

### Zeus

### Locality-aware distributed transactions



A. Katsarakis<sup>\*</sup>, Y. Ma<sup>†</sup>, Z. Tan<sup>§</sup>, A. Bainbridge, M. Balkwill,
 A. Dragojevic, B. Grot<sup>\*</sup>, B. Radunovic, Y. Zhang
 \*University of Edinburgh, <sup>†</sup>Fudan University, <sup>§</sup>UCLA, Microsoft Research

zeus-protocol.com

Thanks to:

**EPSRC** Microsoft Research



Keep data in-memory, replicated, sharded across nodes of a datacenter



distributed datastore

















Traditional distributed txs well-known as expensive



#### Many tx applications exhibit **dynamic locality** network functions, peer-to-peer payments ...



Many tx applications exhibit **dynamic locality** network functions, peer-to-peer payments ...

Example: cellular control plane manages phone connectivity and handovers among base stations



base-station B



Many tx applications exhibit **dynamic locality** network functions, peer-to-peer payments ...

Example: cellular control plane manages phone connectivity and handovers among base stations

#### Locality

every phone user repeats txs: same phone & nearest base-station







Many tx applications exhibit **dynamic locality** network functions, peer-to-peer payments ...

Example: cellular control plane manages phone connectivity and handovers among base stations

#### Locality

every phone user repeats txs: same phone & nearest base-station

#### But locality is dynamic

changes at run-time e.g., user commutes  $\rightarrow$  base-station changes





Many tx applications exhibit **dynamic locality** network functions, peer-to-peer payments ...

Example: cellular control plane manages phone connectivity and handovers among base stations

#### Locality

every phone user repeats txs: same phone & nearest base-station

#### But locality is dynamic

changes at run-time







Static sharding (e.g., consistent hashing)

Objects placed randomly on fixed nodes

- b easy to locate and access objects
- C reliable txs regardless of access pattern





tx coordinator



Static sharding (e.g., consistent hashing)

Objects placed randomly on fixed nodes

- b easy to locate and access objects
- 🕐 reliable txs regardless of access pattern





tx coordinator

expensive reliable txs

mostly **remote accesses** some blocking (control flow, pointer chasing)





**Static sharding** (e.g., consistent hashing)

Objects placed randomly on fixed nodes

easy to locate and access objects

reliable txs regardless of access pattern



tx coordinator

#### expensive reliable txs

mostly remote accesses some blocking (control flow, pointer chasing) related objects on different shards

costly distributed commit





**Static sharding** (e.g., consistent hashing)

Objects placed randomly on fixed nodes

- easy to locate and access objects
- reliable txs regardless of access pattern



tx coordinator



#### expensive reliable txs

mostly remote accesses some blocking (control flow, pointer chasing)

related objects on different shards costly distributed commit



Cannot exploit locality  $\rightarrow$  expensive reliable txs





Inspired by multiprocessor's hardware transactional memory





Inspired by multiprocessor's hardware transactional memory



Each object has a single node **owner** = **data** + **exclusive write access** the **owner changes dynamically** and is tracked by replicated directory





Inspired by multiprocessor's hardware transactional memory



Each object has a single node **owner** = **data** + **exclusive write access** the **owner changes dynamically** and is tracked by replicated directory

Coordinator executes a tx by **acquiring ownership** of all its objects → **single-node commit** 





Inspired by multiprocessor's hardware transactional memory



Each object has a single node **owner** = **data** + **exclusive write access** the **owner changes dynamically** and is tracked by replicated directory

Coordinator executes a tx by **acquiring ownership** of all its objects → **single-node commit** 

Ownership stays with coordinator

 $\rightarrow$  future txs on these objects enjoy local accesses





Inspired by multiprocessor's hardware transactional memory



Each object has a single node **owner** = **data** + **exclusive write access** the **owner changes dynamically** and is tracked by replicated directory

Coordinator executes a tx by **acquiring ownership** of all its objects → **single-node commit** 

Ownership stays with coordinator → future txs on these objects enjoy local accesses

What are the exact steps?



2.

1.



#### **1.** Execute as the owner

a) at object access: if (not owner) get ownership

2.



#### 1. Execute as the owner

a) at object access: if (not owner) get ownership

How to get ownership reliably?

2.







1) Coordinator gets ownership from current owner







Coordinator gets ownership from current owner
 Keeps consistent directory replicas \_\_\_\_\_\_arbiters

D

directory

replica



- Coordinator gets ownership from current owner
   Keeps consistent directory replicas
  - 1. Coordinator sends object ownership **invalidations** (through a directory replica) to all arbiters





- 1) Coordinator gets ownership from current owner
- 2) Keeps consistent directory replicas -
  - 1. Coordinator sends object ownership invalidations (through a directory replica) to all arbiters
  - 2. Arbiters acknowledge the coordinator directly





- 1) Coordinator gets ownership from current owner
- 2) Keeps consistent directory replicas
  - 1. Coordinator sends object ownership invalidations (through a directory replica) to all arbiters
  - 2. Arbiters acknowledge the coordinator directly

Ownership is acquired & coordinator proceeds with tx 🗸





- 1) Coordinator gets ownership from current owner
- 2) Keeps consistent directory replicas ———
  - 1. Coordinator sends object ownership **invalidations** (through a directory replica) to all arbiters
  - 2. Arbiters acknowledge the coordinator directly
  - 3. Coordinator sends validations informing arbiters for acquisition





arbiters

- 1) Coordinator gets ownership from current owner
- 2) Keeps consistent directory replicas ———
  - 1. Coordinator sends object ownership invalidations (through a directory replica) to all arbiters
  - 2. Arbiters acknowledge the coordinator directly
  - 3. Coordinator sends validations informing arbiters for acquisition

# **Conflicts**: logical timestamps, **fault tolerance**: idempotent replays **as in Hermes** [ASPLOS'20]





arbiters

- 1) Coordinator gets ownership from current owner
- 2) Keeps consistent directory replicas
  - 1. Coordinator sends object ownership **invalidations** (through a directory replica) to all arbiters
  - 2. Arbiters acknowledge the coordinator directly
  - 3. Coordinator sends validations informing arbiters for acquisition

Conflicts: logical timestamps, fault tolerance: idempotent replays
Correctness verified under conflicts and faults



#### 1. Execute as the owner

a) at object access: if (not owner) get ownership

Locality + ownership stays with coordinator



#### 1. Execute as the owner

a) at object access: if (not owner) get ownership

b) local access

common

case

Locality + ownership stays with coordinator



#### 1. Execute as the owner

a) at object access: if (not owner) get ownership

b) local access

common

case



#### 1. Execute as the owner

a) at object access: if (not owner) get ownership

🍗 b) **local access** 

common

case

#### 2. Local commit

commits tx: traditional single-node commit



#### 1. Execute as the owner

a) at object access: if (not owner) get ownershipb) local access

common

case

#### 2. Local commit

**commits tx**: traditional single-node commit

#### 3. Reliable commit

**completes tx**: updating replicas



#### 1. Execute as the owner

a) at object access: if (not owner) get ownership

≽ b) **local access** 

common

case

#### 2. Local commit

commits tx: traditional single-node commit

#### 3. Reliable commit

completes tx: updating replicas

#### Great! But how efficient is reliable commit?





#### **1.** Committed tx $\rightarrow$ no conflicts $\rightarrow$ fast tx completion

- coordinator sends updates to replicas and waits for ACKs
- read-only txs: no updates  $\rightarrow$  no reliable commit





#### **1.** Committed tx $\rightarrow$ no conflicts $\rightarrow$ fast tx completion

- coordinator sends updates to replicas and waits for ACKs
- read-only txs: no updates  $\rightarrow$  no reliable commit



#### 2. No conflicts → no aborts → pipelined txs (no waiting for replication)

- subsequent txs use local state with certainty & issue updates
- coordinator sequences updates, which replicas apply in order





#### **1.** Committed tx $\rightarrow$ no conflicts $\rightarrow$ fast tx completion

- coordinator sends updates to replicas and waits for ACKs
- read-only txs: no updates  $\rightarrow$  no reliable commit



2. No conflicts → no aborts → pipelined txs (no waiting for replication)

- subsequent txs use local state with certainty & issue updates
- coordinator sequences updates, which replicas apply in order

Fault tolerance: idempotent replays





#### **1.** Committed tx $\rightarrow$ no conflicts $\rightarrow$ fast tx completion

- coordinator sends updates to replicas and waits for ACKs
- read-only txs: no updates  $\rightarrow$  no reliable commit



2. No conflicts → no aborts → pipelined txs (no waiting for replication)

- subsequent txs use local state with certainty & issue updates
- coordinator sequences updates, which replicas apply in order

Very efficient! Correctness verified under faults





















#### Locality-aware, distributed and reliable



#### Awesome! Does it translate into performance?











Zeus: within 9% of ideal





Zeus: within 9% of ideal







6 nodes, 3-way replication, Zeus 40Gb (no RDMA)



Paper: more benchmarks, ownership, latency ...

3



**State-of-the-art** reliable txs over **static sharding**: cannot exploit dynamic locality

remote accesses

costly distributed commit



State-of-the-art reliable txs over static sharding: cannot exploit dynamic locality remote accesses costly distributed commit

Zeus' reliable txs exploit locality via dynamic ownership: local accesses in the common case single-node commit

- local for read-only txs
- fast and pipelined for write txs





**zeus-protocol.com** TLA<sup>+</sup> specification, Q&A ...



State-of-the-art reliable txs over static sharding: cannot exploit dynamic locality remote accesses costly distributed commit

Zeus' reliable txs exploit locality via dynamic ownership: local accesses in the common case single-node commit

- local for read-only txs
- fast and pipelined for write txs

Performance 10s millions txs/second up to 2x state-of-the-art



**zeus-protocol.com** TLA<sup>+</sup> specification, Q&A ...



State-of-the-art reliable txs over static sharding: cannot exploit dynamic locality remote accesses costly distributed commit

Zeus' reliable txs exploit locality via dynamic ownership: local accesses in the common case single-node commit

- local for read-only txs
- fast and pipelined for write txs

Performance 10s millions txs/second up to 2x state-of-the-art

**Bonus: programmability!** 







**<u>zeus-protocol.com</u>** TLA<sup>+</sup> specification, Q&A ...



State-of-the-art reliable txs over static sharding: cannot exploit dynamic locality remote accesses costly distributed commit

Zeus' reliable txs exploit locality via dynamic ownership: local accesses in the common case single-node commit

- local for read-only txs
- fast and pipelined for write txs

Performance 10s millions txs/second up to 2x state-of-the-art

Bonus: programmability!





zeus-protocol.com TLA<sup>+</sup> specification, Q&A ...

Reliable txs with locality? Use Zeus!

