施闻轩 · TiKV Engineer
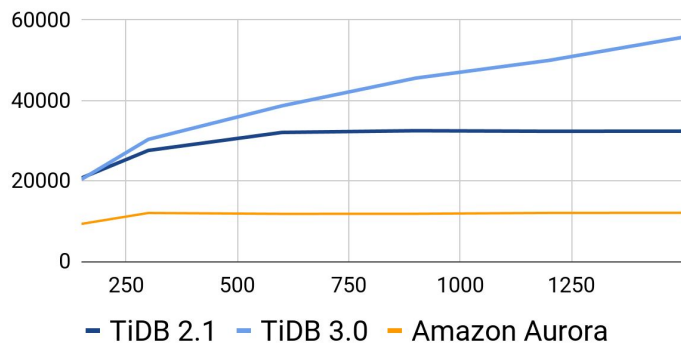
TiKV Optimization & Development

TiKV Coprocessor SIG Tech Lead

# Sysbench

2.1 → 3.0

3.0 → 4.0 (not released)

### Sysbench - Update Non-Index



Legend: TiDB 2.1, TiDB 3.0, Amazon Aurora

Right chart legend: 3.0.1 QPS, 4.0 QPS

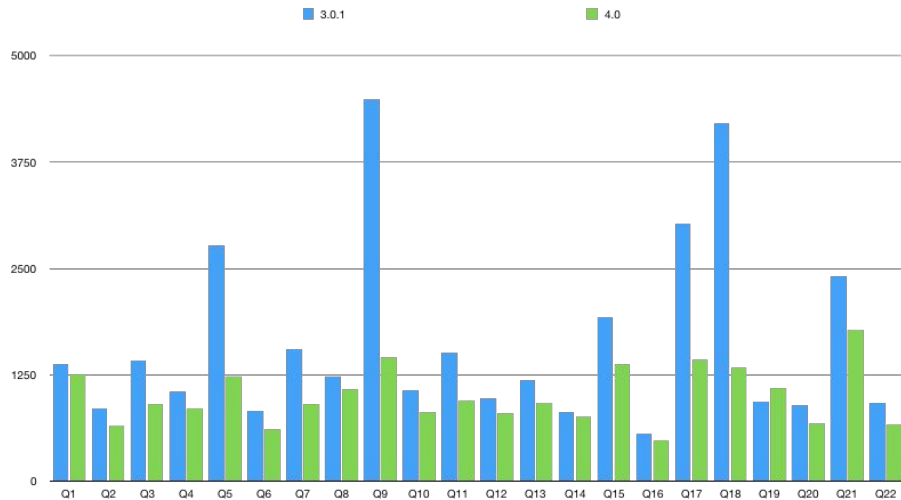| Benchmark | 3.0.1 QPS | 4.0 QPS |
|---|---|---|
| oltp_insert | 24838 | 34330.48 |
| oltp_update_non_index | 33849 | 48808.18 |
| oltp_update_index | 21598 | 29676.33 |
| oltp_read_write | 47166 | 59668.67 |
| oltp_point_select | 180115 | 241559.69 |

PingCAP

PingCAP.com

# TPC-H

## 2.1 → 3.0

## 3.0 → 4.0 (not released)

# Performance boost comes from…

1. New Feature / Algorithm / Data Structure

   ○ v3.0: Multi-thread Raft

   ○ v3.0: Batch gRPC messages

   ○ v3.0: Titan Engine

   ○ v3.0: Vectorization

   ○ v4.0: Will be officially announced in future

2. **Engineering Improvement / Constant Optimization**

# Rust Claims...

**Rust**

**GET STARTED**

Version 1.39.0

A language empowering everyone
to build reliable and efficient software.

## Why Rust?

### Performance

Rust is blazingly fast and memory-efficient: with no runtime or garbage collector, it can power performance-critical services, run on embedded devices, and easily integrate with other languages.

### Reliability

Rust's rich type system and ownership model guarantee memory-safety and thread-safety — and enable you to eliminate many classes of bugs at compile-time.

### Productivity

Rust has great documentation, a friendly compiler with useful error messages, and top-notch tooling — an integrated package manager and build tool, smart multi-editor support with auto-completion and type inspections, an auto-formatter, and more.

PingCAP

# Rust Claims...

**Rust**

doesn't mean that you can

write efficient code unthinkingly

Rust is blazingly fast and memory-efficient: with no runtime or garbage collector, it can power performance-critical services, run on embedded devices, and easily integrate with other languages.

Rust's rich type system and ownership model guarantee memory-safety and thread-safety — and enable you to eliminate many classes of bugs at compile-time.

Rust has great documentation, a friendly compiler with useful error messages, and top-notch tooling — an integrated package manager and build tool, smart multi-editor support with auto-completion and type inspections, an auto-formatter, and more.

PingCAP

PingCAP.com

# Find Hotspots
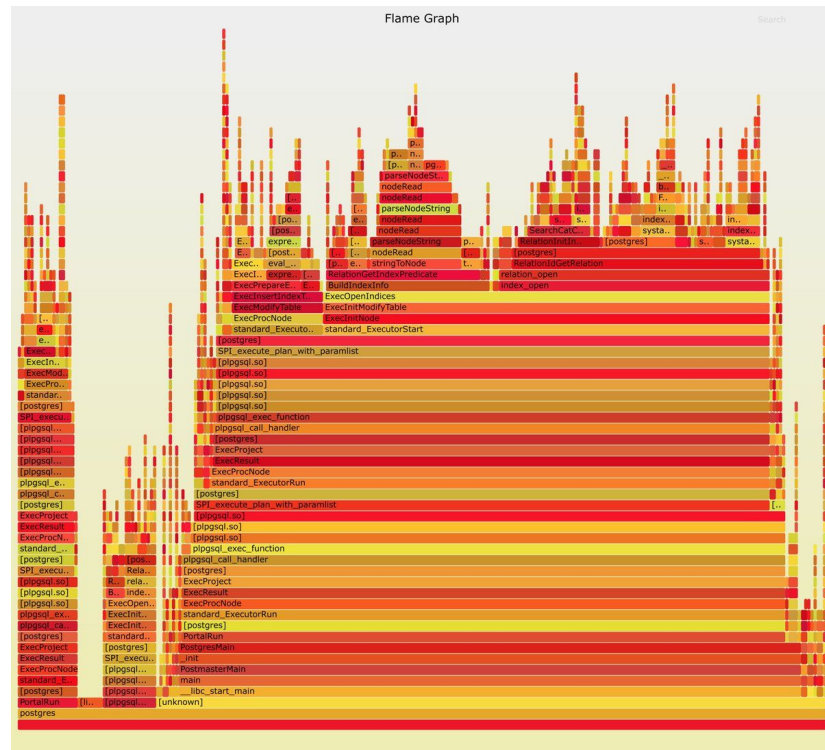
# Why?

Amdahl's law

$$S_{\text{latency}}(s) = \frac{1}{(1-p) + \frac{p}{s}}$$

PingCAP

# perf + FlameGraph

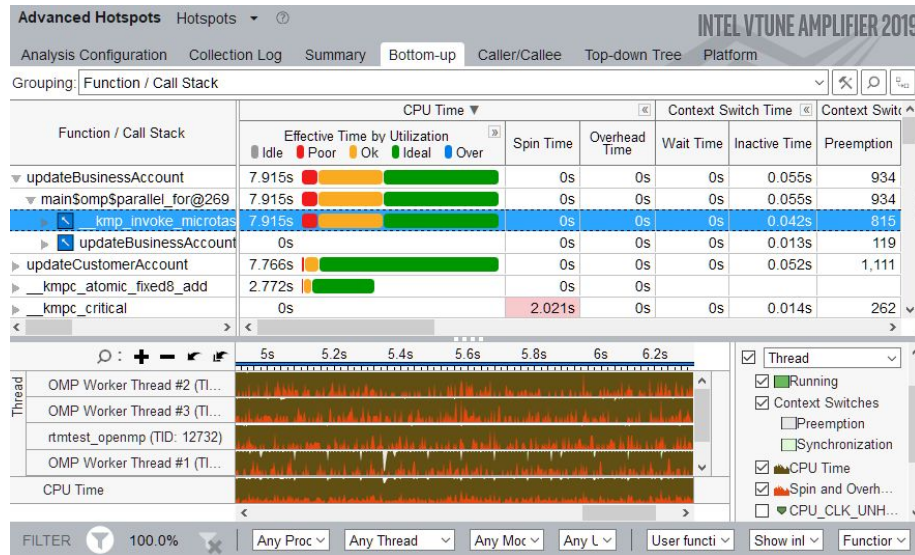https://github.com/brendangregg/FlameGraph

```
$ perf record …
$ perf script |
  ./stackcollapse-perf.pl |
  ./flamegraph.pl
```

Even very useful in production!



Flame Graph

# Intel VTune

- It is free now

- Easy to use GUI

- Top-down Microarchitecture Analysis

# Low-hanging Fruits

# Replace Memory Allocators

Use efficient memory allocators like [jemalloc](jemalloc):

```
#[global_allocator]
static GLOBAL: Jemalloc = Jemalloc;
```

Improve ~25% for TPC-H Q1.

# Link Time Optimization

LTO can apply cross-crate and even cross-language optimizations!

```
[profile.release]
lto = true
```

Improve ~15% for TPC-H Q1.

# Set Target CPU

If not set (default):

- Portable release (e.g. x86_64)
- Modern instructions like AVX is not used

# (Almost) Drop-in Replacement

- Hash Collections like `std::collections::HashMap` → hashbrown *(already in std)*
  - 2.6x faster insert

- Concurrent utilities like `std::sync::mpsc` → crossbeam
  - It is even faster than Go channel

- Synchronization primitives like `std::sync::Mutex` → parking_lot
  - 5x faster Mutex

PingCAP

# Micro Optimizations

# Compiler Can Be Stupid

```rust
1  pub fn foo(v: Vec<u64>) -> Vec<u64> {
2      v.into_iter().map(|v| v).collect()
3  }
4
```

```asm
34  example::foo:
35      push    rbp
36      push    r15
37      push    r14
38      push    r13
39      push    r12
40      push    rbx
41      sub     rsp, 72
42      mov     r12, rdi
43      mov     rbx, qword ptr [rsi]
44      mov     rax, qword ptr [rsi + 8]
45      mov     r15, qword ptr [rsi + 16]
46      lea     r14, [rbx + 8*r15]
47      mov     qword ptr [rsp], 8
48      xorps   xmm0, xmm0
49      movups  xmmword ptr [rsp + 8], xmm0
50      mov     qword ptr [rsp + 40], rbx
51      mov     qword ptr [rsp + 48], rax
52      mov     qword ptr [rsp + 56], rbx
53      mov     qword ptr [rsp + 64], r14
54      test    r15, r15
55      je      .LBB2_1
56      mov     qword ptr [rsp + 24], rax
57      mov     qword ptr [rsp + 32], r12
58      lea     r12, [8*r15]
59      sar     r12, 3
60      mov     ecx, 8
61      xor     r13d, r13d
62      mov     rax, r12
63      mul     rcx
64      mov     rbp, rax
65      setno   al
66      jo      .LBB2_5
67      mov     r13b, al
68      shl     r13, 3
69      mov     rdi, rbp
```

# Less Allocation, Less Copy

Pre-allocate collections when possible:

- `Vec::with_capacity()`

# Less Allocation, Less Copy

Use memory pool techniques to reduce allocation:

- Dynamic collections over pre-allocated memory: [bumpalo](#)
- Dynamically allocate same object: [TypedArena](#)
- 1d `Vec<Vec<T>>`: [nested](#)

# Less Allocation, Less Copy

Prefer stack allocation instead of heap allocation:

- Static array allocates on stack: `[T; N]`
- `Box<T>` allocates on heap
- `Vec<T>` allocates data on heap (obviously)
- SmallVec can be very useful
  - allocates N bytes on stack but is also growable (on heap for len > N).

# Less Allocation, Less Copy

Move less data

- `Box<[T]>` is smaller than `Vec<T>`

- Wrap *big* and *frequently moved* structures with `Box<T>` to keep it small

# Less Allocation, Less Copy

Error size matters:

- `Result<T, E>` is an enum: `sizeof(Result<T, E>) == max(sizeof(T), sizeof(E))`
- The size of `E` is the minimum size of your `Result`.
- Try `Result` is not zero cost sometimes!

TiKV TPC-H improved by 14% by reducing errors from 200 bytes to 8 bytes:

- Reduced 3×200 bytes memory copy each KV iterate
- `Result<T, Error>` → `Result<T, Box<Error>>`

# Less Allocation, Less Copy

Use reference where possible:

- Annotate the structure with a lifetime parameter to contain references.

```rust
pub struct WriteRef<'a> {
    pub write_type: WriteType,
    pub start_ts: TimeStamp,
    pub short_value: Option<&'a [u8]>,
}
```

# Less Allocation, Less Copy

Use reference where possible:

- Use `Arc<T>` instead of a raw reference when lifetime is hard to write.

# Less Allocation, Less Copy

Use reference where possible:

- Use DSTs when you only want to wrap a DST reference (like `&[u8]`).

Example: Need `&[u8]` internally, but encapsulate with a type to provide extra feature.

```rust
struct Key(Vec<u8>);
fn get(k: &Key);
```

```rust
struct KeyRef<'a>(&'a [u8]);
fn get(k: KeyRef<'_>);
```

```rust
struct KeyRef([u8]);
fn get(k: &KeyRef);
```

| ✕ Not zero cost | ✓ Zero cost, but interface doesn't accept reference any more | ✓ Best: DST wrapper |

# Less Allocation, Less Copy

Use reference where possible:

- Unbounded lifetime / raw pointer:

```rust
unsafe fn erase_lifetime<'a, T: ?Sized>(v: &T) -> &'a T {
    &*(v as *const T)
}
```

- Useful to simulate a `'Self` lifetime.

- Hate `unsafe{}` ? [owning_ref](#) or [rental](#) can help.
  - *However sometimes "unsafe" but **clear & simple** code is better.*

# Prefer Static Dispatch

```rust
fn foo(callback:
    Box<dyn FnOnce()>
);
```

```rust
fn foo<F: FnOnce()>
    (callback: F);
```

```rust
fn foo(callback:
    impl FnOnce()
);
```

| ✕ Box & Dynamic Dispatch | ✓ Static Dispatch | ✓ Syntax Sugar |
|---|---|---|

Use `Box<dyn T>` only when you want different base types in one type:

- Return different base types in one function: `fn foo() -> Box<dyn T>`

- Store different base types in one container: `Vec<Box<dyn T>>`

PingCAP

PingCAP.com

# Reduce Branches

- To reduce instructions

- To reduce misprediction

Example: Memory-comparable encoding bytes (padding = 4).

| F | o | o |
|---|---|---|

Source

| F | o | o | \1 | \FE |
|---|---|---|----|-----|

Asecnding

| F | o | o | \FE | \1 |
|---|---|---|-----|----|

Descending

```
fn memcmp_encode(bytes: &[u8], is_desc: bool);
```

✕ Branch

PingCAP

# Reduce Branches

- To reduce instructions
- To reduce misprediction

Example: Memory-comparable encoding bytes (padding = 4).

| F | o | o |
|---|---|---|

Source

| F | o | o | \1 | \FE |
|---|---|---|----|-----|

Asecnding

| F | o | o | \FE | \1 |
|---|---|---|-----|----|

Descending

```rust
fn memcmp_encode_asc(bytes: &[u8]);

fn memcmp_encode_desc(bytes: &[u8]);
```

✓ No branch, some duplicate code

# Reduce Branches

- To reduce instructions

- To reduce misprediction

Example: Memory-comparable encoding bytes (padding = 4).

| | | | | |
|---|---|---|---|---|
| F | o | o | | |

Source

| | | | | |
|---|---|---|---|---|
| F | o | o | \1 | \FE |

Asecnding

| | | | | |
|---|---|---|---|---|
| F | o | o | \FE | \1 |

Descending

```rust
trait Mode { const FLIP_BITS: bool; }

struct ModeAsc;
impl Mode for ModeAsc {
    const FLIP_BITS: bool = false;
}
struct ModeDesc;
impl Mode for ModeDesc {
    const FLIP_BITS: bool = true;
}

fn memcmp_encode(bytes: &[u8], mode: impl Mode);
```

✓ No branch (generics)

PingCAP

PingCAP.com

# Reduce Branches

- To reduce instructions

- To reduce misprediction

Example: Memory-comparable encoding bytes (padding = 4).

| F | o | o |
|---|---|---|

Source

| F | o | o | \1 | \FE |
|---|---|---|----|-----|

Asecnding

| F | o | o | \FE | \1 |
|---|---|---|-----|----|

Descending

```
fn memcmp_encode<const T: bool>
    (bytes: &[u8], is_desc: T);
```

✓ No branch (const generics) NIGHTLY

Useful Friends

# Micro Benchmark: cargo bench

```rust
#[bench]
fn bench_foo(b: &mut Bencher) {
    b.iter(|| foo());
}
```

Run benchmark:

```
$ cargo bench
```

# Micro Benchmark: cargo bench

Remember to add `black_box`:

```rust
#[bench]
fn bench_add(b: &mut Bencher) {
    b.iter(|| add(1, 2));
}
```

```rust
#[bench]
fn bench_add(b: &mut Bencher) {
    b.iter(|| add(
        black_box(1),
        black_box(2),
    ));
}
```

✕ Constants are likely to be optimized away

✓ Wrap constants with `black_box`

# Micro Benchmark: criterion

criterion is a replacement for cargo bench.

Features:

- Statistics-driven, avoid measurement noise
- Support different measurements (e.g. number of instructions instead of wall time)
- Generate charts
- Support parameterized benchmark

# Compiler Explorer

[https://godbolt.org/](https://godbolt.org/)

# Callgrind + KCacheGrind

- Heavyweight
- Not a sample based profiler
- Run in Valgrind
- **Precise** function calls
- You can use the callgrind crate to control start / stop (to skip counting unwanted part)
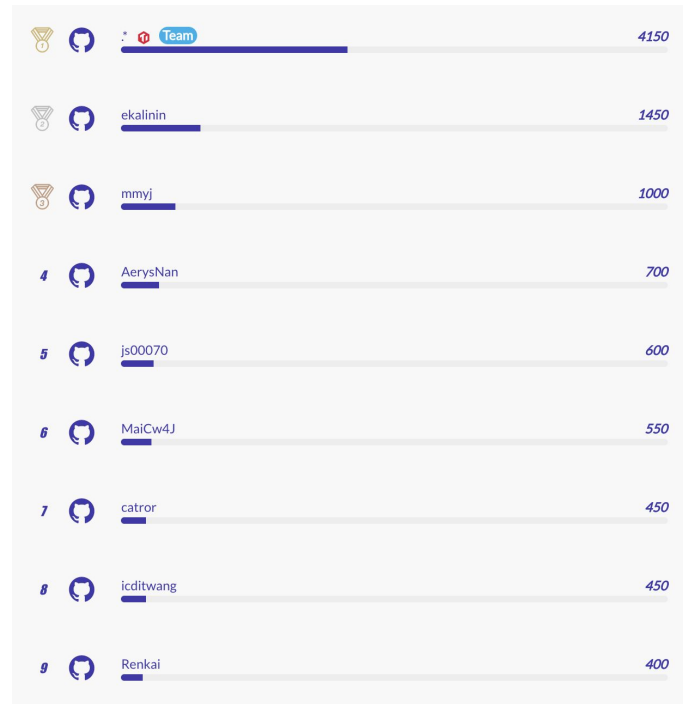


PingCAP.com

# Want to Contribute?

# TiKV Community

- [TiKV Help Wanted Issues](#)

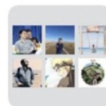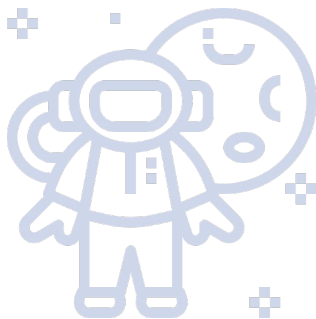- [TiKV PCP Issues](#)

- TiDB Hackathon

# PCP Season 1

- [TiDB Performance Challenge Program](#)

- Focus on *performance improvement*

- Season 1: 2019.11.04 - 2020.02.04

- You will be mentored

- You can be rewarded

| | | | |
|---|---|---|---|
| 🥇 | ⬡ | ⬣ Team | 4150 |
| 🥈 | ⬡ | ekalinin | 1450 |
| 🥉 | ⬡ | mmyj | 1000 |
| 4 | ⬡ | AerysNan | 700 |
| 5 | ⬡ | js00070 | 600 |
| 6 | ⬡ | MaiCw4J | 550 |
| 7 | ⬡ | catror | 450 |
| 8 | ⬡ | icditwang | 450 |
| 9 | ⬡ | Renkai | 400 |

# Special Interest Group

- Focus on specific module.

- Promote to SIG Reviewer / SIG Committer.

- You can receive SIG credit and PCP credit at the same time!

- [TiKV Coprocessor SIG](#)

- [TiKV Engine SIG](#)

- [TiDB Expression SIG](#)

Thank You !

Infra Meetup No. 120 活动交流群

该二维码 7 天内 (12 月 13 日前) 有效，重新进入将更新