



**T** TigerBeetle

# 1000X WORLD TOUR

13 CITIES | 6 DAYS

# The Tale of Taming TigerBeetle's Tail Latency



## Act I: The Protagonist

---

# TigerBeetle

TigerBeetle is a high-performance database for financial accounting.

# TigerBeetle

TigerBeetle is a high-performance database for financial accounting.

## Debit/Credit Interface

```
Transfer = struct {
    id: u128,
    debit_account_id: u128,    // Who
    credit_account_id: u128,   // Who
    amount: u128,              // How much
    timestamp: u64,            // When
    //...                  What, Why, Where
}
```

128 bytes (2 cache lines).

## Debit/Credit Interface

```
Transfer = struct {
    id: u128,
    debit_account_id: u128, // Who
    credit_account_id: u128, // Who
    amount: u128,           // How much
    timestamp: u64,          // When
    //...                   // What, Why, Where
}
```

128 bytes (2 cache lines).

Batching, batching, batching.

# TigerBeetle's Design Decisions

- strictly serializable
- highly durable and available<sup>1</sup>
  - network and node faults<sup>2</sup>
  - storage fault model<sup>3</sup> <sup>4</sup>
  - deterministic simulation testing and extensive fuzzing

---

<sup>1</sup><https://jepsen.io/analyses/tigerbeetle-0.16.11>

<sup>2</sup>Viewstamped Replication

<sup>3</sup>A Study of SSD Reliability in Large-Scale Enterprise Storage Deployments

<sup>4</sup>Protocol Aware Recovery

# TigerBeetle's Design Decisions

- strictly serializable
- highly durable and available<sup>1</sup>
  - network and node faults<sup>2</sup>
  - storage fault model<sup>3</sup> <sup>4</sup>
  - deterministic simulation testing and extensive fuzzing

*"I want to emphasize, TigerBeetle has probably the strongest testing culture of any database I have worked with." Kyle Kingsbury @ SD25*

---

<sup>1</sup><https://jepsen.io/analyses/tigerbeetle-0.16.11>

<sup>2</sup>Viewstamped Replication

<sup>3</sup>A Study of SSD Reliability in Large-Scale Enterprise Storage Deployments

<sup>4</sup>Protocol Aware Recovery

# TigerBeetle's Design Decisions

- strictly serializable
- highly durable and available<sup>5</sup>
  - network and node faults<sup>6</sup>
  - storage fault model<sup>7</sup>
  - deterministic simulation testing and fuzzing.
- no dependencies, static memory allocation, single core, asserts in production
- optimize for write-heavy and skewed workloads

---

<sup>5</sup><https://jepsen.io/analyses/tigerbeetle-0.16.11>

<sup>6</sup>Viewstamped Replication

<sup>7</sup>Protocol Aware Recovery

Why does tail latency matter?

tail latency = user experience

## Tail latency and parallel calls

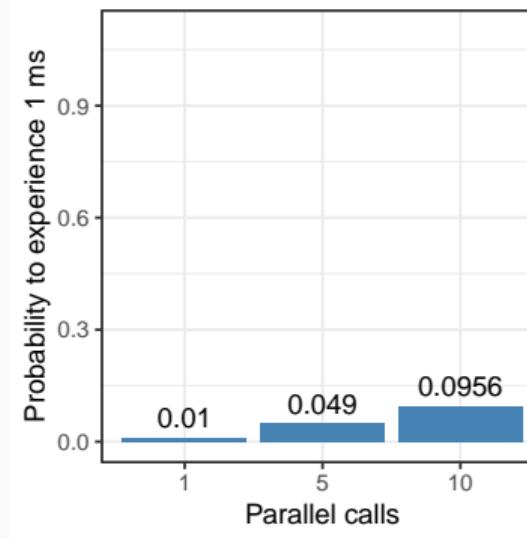
Assume P99 of a system is 1ms.

How high is the probability of hitting this latency with 1, 5, or 10 parallel requests?

## Tail latency and parallel calls

Assume P99 of a system is 1ms.

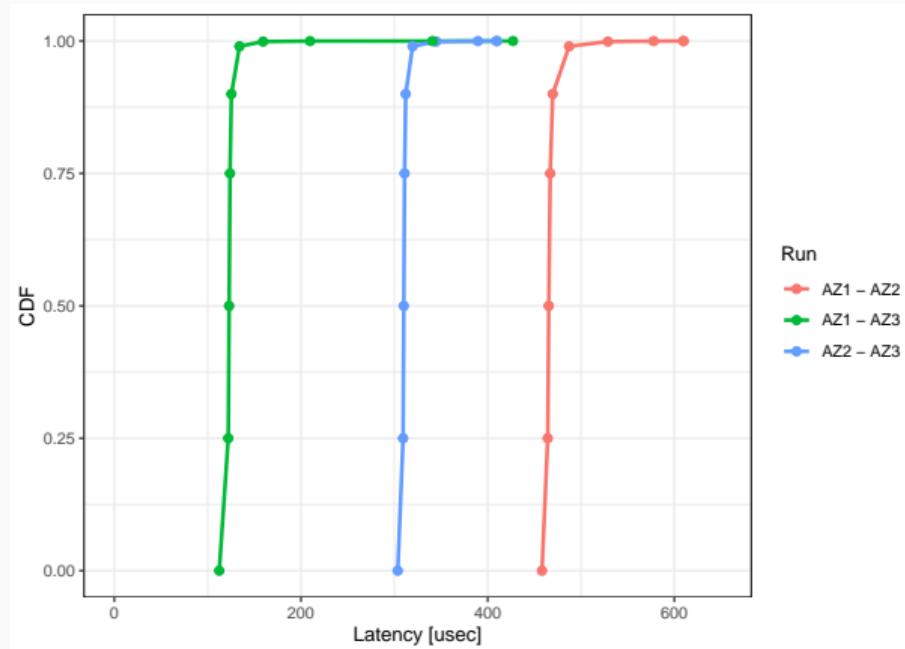
How high is the probability of hitting this latency with 1, 5, or 10 parallel requests?



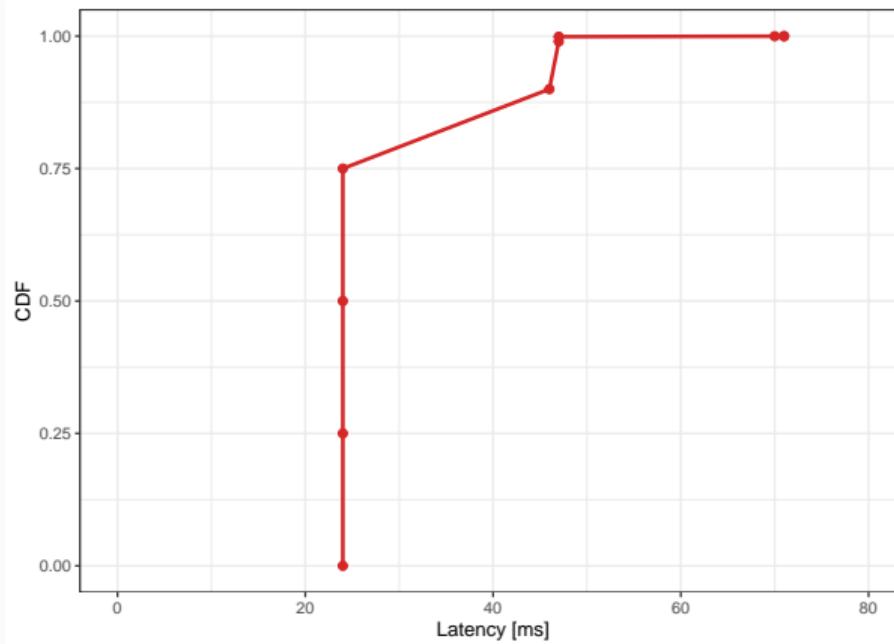
## Act II: The Big Tail in Networks

---

# The Network Jitters - Multi AZ



# The Network Jitters - Cross Region



# The Fallacies of Distributed Computing

1. The network is reliable
2. Latency is zero
3. Bandwidth is infinite
4. The network is secure
5. Topology doesn't change
6. There is one administrator
7. Transport cost is zero
8. The network is homogeneous

## *Viewstamped Replication*

VR provides state machine replication.

# Star Replication

## Quorum Latency

Latency	$2^8$	3	4	5
Mean	49	53	60	71
p50	48	48	60	71
p90	52	66	71	87
p95	58	71	71	94
p99	69	71	81	94
p99.9	71	75	93	105

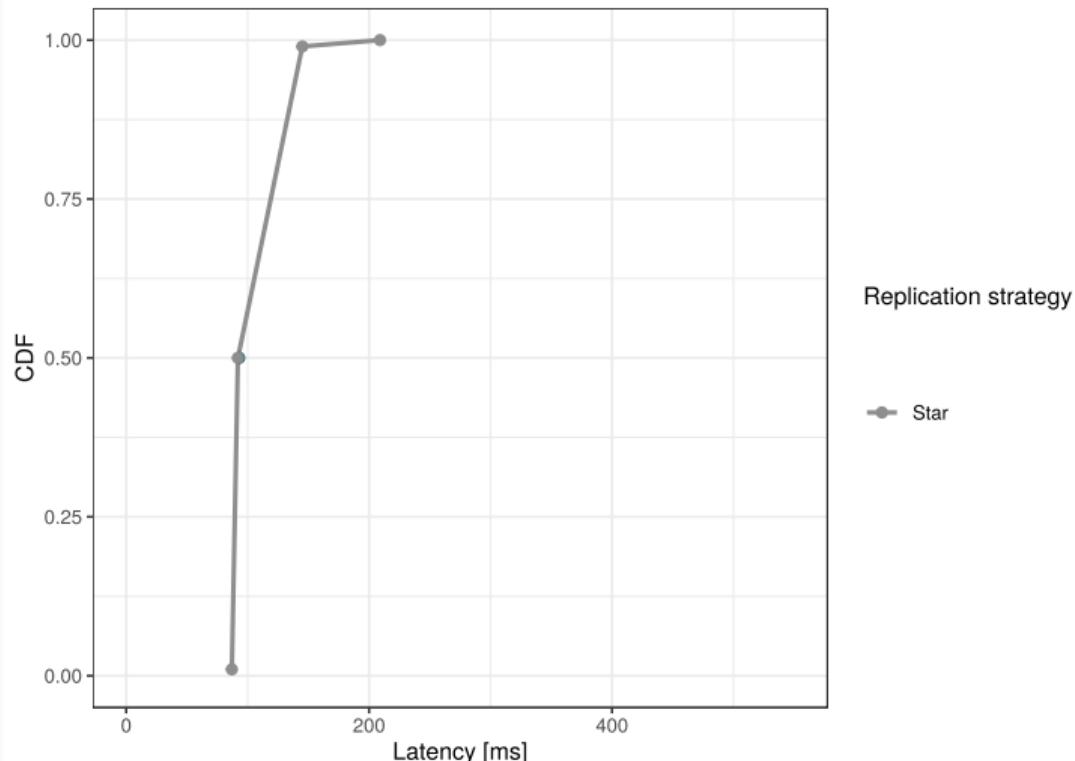
Round-trip latency [ms] statistics for waiting on first N responses.

---

<sup>8</sup>H. Howard, Flexible Paxos: Quorum intersection revisited

# Star Replication [3 regions, 6 x i4i.4xlarge]

Execute 10M transfers.

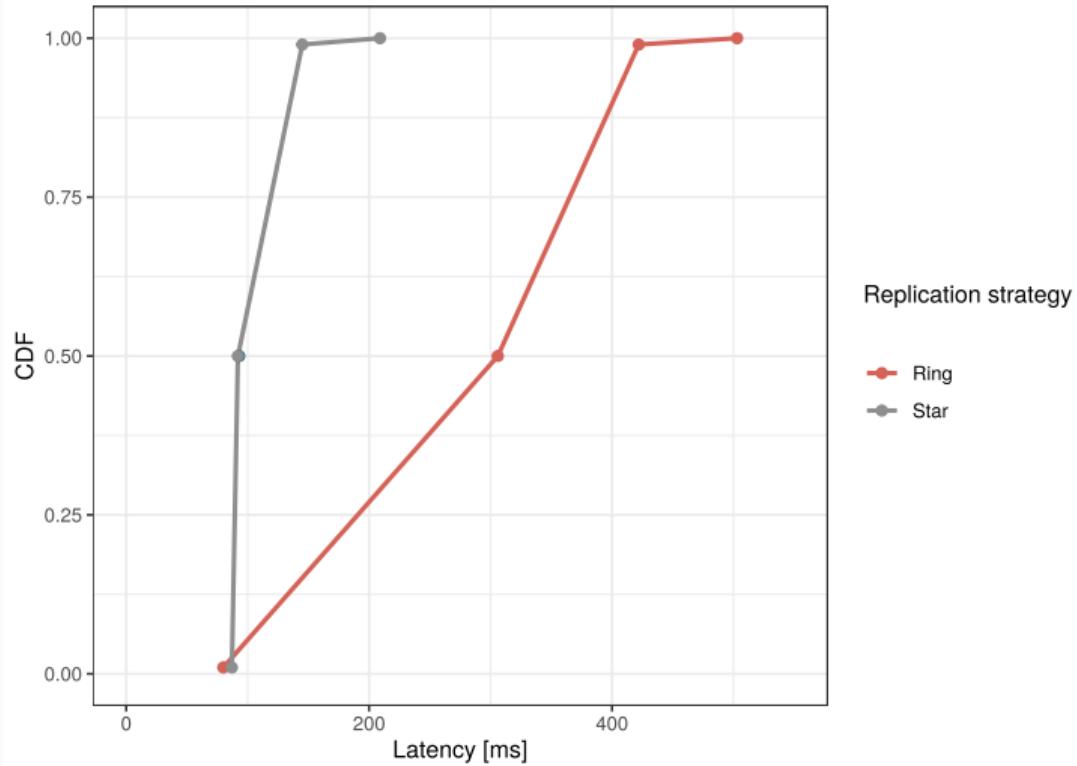


## Star Replication

- + 1 RTT latency (for prepare)
- + masks failures (quorum)
- bandwidth consumption

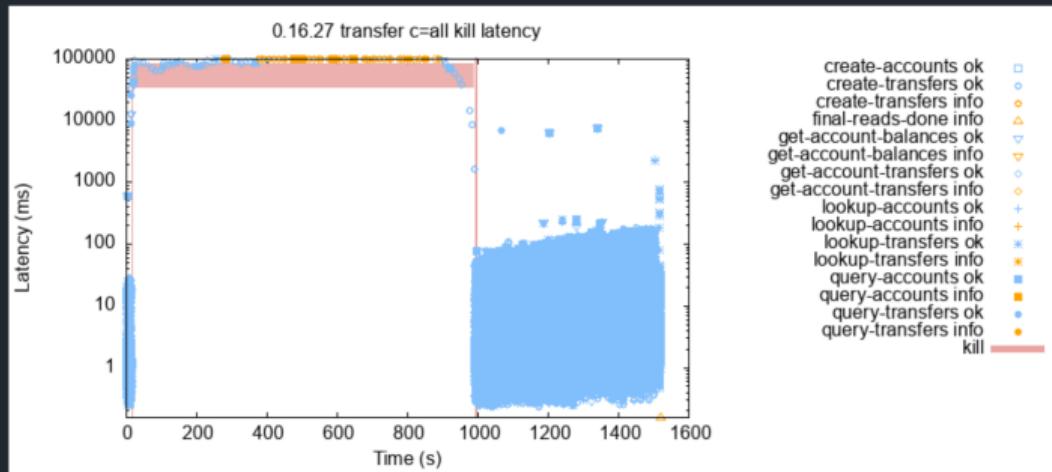
# Ring Replication

# Ring Replication



# Ring Replication and Failures [Issue #2749]

In TigerBeetle 0.16.17 through 0.16.27, single-node failures are frequently associated with higher latencies across all clients--often from three to five orders of magnitude. For instance, consider [this test](#). We killed one of three nodes, and saw latencies jump from 1-50 ms up to ~100 seconds. Latencies remained elevated until the node was restarted, almost a thousand seconds later.



High latencies may recover spontaneously after tens to hundreds of seconds, or persist for thousands of seconds.

This may involve the failure of a non-leader node. It also seems more likely in smaller clusters, rather than large ones.

# Ring Replication

- + bandwidth optimized
- latency is not great
- masking faults is harder and requires retries<sup>9</sup>

---

<sup>9</sup>Try again: The tools and techniques behind resilient systems

Can we find a better trade-off?

	Star	Ring	?
Latency	++	-	
Bandwidth	-	++	
Failures	++	-	

## Hybrid Replication

Goal: Fast common case latency, no message critical, bandwidth more balanced.

Good, right?



## The Fallacies strike back!

1. The network is reliable
2. Latency is zero
3. Bandwidth is infinite
4. The network is secure
5. Topology doesn't change
6. There is one administrator
7. Transport cost is zero
8. The network is homogeneous

## Why is it bad?

Gray Failure: The Achilles' Heel of Cloud-Scale Systems

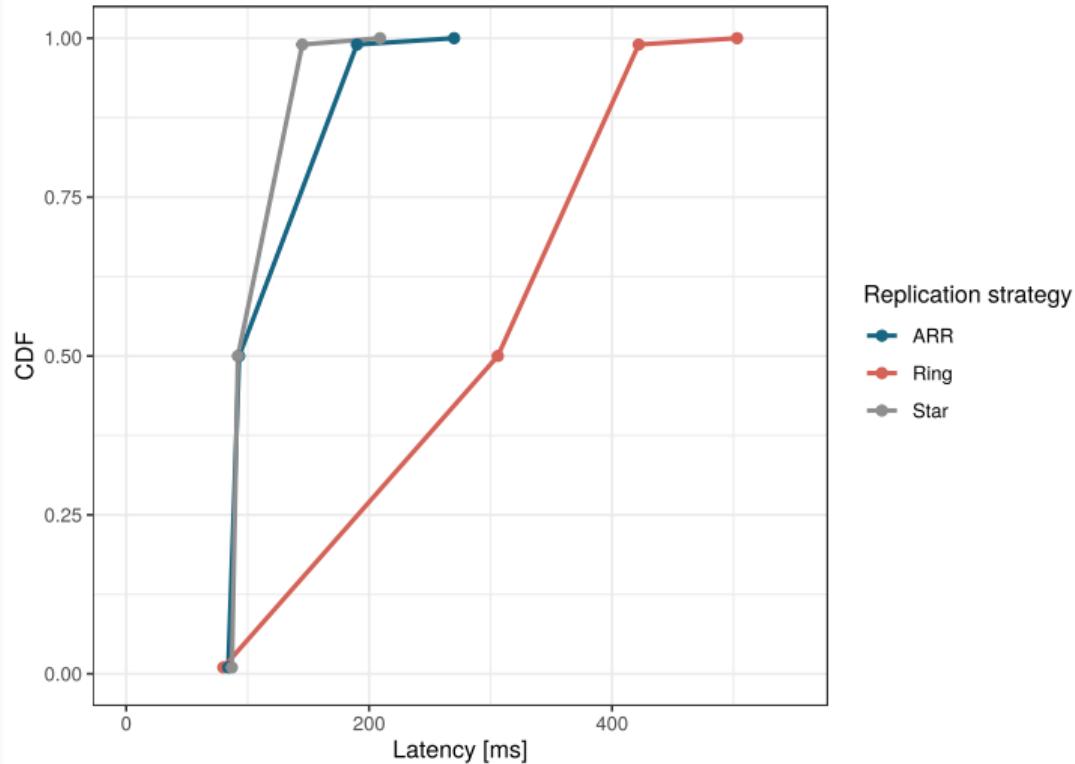
## *Adaptive Replication Routing*

Inspired laterally by congestion control<sup>10</sup>

---

<sup>10</sup>PCC: Re-architecting Congestion Control for Consistent High Performance

# Adaptive Replication Routing



# Taming The Tail But at What Cost?

	Star	Ring	ARR <sup>11</sup>
Latency	++	-	+
Bandwidth	-	++	+
Failures	++	-	+
Topology	++	-	++
Transport Cost	-	0	++

---

<sup>11</sup>Fuzzer master class: <https://tigerbeetle.com/blog/2025-11-28-tale-of-four-fuzzers/>

Remember the Fallacies when building distributed systems.

## Act III: Tail Latency Inside One Beetle

---



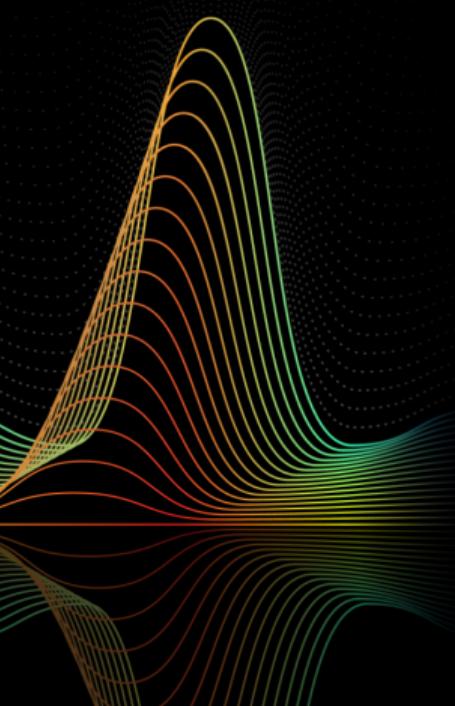
# The Tale of Taming TigerBeetle's Tail Latency



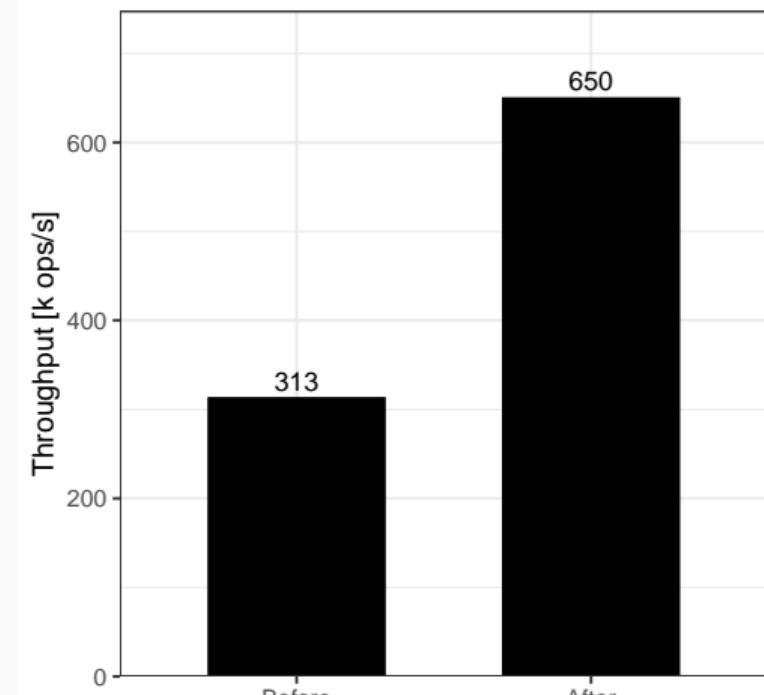
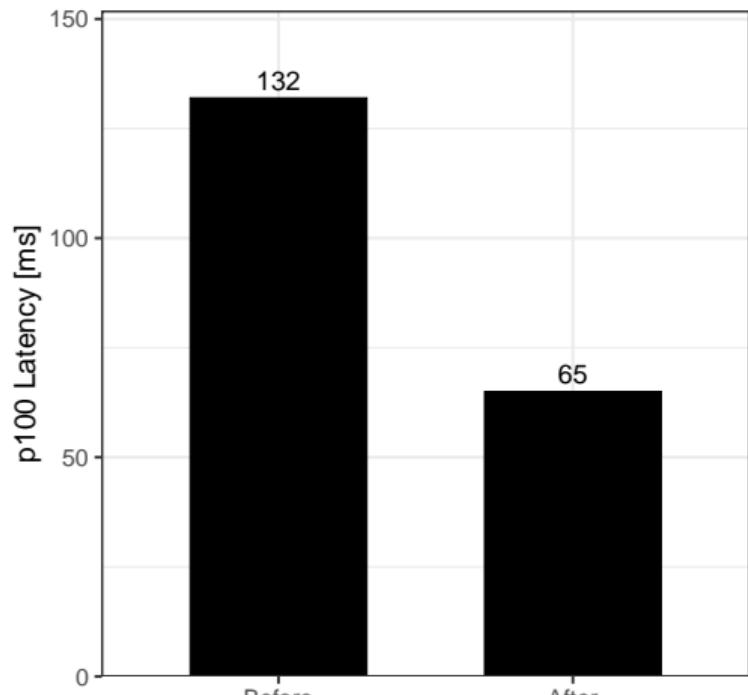
Tobias Ziegler  
*Software Engineer*



*TigerBeetle*



# Results



# Recording

P99 CONF 2025 | The Tale of Taming TigerBeetle's Tail Latency by Tobias Ziegler

