCG_NR_STAT; i++) {  
+     for (i = 0; i < MASK_MEMCG_NR_STAT; i++) {
```

- ▶ workload

```
systemd-run --scope -u measure.scope \  
-p MemoryMax=1600M \  
make -C $LINUX_TREE -j $NR_CPUS bzImage
```

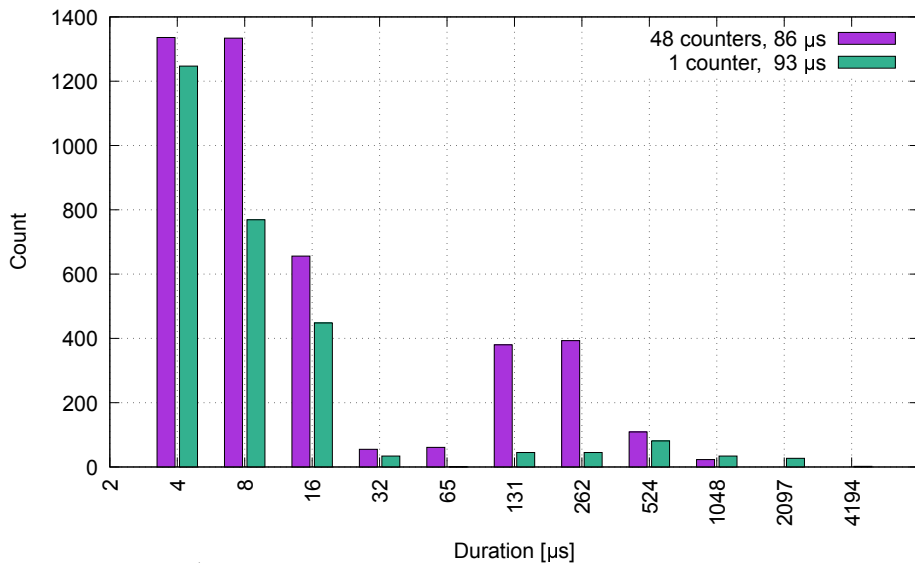
# Measurement results



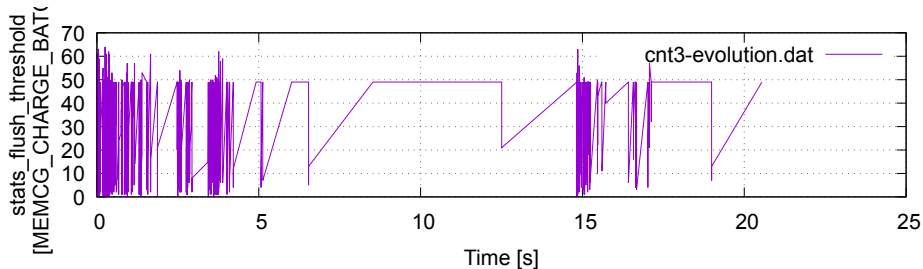
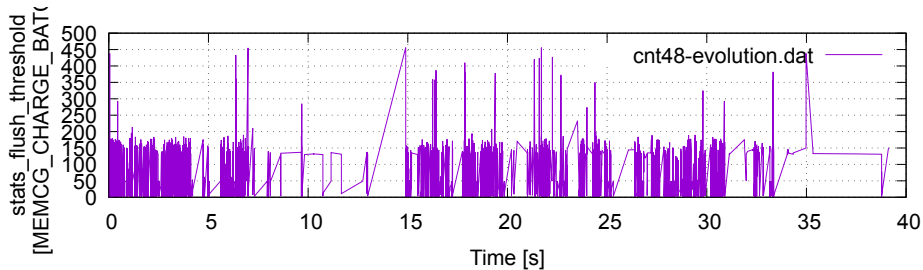
► basically noise



## Measurement results 2



# Measurement results 3





## Related work

- ▶ memcg: unify memcg stat flushing
  - ▶ piggy-back on running flushes in writeback calculation
  - ▶ bot: “we noticed a 47.8% improvement of fio.write\_iops”
- ▶ [PATCH 3/3] memcg: reduce size of memcg vmstats structures
  - ▶ flushing less entries
- ▶ [PATCH v6 3/3] blk-cgroup: Optimize blkcg\_rstat\_flush()
  - ▶ flushing less devices

## Related work

- ▶ memcg: unify memcg stat flushing
  - ▶ piggy-back on running flushes in writeback calculation
  - ▶ bot: “we noticed a 47.8% improvement of fio.write\_iops”
- ▶ [PATCH 3/3] memcg: reduce size of memcg vmstats structures
  - ▶ flushing less entries
- ▶ [PATCH v6 3/3] blk-cgroup: Optimize blkcg\_rstat\_flush()
  - ▶ flushing less devices
- ▶ cgroup: bpf: enable bpf programs to integrate with rstat
  - ▶ blank cheque in flush path
- ▶ [PATCH v7 1/4] mm: add NR\_SECONDARY\_PAGETABLE to count secondary page table uses.
  - ▶ yet another stat item

# Areas for improvement

- ▶ per controller update subtrees
  - ▶ smaller flushes (but less effective)
  - ▶ bigger memory footprint
- ▶ subtree flushing (memcg)
  - ▶ subtree locking?
  - ▶ harder error estimate tracking
- ▶ hardcode `rstat_css_list` iteration
- ▶ selective flush (required fields only)
- ▶ selective flush 2 (generalize memcg error-based flushing)
- ▶ detach rstatc from update tree *after* it's processed
  - ▶ frees time to writers

# The end

- ▶ Summary
  - ▶ The worry about  $O(nr\_items^2)$  was not fulfilled.
  - ▶ Flushing still susceptible to take long (needs measurement).
- ▶ Thank you for attention!
- ▶ Q & A