# Outline

- Problem Anatomy

- Measurements

- Designs an Implementations

- Iterations and Results

- Beyond the Results

- Innovations

- Appendix

- Basic Constraints

  - Memory limited

    - State limited (2GB for C++)

    - Input must be stored into disk

      - EXT4 needed(Raw disk is not allowed)

      - DIO wanted(Buffered IO is globally blocking and not controllable)

  - No compression

    - Good randoms hard to be compressed

    - Computing power overflow

- Basic Constraints

  - Arbitrary Crash resisted

    - WAL somewhere?

      - Sync

      - Kernel managed buffer (file mmaped memeory)

- Basic Relaxations
  - Fixed length keys/values
  - Write and Read separated

- *Goal: Striving to run at full IO bandwidth*

  - Required to implement: *Write/Read/Range*

  - (point) *Read*: special case of *Range*

- *Idea*

  - Design to make *Range* fastest

  - And then that design can be fastest in *Write/Read* as well

- *Idea(cont.)*
  - for *Range*
    - full range iteration - global order
    - How order?
    - Order of  what?
    - When to order?

- *Idea(cont.)*

  - How order?

    - Global in one store(file)

      - expensive in (ext4) file system

    - inter-store order - range partition

      - + intra-store order -> global order

- *Idea(cont.)*
    - Order of what?
        - Key (as least)
            - value is seq appended
        - Key + value
            - sub-optimal for modern (nvme) storage

- *Idea(cont.)*
  - When to order?
    - at building(Write) time
      - Take use of the redundant computing power @ write
    - at querying(Read/Range) time

- *Start to code?*

- Measure, Measure, Measure

  - Not just blindly benchmark

  - Measurements make optimizations scientific (and help you further)

*Basic Spec*

| Item | official spec[a] | public media |
|------|------------------|--------------|
| 4k random write | 500k | 562,619[b-d] |
| 4k random read | 550k | 580K[b-d] |
| sequential write | 2000 MiB/s | 2.17 GB/s[b-d] |
| sequential read | 2400 MiB/s | 2.53 GB/s(64K)[b-d] |

a. https://ark.intel.com/products/97162/Intel-Optane-SSD-DC-P4800X-Series-375GB-1-2-Height-PCIe-x4-3D-XPoint-
b. https://www.tomshardware.com/reviews/intel-optane-3d-xpoint-p4800x,5030-5.html
c. https://www.storagereview.com/intel_optane_ssd_dc_p4800x_review
d. https://www.anandtech.com/show/11930/intel-optane-ssd-dc-p4800x-750gb-handson-review/3

POLARDB
数据库性能大赛总决赛

- *Basic Infos*

  - 4k random write/read throughput seem awesome if enough parallels

  - 4k random write on par with 4k sequential write

  - Sequential read faster than random read

  - Intel Device 2701(P4800x) controller is poor

    - 31 hardware IO queues

    - max 128k message

  - Intel Device 2700(900P) is cheaper version of  2701(P4800x)

    - Measurements from 900P will be shown later

  - Several configs interesting but can not be tweaked

    - IRQ affinity/polling/IO scheduler...

- *Hints*
  - *Write*
    - 4k-128k
    - 4k is unlikely(in that the contest)
  - *Read*
    - 64threads 4k random read may be enough
  - *Range*
    - Sequential read favored

*Is 4k enough if enough parallels(queue depth) (for Optane) in my work station with centos 7.2?*

| Item | official | public | Arch/4.19.x/EXT4 | Cent OS 7.2/3.10.0-327/EXT4 |
|---|---|---|---|---|
| 4k random write | 500k | 562,619 | 552k | 553k |
| 4k random read | 550k | 580k | 587k | 587k |
| sequential write | 2000 MiB/s | 2.17 GB/s | 2204 MiB/s | 2207 MiB/s(16k bs, 564,992 4k eq) |
| sequential read | 2400 MiB/s | 2.53 GB/s | 2607 MiB/s | 2613 MiB/s(1M bs) |

1. Hardware: 2*Xeon Gold 5120, Intel Device 2700(900P)/280G
2. peak records of fio 3.1x runs under all conditions, detailed conditions seen later and bench working flow seen in the crack guide

POLARDB
数据库性能大赛总决赛

阿里云   (intel)

## Measurements

- *Sum up*

  - Online benchmark environment(Cent OS 7.2) may have slightly better performance than my workstation.

    - Modern hardware+kernel greatly improve the perf of MMU activity (and 4k random write with libaio)

  - File system(EXT4) imposes more or less performance tax on disk device.

  - IO side decision:

    - *Write:* 4*4k bs, seq

    - *Read:* 4k bs, rand

    - *Range:* 1M bs, seq, 4Threads

- *Sum-up*
  - Lower Bound of Rank Scores from IO view
    - All peaks counts from measurements:

      64000000/564992+62000000/587000+2*64000000*4/(2613*1024)+0.7+0.9=411.85s

    - Relaxed to common top observations in online runs:

      64000000/562000+62000000/587000+2*64000000*4/(2601*1024)+0.7+0.9=413.41s

      - Write phase: only get one 562k(equ, online), common iops value is 558k

      - 64000000/558000-64000000/562000=0.816s

      - Range, in fact, has other tricks for improvement but which is considered anti-rule and meanless in engineering.
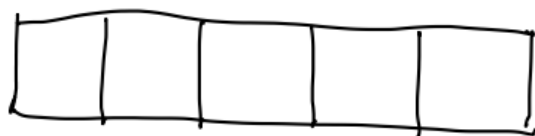
System design sketch

*Keep It Simple Stupid*

Part structure sketch

Put — Insertion of Sparse Sorted Array (SSA)

example: 0 ~ 9

gen 4 rands: ① Most Case  1  2  4  7

$O(1)$   | 1 | 2 | 4 | 7 | void |

3 ins: shift + load + store (cache hit)

0 1  2 3  4 5  6 7  8 9

"Random" Input:   n is small const

BST → BAD

SIMD → HARD TO HELP

alignment......

← $O(n)$ →

② Local clustering   1   5

4  | 1 | V | 5 | V | V |

   | 1 | V | 4 | 5 | V |

△ totally $O(1)$ and may be faster than (concurrent) hashtable

SSA insertion sketch

Logic of Write
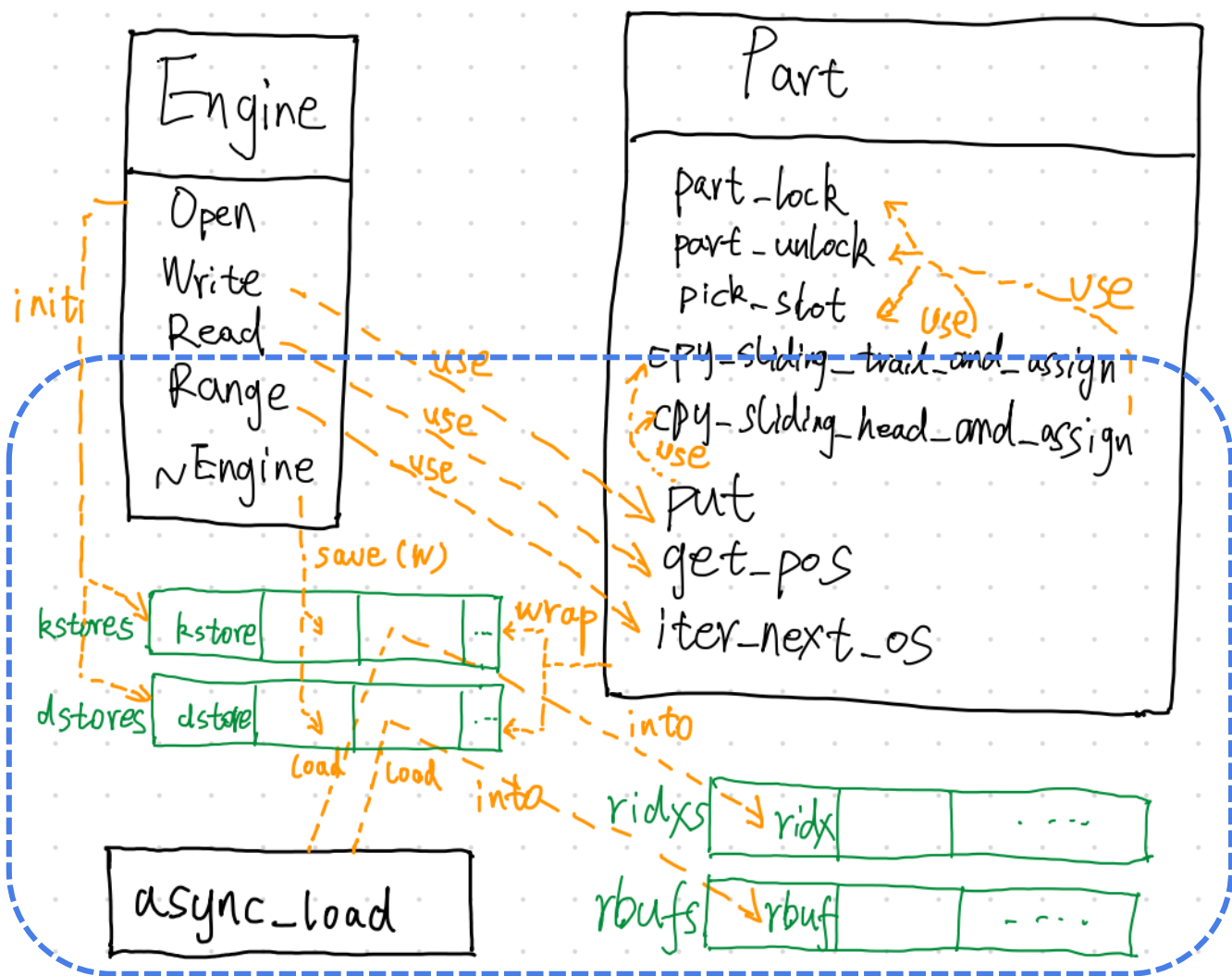
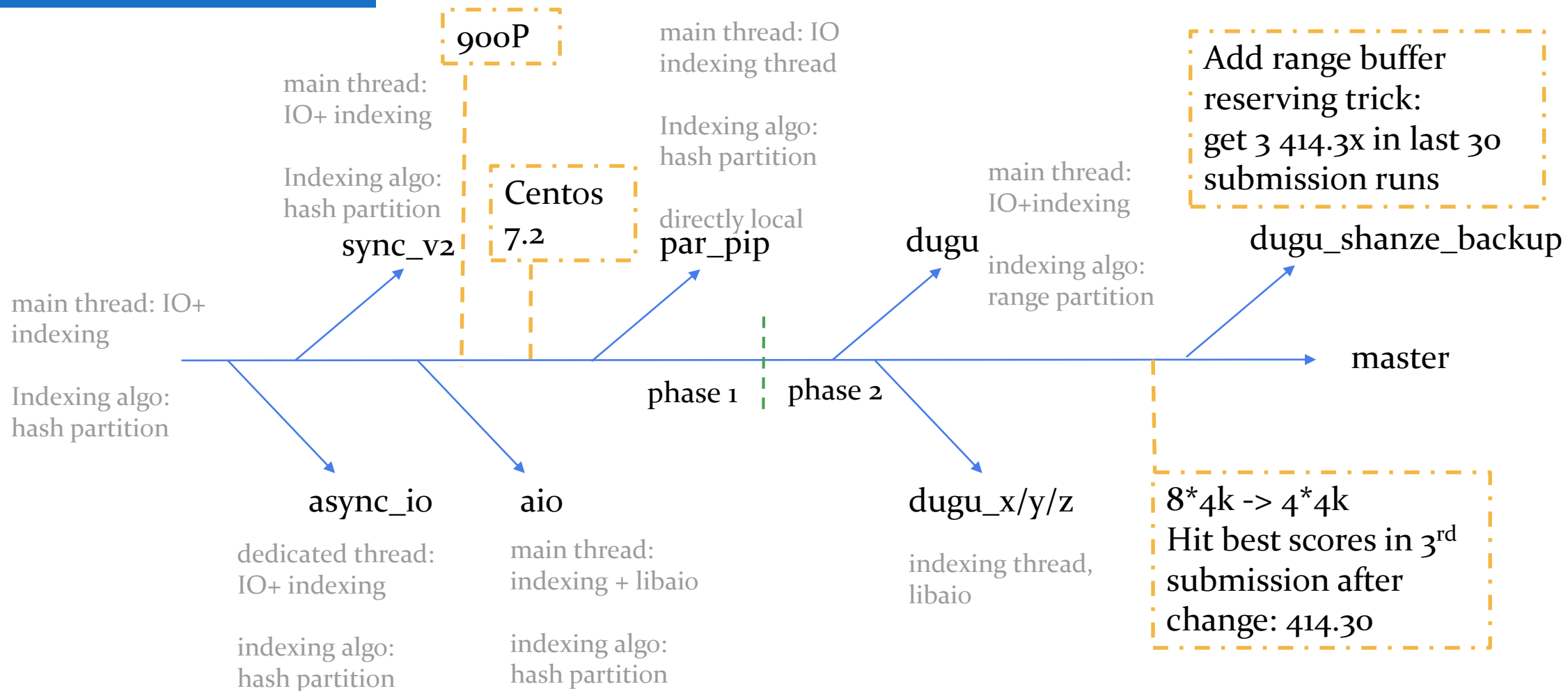Take use of pre-taken offset build a dump logic which is better than that of general lock-free queue

Logic of Read

Logic of Range

900P

main thread: IO
indexing thread

Add range buffer
reserving trick:
get 3 414.3x in last 30
submission runs

main thread:
IO+ indexing

Indexing algo:
hash partition

Indexing algo:
hash partition

Centos
7.2

main thread:
IO+indexing

directly local

sync_v2

par_pip

dugu

indexing algo:
range partition

dugu_shanze_backup

main thread: IO+
indexing

master

Indexing algo:
hash partition

phase 1    phase 2

async_io

aio

dugu_x/y/z

8*4k -> 4*4k
Hit best scores in 3rd
submission after
change: 414.30

dedicated thread:
IO+ indexing

main thread:
indexing + libaio

indexing thread,
libaio

indexing algo:
hash partition

indexing algo:
hash partition

- *Trick#1* range buffer reserving

  - 2 pass iterations

  - only use 1st part of range buffer slots in 1st pass to iterate in 2nd pass

  - only 2-4 minimum Part slots needed(4 slots locally tested, NOTE: Part obj size ~120MB )(if trick#2 used, then we only need 1 Part-sized slot)

  - definitely improve  0.7-0.9 seconds upon above general design(NOTE: mem limit of Java is 3G) if IO is assumed to be constant

  - done in dugu_shanze_backup branch

  - Outcome: 6/12 or 7/15, get 3 414.3x in last 30 submission runs

  - Ugly for engineering

- *Trick#2*  single iterator

- iterates the values by single thread, then dispatches the values into a ring buffer, all clients (second-hand) consume extracted values

- crucial difference:  the real time point to produce (and consume)

  - Idea of previous design works for arbitrary concurrent range iteration

  - But this trick can not in that it is just single iterator

    - advantage: no cache trashing

- owned online peak 42521287, others' 426xxxxx (1/425*200 ~= 0.47s)

  - this tiny diff also observed locally(by comparing fio and program runs)

- Not implemented (suspected violation to rule, nonsense for engineering)

- *Too many things beyond the results*

  - What's the effect of file system to fast storage?

  - How many modern features have we embraced?

  - What's better designs if more complex loads met?

  - How does the fast data(base) system evolve to its destination?

- File system

  - *"I can not reproduce the 4k random IOPS from public media."*

  - Investigation#1:  Random IO IOPS is very poor for my Optane SSD when I mock the proposed benchmark flow via fio run.

    - Clean the test directory (and drop the cache) -> run -> …

POLARDB
数据库性能大赛总决赛

random write
slow down

# Beyond the Results

normal
seq write

- *Why my online benchmark immune from this?*

  - *Random writes show huge overhead for un-pre-filled files writes. Forget it in the contest*

  - *fallocate cannot help for this*

  - *Just the sequential write*

- *Investigation#2:*
  - Program has a high probability (~2/3) to show an small unstable run

- Adaptive sleep
  - linear backoff algorithm
  - idiomatic rate regulation policy for dismatch-varied P/C systems

- *Conclusion #1*
  - *FS/EXT4 has obvious overheads for fast disks although some FS like XFS may pay more attention on these*
  - *Modern kernel still has many places to improve*
    - *spin lock*
    - *lack resilient configurations*
  - *So, finally, if you need to squeeze out all the performance from fast storages, you may want your own FS or directly work on raw devices*

- Modern features

  - *"Have we fully benefited from the dramatically evolving modern features?"*

    - Optane like are game changers (but maybe just for Optane DIMM)

  - Software

    - Centos 7.2 is ancient

  - Hardware

    - There are pearls out of toothpastes in recent-generation Intel processors

- Investigation#3（Software）new kernel IO-side parameters

  - RWF_HIPRI(kernel 4.6+)

  - New blk-mq(kernel 3.13+) schedulers

    - *Kyber(Facebook, kernel 4.12+)*

# Beyond the Results

## Some Interesting Kernel IO-side Parameters Investigation

| Category | bfq | mq-deadline | Kyber | none |
|---|---|---|---|---|
| 16k(prefilled/t31/4k equ/polling) | 0 | 0 | 562,432 | 562,432 |
| 16k(prefilled/t31/4k equ) | 0 | 0 | 562,688 | 562,432 |
| 4k(prefilled/t56/polling) | 0 | 0 | 559,104 | 558,592 |
| 4k(prefilled/t56) | 121,000 | 280,000 | 558,592 | 557,824 |
| 4k(no prefilled) | 120,000 | 274,000 | 555,000 | 552,000 |

Legend: bfq, mq-deadline, Kyber, none

*Findings:*
1. Polling mode slightly positive in 4k but not for >4k
2. Kyber scheduler helps a little in all spectrum

Benchmark method: running (fio) program interleavingly untill the nexting two runs can not hit higher results, pick up the highest results: none -> kyber -> none -> kyber -> ...

- Investigation#4 (Hardware) Intel TSX-NI
  - Parallel (speculative) execution on multi cores without sync (when no data race)
    - little overhead in the optimistic case
    - This competition just provides a natural optimistic case
    - Broadwell+(GA)
    - Lock-free is hard for complex data structures, and even not necessarily faster than lock
  - NOTE: competition's CPU power overflow, so:
  - *Conclusion#4*
    - *TSX should be first-try alternative for 'old' lock(std::mutex) with basic guards*

- **Better designs**
  - *"Are the patterns and designs found in the competition valid for the real world engineering?"*
  - Mmap for process crash protection is homework-level simplification
    - Smaller block size(4k) is still best choice for OLTP
  - IO stuffs be tackled at the IO layer
    - sequential loads merged at the IO layer: Kyber like IO scheduler
    - Make your own FS (abcFS ... PolarFS)
  - NUMA-ware rather than NUMA-hidden(UMA) – however PolarDB is UMA style...
  - Sync mode is definitely the performance killer, async io is favored in wise practices
    - libaio is truly a good to do async io when play with modern kernel

- System evolve

  - *"Your Part is naïve. So, in practice, we still use B-tree like data structures.  Really?"*

  - If distribution of key known well, flat array like indexing (base)structure is enough

  - The Part algorithm works well with recent hot learned index structure* as base structures with extended actions

    - query*, insertion, update(remove) and even concurrent

    - hard part shifted to modelling (via ANN* ...)

  * *The Case for Learned Index Structures,*

  *2017-2018, Tim(MIT), Alex, Ed, Jeff Dean(Google) et al.*

The Case for Learned Index Structures

| Tim Kraska* | Alex Beutel | Ed H. Chi |
|---|---|---|
| MIT | Google, Inc. | Google, Inc. |
| Cambridge, MA | Mountain View, CA | Mountain View, CA |
| kraska@mit.edu | alexbeutel@google.com | edchi@google.com |

| | Jeffrey Dean | Neoklis Polyzotis | |
|---|---|---|---|
| | Google, Inc. | Google, Inc. | |
| | Mountain View, CA | Mountain View, CA | |
| | jeff@google.com | npolyzotis@google.com | |

POLARDB
数据库性能大赛总决赛

阿里云    (intel)

- *Proposed patterns* *to full IO bandwidth*

    - Observations: 562k+iops(w)，605k+iops(r)，42.5M+iops(range)

- *Concise but efficient* *designs and implementations*

    - Simple sync write taking full use of prebuilt offset (faster than general lock-free queue)

    - O(1) range indexing structure

    - customized built-in replacements are confirmed to be effective for more stable peak (pdbr_cpy/mov/bswap)

    - Tens of designs and options(IO pattern/Comp pattern/core affinity/numa...) have been implemented and compared

- *More* *exploits to modern hardware*

    - TSX used wisely - Thread safety is guaranteed without performance loss

- *Proposed high performance designs* *works better under more balanced engineering scenarios*
  - ugly tricks that are not engineering discarded
- *Correctness First*
  - Lines of test sources are more than that of implementation sources.
  - Concurrent correctness done even in low contention
- *Measurements* *based explorations* *makes optimization scientific*
  - Adaptive sleep
  - TSX optimization
  - Better understandings to kernel and modern engineering

# Wooooooooooo

Salute to Aliyun's Ten years

Thanks for Aliyun, Tianchi, 凝岚

# 金明剑(Jin Mingjian)

Director, Data Department at Tigerjoys

Doctor, University of Chinese Academy of Sciences

- Scala@Google Summer of Code, 2010
- Java8@Landz(battle tested for Netty replacement...), 2013

Interests: Data Engineering with some Data science

# Appendix#1



before slow down

| Function / Call Stack | CPU Time ▼ |
|---|---|
| ▼ __read_once_size | 0.501s |
| ▼ ↖ native_queued_spin_lock_slowpa | 0.501s |
| ▶ ↖ pv_queued_spin_lock_slowpath | 0.493s |
| ▶ ↖ pv_queued_spin_lock_slowpath | 0.004s |
| ▶ ↖ pv_queued_spin_lock_slowpath | 0.004s |

after slow down

| Function / Call Stack | CPU Time ▼ |
|---|---|
| ▼ __read_once_size | 0.464s |
| ▼ ↖ native_queued_spin_lock_slowpa | 0.464s |
| ▶ ↖ ____versions ← jbd2_journal_g | 0.265s |
| ▶ ↖ jbd2_journal_dirty_metadata ← | 0.194s |
| ▶ ↖ __jbd2_log_wait_for_space ← j | 0.005s |

before
slow
down

after
slow
down

vmlinux!__read_once_size - compiler.h
vmlinux!native_queued_spin_lock_slowpath+
vmlinux!native_queued_spin_lock_slowpath+
vmlinux!pv_queued_spin_lock_slowpath+0xa
vmlinux!queued_spin_lock_slowpath - qspinl
vmlinux!queued_spin_lock+0x11 - qspinlock.l
vmlinux!do_raw_spin_lock - spinlock.h:180
vmlinux!__raw_spin_lock+0xa - spinlock_api_
vmlinux!_raw_spin_lock+0x9 - spinlock.c:144
jbd2!____versions+0x1f2f - jbd2.mod.c:27
jbd2!jbd2_journal_get_write_access+0x2c - t
ext4!__ext4_journal_get_write_access+0x2c
ext4!ext4_split_extent_at+0x201 - extents.c:3
ext4!ext4_split_extent+0xc2 - extents.c:3378
ext4!ext4_split_convert_extents+0xac - exter
ext4!ext4_ext_handle_unwritten_extents+0x9
ext4!ext4_ext_map_blocks+0x4ff - extents.c:
ext4!ext4_map_blocks+0xed - inode.c:636
ext4!_ext4_get_block+0x8e - inode.c:785
ext4!ext4_get_block_trans+0x8b - inode.c:84
ext4!ext4_dio_get_block_unwritten_sync+0x3

vmlinux!__read_once_size - compiler.h
vmlinux!native_queued_spin_lock_slowpath+0xcf - qspinlock.c:
vmlinux!native_queued_spin_lock_slowpath+0x2d - qspinlock.c
vmlinux!pv_queued_spin_lock_slowpath+0xa - paravirt.h:679
vmlinux!queued_spin_lock_slowpath - qspinlock.h:32
vmlinux!queued_spin_lock+0x11 - qspinlock.h:88
vmlinux!do_raw_spin_lock - spinlock.h:180
vmlinux!__raw_spin_lock+0xa - spinlock_api_smp.h:143
vmlinux!_raw_spin_lock+0x9 - spinlock.c:144
jbd2!jbd2_journal_dirty_metadata+0x14d - transaction.c:1465
ext4!trace_event_define_fields_ext4_find_delalloc_range+0x65
ext4!ext4_split_extent_at+0x117 - extents.c:3234
ext4!ext4_split_extent+0x10a - extents.c:3339
ext4!ext4_split_convert_extents+0xac - extents.c:3703
ext4!ext4_ext_handle_unwritten_extents+0x9be - extents.c:403
ext4!ext4_ext_map_blocks+0x4ff - extents.c:4345
ext4!ext4_map_blocks+0xed - inode.c:636
ext4!_ext4_get_block+0x8e - inode.c:785
ext4!ext4_get_block_trans+0x8b - inode.c:845
ext4!ext4_dio_get_block_unwritten_sync+0x30 - inode.c:924
vmlinux!get_more_blocks+0xe6 - direct-io.c:711

**Top Hotspots**

This section lists the most active functions in your application.
performance.

| Function | Module | CPU Time |
|---|---|---|
| polar_race::pdbr_cpy | gtestAll | 0.085s |
| __mcount_internal | libc-2.28.so | 0.035s |
| polar_race::EngineRace::Write | gtestAll | 0.030s |
| queued_read_lock | vmlinux | 0.020s |
| ext4_mark_iloc_dirty | ext4 | 0.015s |
| [Others] | | 0.301s ⚑ |

**Top Hotspots**

This section lists the most active functions in your application.
performance.

| Function | Module | CPU Time |
|---|---|---|
| __read_once_size | vmlinux | 0.200s |
| __mcount_internal | libc-2.28.so | 0.040s |
| polar_race::pdbr_cpy | gtestAll | 0.025s |
| ext4_mark_iloc_dirty | ext4 | 0.020s |
| [crc32c_intel] | crc32c_intel | 0.020s |
| [Others] | | 0.271s |

```
vmlinux!__read_once_size - compiler.h
vmlinux!native_queued_spin_lock_slowpath+0xcf
vmlinux!native_queued_spin_lock_slowpath+0x2(
vmlinux!pv_queued_spin_lock_slowpath+0xa - pa
vmlinux!queued_spin_lock_slowpath - qspinlock.h
vmlinux!queued_spin_lock+0x45 - qspinlock.h:88
vmlinux!queued_read_lock_slowpath+0x1f - qrwl(
jbd2!trace_event_define_fields_jbd2_handle_stat
jbd2!jbd2__journal_start+0xd9 - transaction.c:439
ext4!ext4_dirty_inode+0x2d - inode.c:311
vmlinux!arch_static_branch - jump_label.h:36
vmlinux!static_key_false - jump_label.h:142
vmlinux!trace_writeback_dirty_inode - writeback.h
vmlinux!__mark_inode_dirty+0x41 - fs-writeback.
vmlinux!generic_update_time+0xb6 - inode.c:165
vmlinux!file_update_time+0xe1 - inode.c:1884
vmlinux!__generic_file_write_iter+0x98 - filemap.
ext4!inode_unlock - fs.h:743
ext4!ext4_file_write_iter+0xc6 - file.c:267
vmlinux!new_sync_write+0xfb - read_write.c:476
vmlinux!vfs_write+0x36 - read_write.c:574
```

```
linux/fs.h: * hopefully graduate it to a proper O_CMTIME flag supported by open(2) soon.
linux/fs.h:#define FMODE_NOCMTIME               ((__force fmode_t)0x800)
linux/fs.h:#define S_NOCMTIME    128      /* Do not update file c/mtime */
linux/fs.h:#define IS_NOCMTIME(inode)    ((inode)->i_flags & S_NOCMTIME)
```

# Beyond the Results

## Part TSX Investigations



Horizontal bar chart categories (top to bottom): "fastest_hashmap"**, std::unordered_map**, no_lock, std::mutex, tsx_cust_mov, tsx_no_fs, tsx_naive

X-axis: 0, 0.2, 0.4, 0.6, 0.8, 1, 1.2, 1.4, 1.6

Legend: ■ microbenchmark (seconds, the shorter is better)

** **fastest_hashmap** and **undereded_map** are single thread 1M k/v insertions(and **no_lock** are not thread safe structure), just added for infos

fastest_hashmap got from https://probablydance.com/2017/02/26/i-wrote-the-fastest-hashtable/

\* microbenchmark infos: 52 threads, 1M K/V insertions per thread, 32k slots per Part, 2048 Parts; pick best in five runs

# Part TSX Optimizaition Steps

Clockticks: 116,861,548,459
Transactional Cycles: 51,544,275,876
Abort Cycles: 14,790,221,004
Abort Cycles (%): 28.694
Total Thread Count: 53
Paused Time: 0s

Clockticks: 93,859,249,077
Transactional Cycles: 46,667,331,412
Abort Cycles: 7,005,122,870
Abort Cycles (%): 15.011
Total Thread Count: 53
Paused Time: 0s

Clockticks: 82,132,208,750
Transactional Cycles: 42,115,776,490
Abort Cycles: 367,151,189
Abort Cycles (%): 0.872
Total Thread Count: 53
Paused Time: 0s

tsx_naive → tsx_no_fs → tsx_cust_mov

TSX Aborts: 13,920,000
Instruction: 4,140,000
Data Conflict: 9,750,000
Capacity: 10,000
Other: 20,000
Total Thread Count: 53
Paused Time: 0s

TSX Aborts: 5,090,000
Instruction: 4,290,000
Data Conflict: 790,000
Capacity: 0
Other: 10,000
Total Thread Count: 53
Paused Time: 0s

TSX Aborts: 140,000
Instruction: 0
Data Conflict: 140,000
Capacity: 0
Other: 0
Total Thread Count: 53
Paused Time: 0s