

ERMIA: Fast Memory-Optimized Database System for Heterogeneous Workloads <https://github.com/ermia-db/ermia>

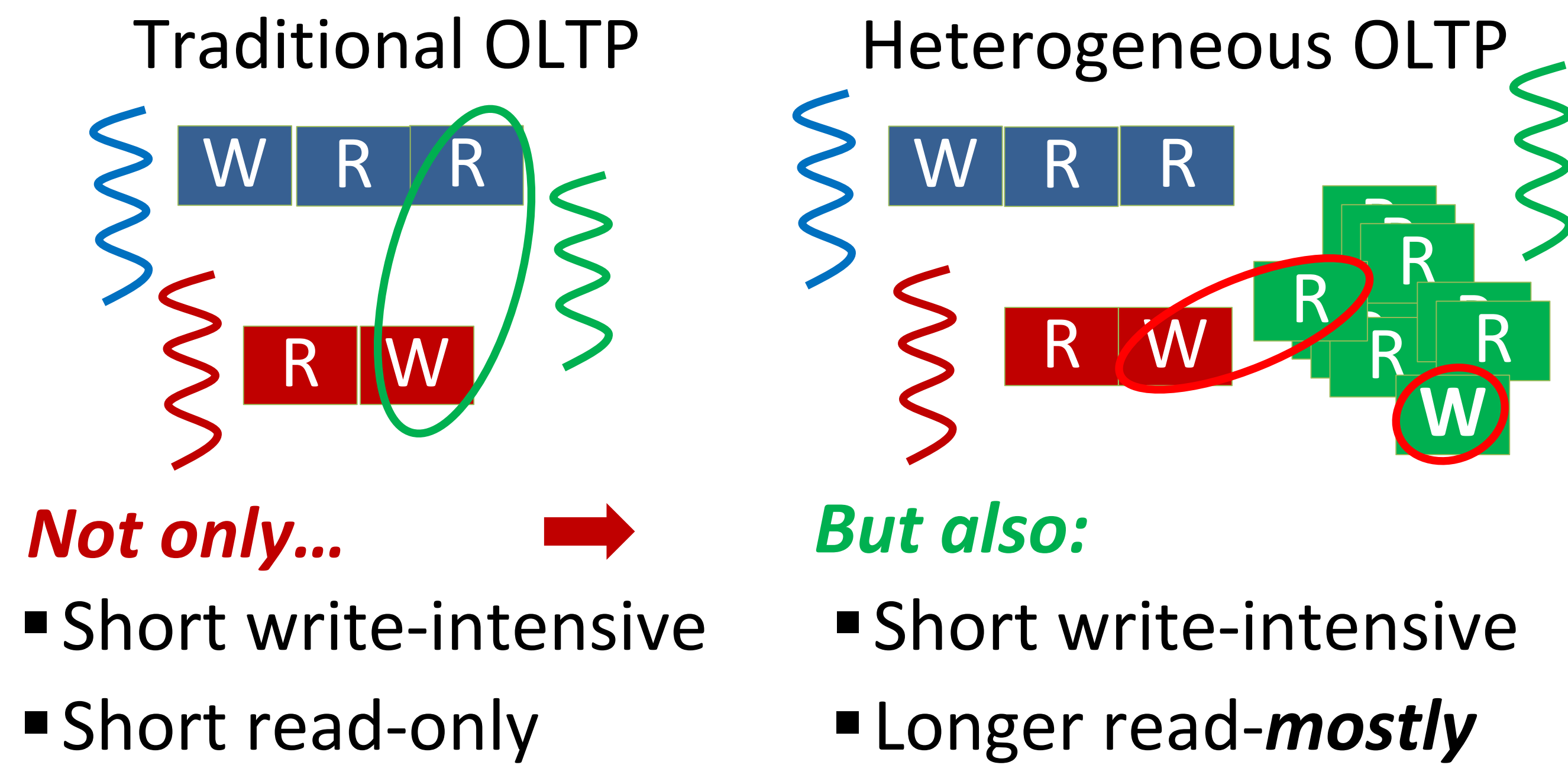
Kangyeon Kim Tianzheng Wang
University of Toronto

Ryan Johnson
LogicBlox

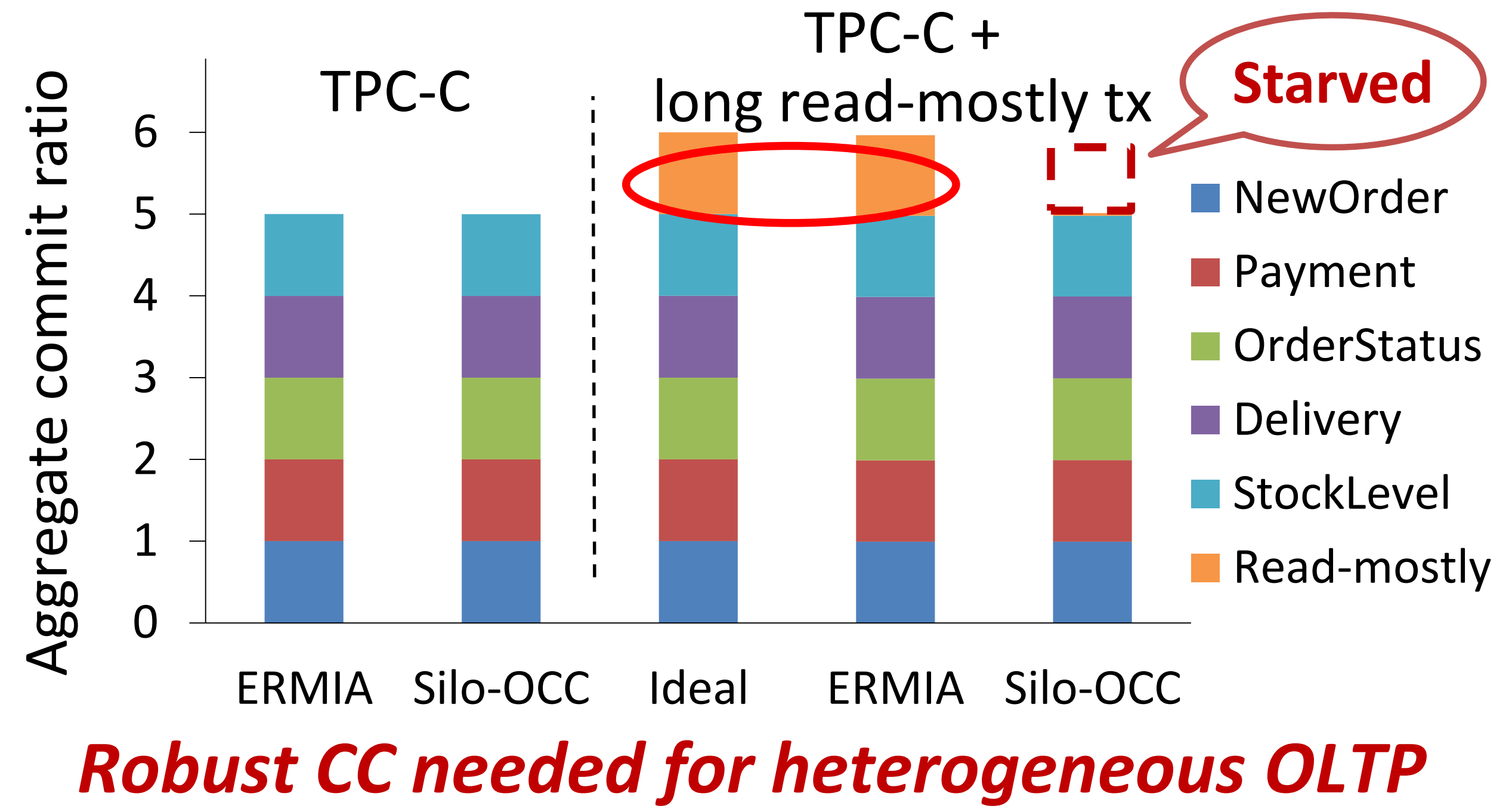
Ippokratis Pandis
Amazon Web Services

What? *Heterogeneous workloads* are coming, but existing MMDBMS isn't good at them
Why? "Wrong" concurrency control in use: *aborts* too many *read-mostly* transactions
How? *Fair and robust CC* (snapshot isolation + certifier) + *scalable physical layer*

Read-mostly analytical components



OCC: not always the best



ERMIA = Snapshot Isolation + Serial Safety Net + Scalable Physical Layer

Fair and robust logical layer

- Read-friendly snapshot isolation**
 - Abort on write-write conflicts, preserves reads
- Serializability with Serial Safety Net ***
 - Cheap certifier on top of any CC \geq RC (e.g. SI)
 - Maintains fairness and robustness

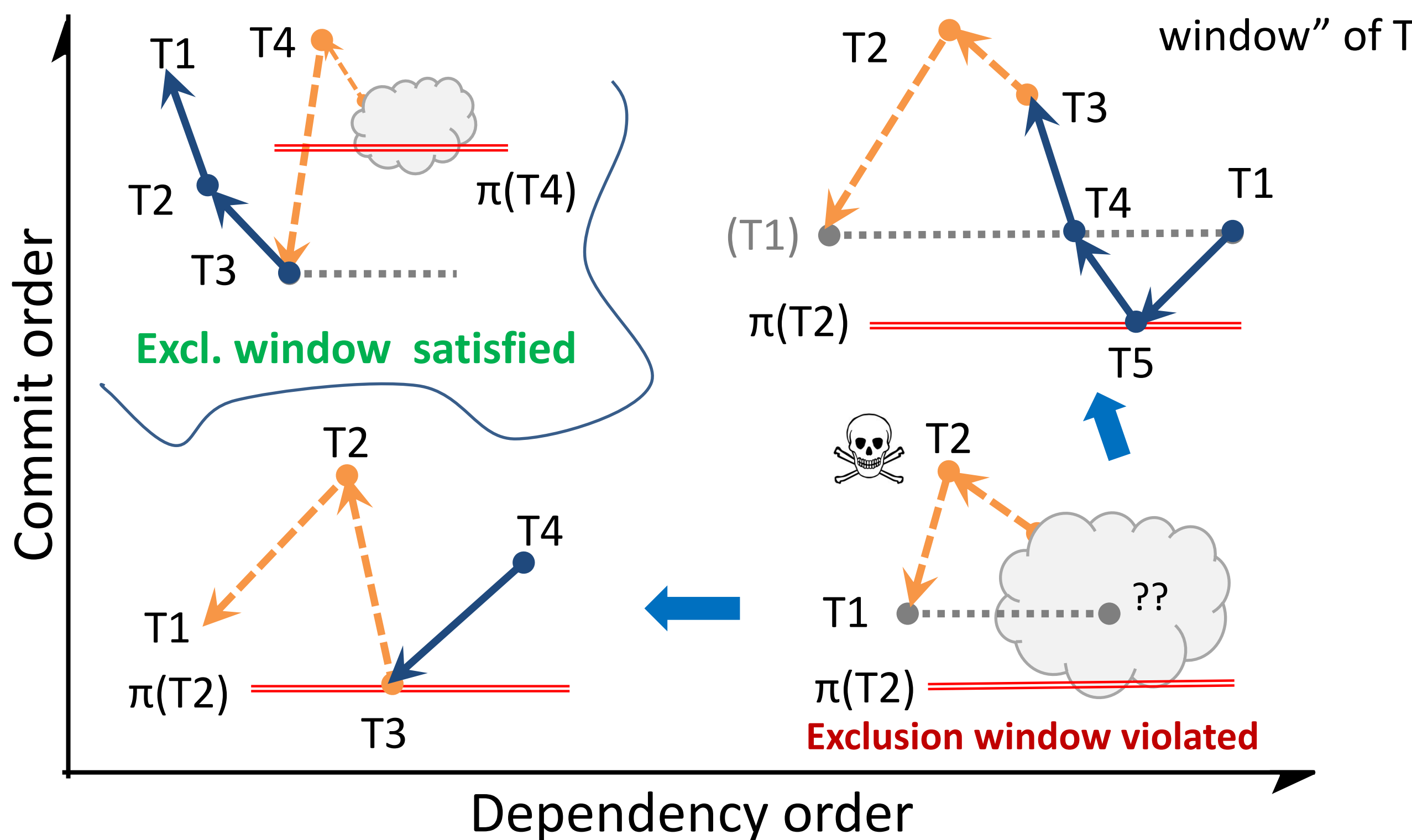
* T. Wang, R. Johnson, A. Fekete, I. Pandis. "Serial Safety Net: Efficient Concurrency Control on Modern Hardware", DaMoN '15

Scalable physical layer

- Minimal global communication**
 - One atomic-fetch-add per tx for global ordering
 - Eases implementation of snapshot isolation
 - Simplifies logging/recovery
- Easy maintenance via indirection**
 - Fast recovery, single-hop index update, etc.

Serial Safety Net

- def:** $c(T)$ = T's commit time, $\pi(T)$ = *earliest* successor
- Forbid** any pred P with $\pi(T) \leq c(P) \leq c(T)$ → "Exclusion window" of T

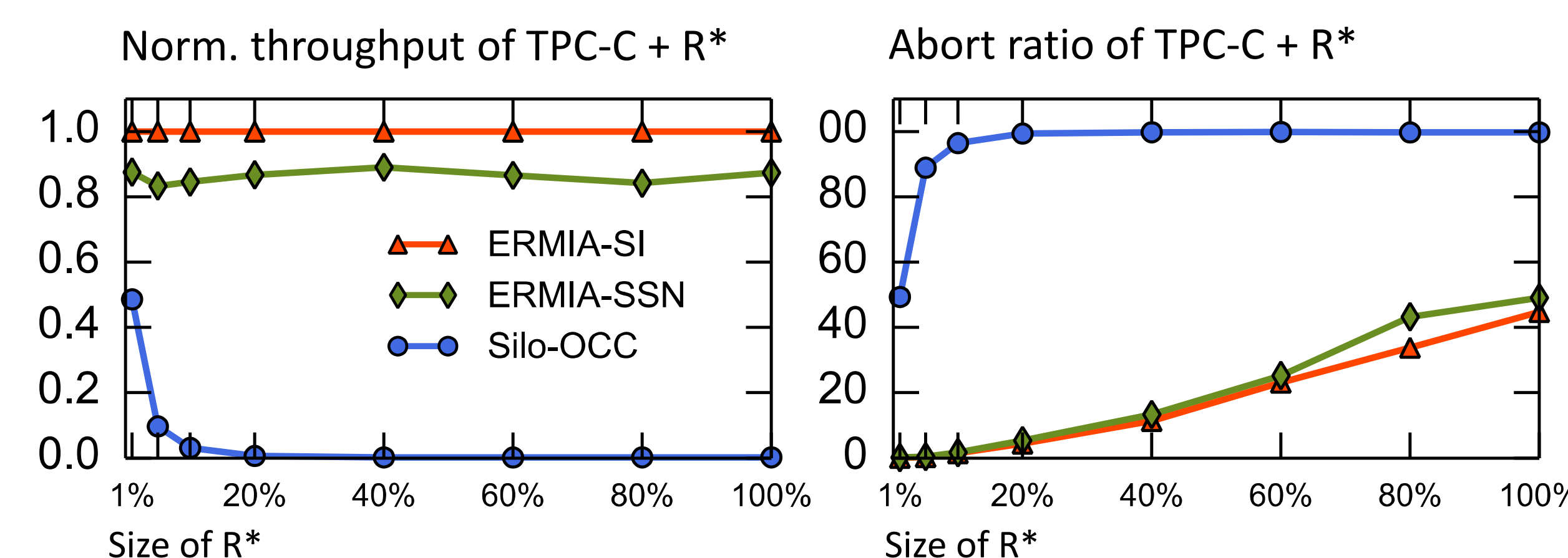
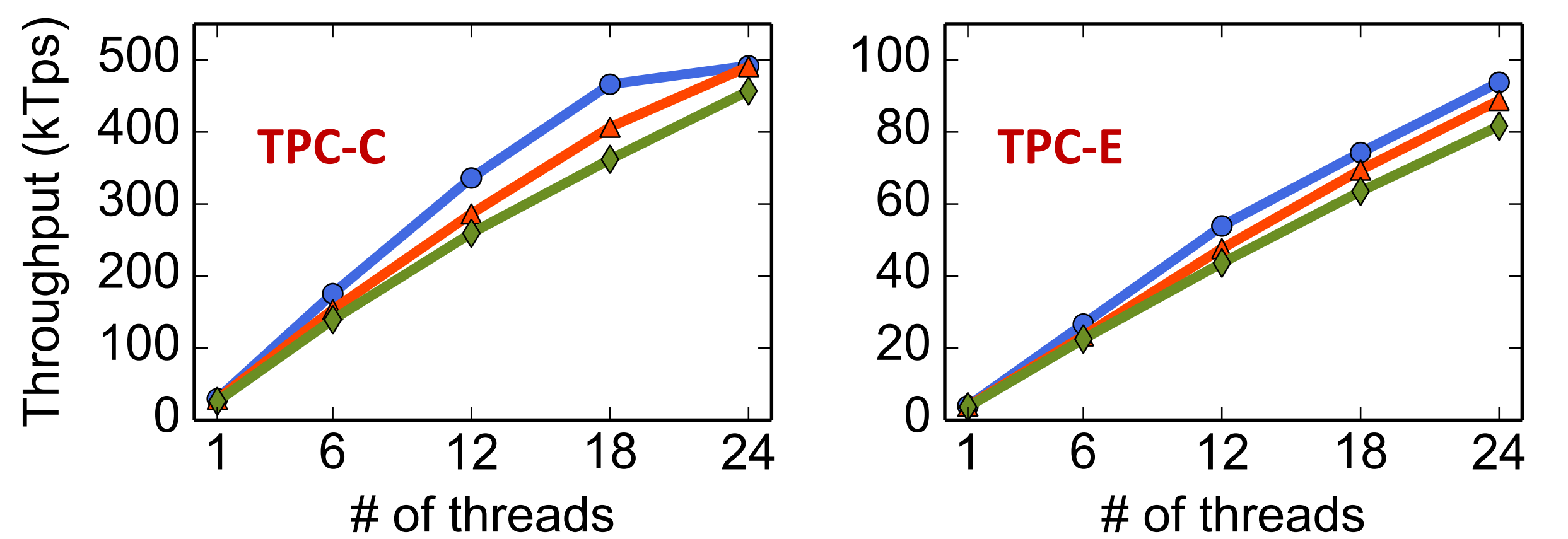


Indirection array

- Object IDs** (instead of pointers) at leaf level
 - Updates:** no index update needed
-
- Install new version: CAS V1 → V2
CAS failure → WW conflict
- Recovery:** load header information only

Robust to "convenient" & real workloads

HW: 4-socket 6-core Intel E7-4807, 64GB RAM



Scalable centralized redo logging

- Private footprint until commit**
 - Carve out space via atomic-fetch-add**
 - Fill in log space asynchronously
 - Abort → log discarded
 - LSN offset = global order**
- Commit LSN (cLSN) = XADD(current LSN, log size)
-