# Deconstructing RDMA-enabled Distributed Transaction Processing: Hybrid is Better!

**Xingda Wei,** Zhiyuan Dong, Rong Chen, Haibo Chen

Institute of Parallel and Distributed Systems (**IPADS**)

Shanghai Jiao Tong University

# Remote Direct Memory Access (RDMA)

Kernel bypassing network

⇨ Ultra **low** latency~(5us)

⇨ Ultra **high** throughput

Offloading technology (one-sided)

⇨ **Bypassing** CPU

⇨ **Read/Write, CAS**[2] server's memory

Gain interests from **Academia** & **Industry**

⇨ **Orders of magnitude improvements** on distributed applications

⇨ Available in the **public cloud**[1]

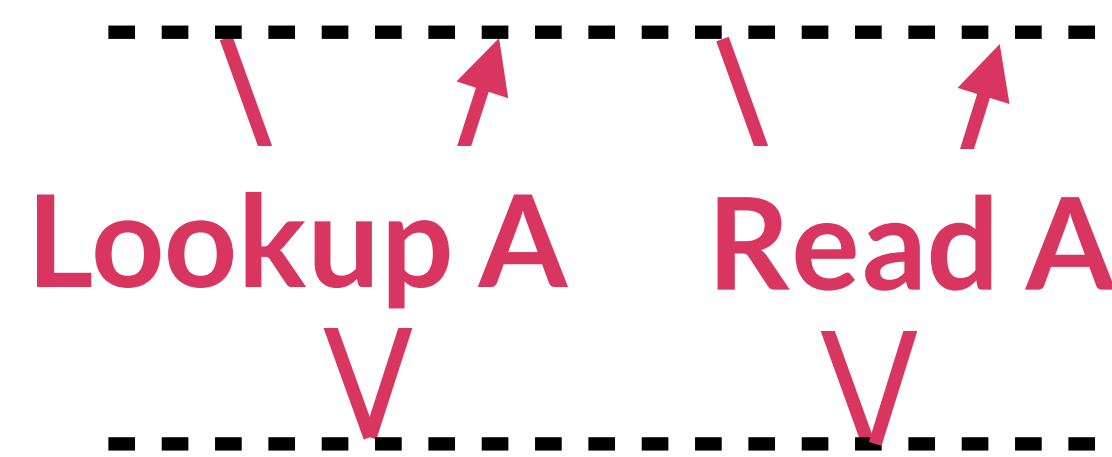[1] https://azure.microsoft.com/en-us/blog/azure-linux-rdma-hpc-available/

[2] Atomic compare and swap

# On-going debate over how to use RDMA for TXs

**Get(A)**

*Coordinator*

*A's store*

**One-sided READ(I)**

Lookup A    Read A

**Two-sided RPC(II)**

RPC request    RPC reply

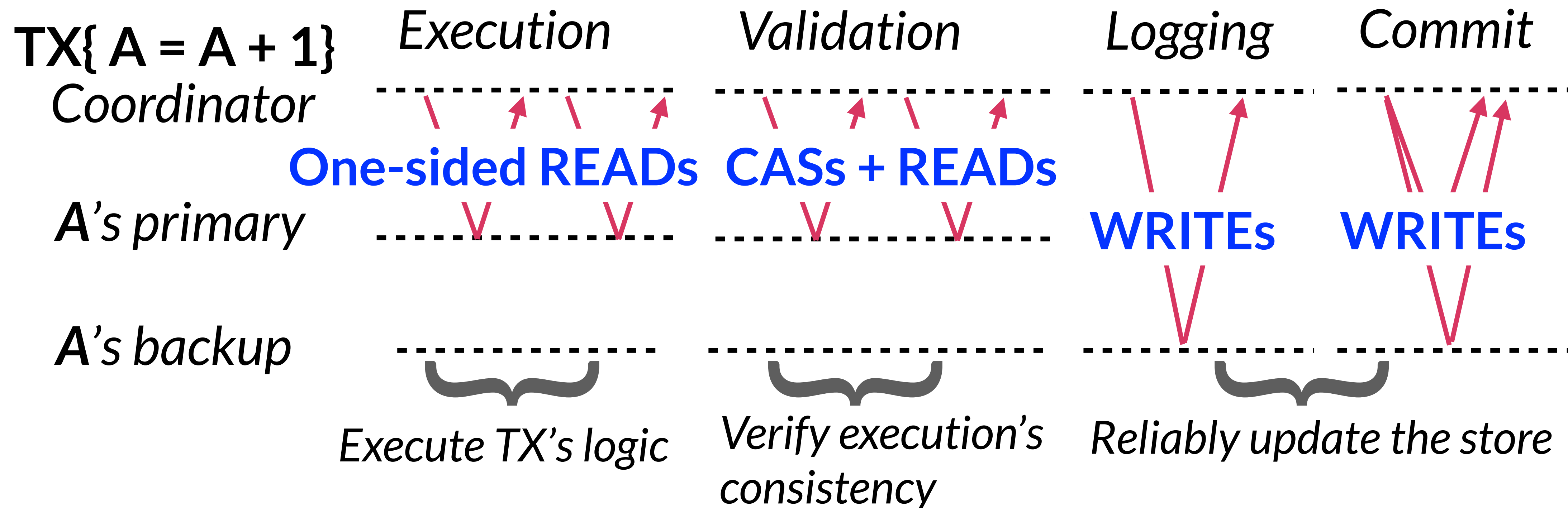| | One-sided READ(I) | Two-sided RPC(II) |
|---|---|---|
| Performance | ✔ | ✘ |
| #Round-trips | >= 2 | 1 |

# Transaction(TX)s are more complex

TX (e.g. OCC[1]) uses multiple **phases** for **serializability** & **availability**

⇨ Each **can be offloaded** w one-sided primitive

| TX{ A = A + 1 } | *Execution* | *Validation* | *Logging* | *Commit* |
|---|---|---|---|---|
| *Coordinator* | | | | |
| | **One-sided READs** | **CASs + READs** | | |
| *A's primary* | | | **WRITEs** | **WRITEs** |
| *A's backup* | | | | |
| | *Execute TX's logic* | *Verify execution's consistency* | *Reliably update the store* | |

[1] Optimistic concurrency control

# Transaction(TX)s are more complex

Protocols

⇨ OCC, 2PL, SI, ....

Impl on hardware devices

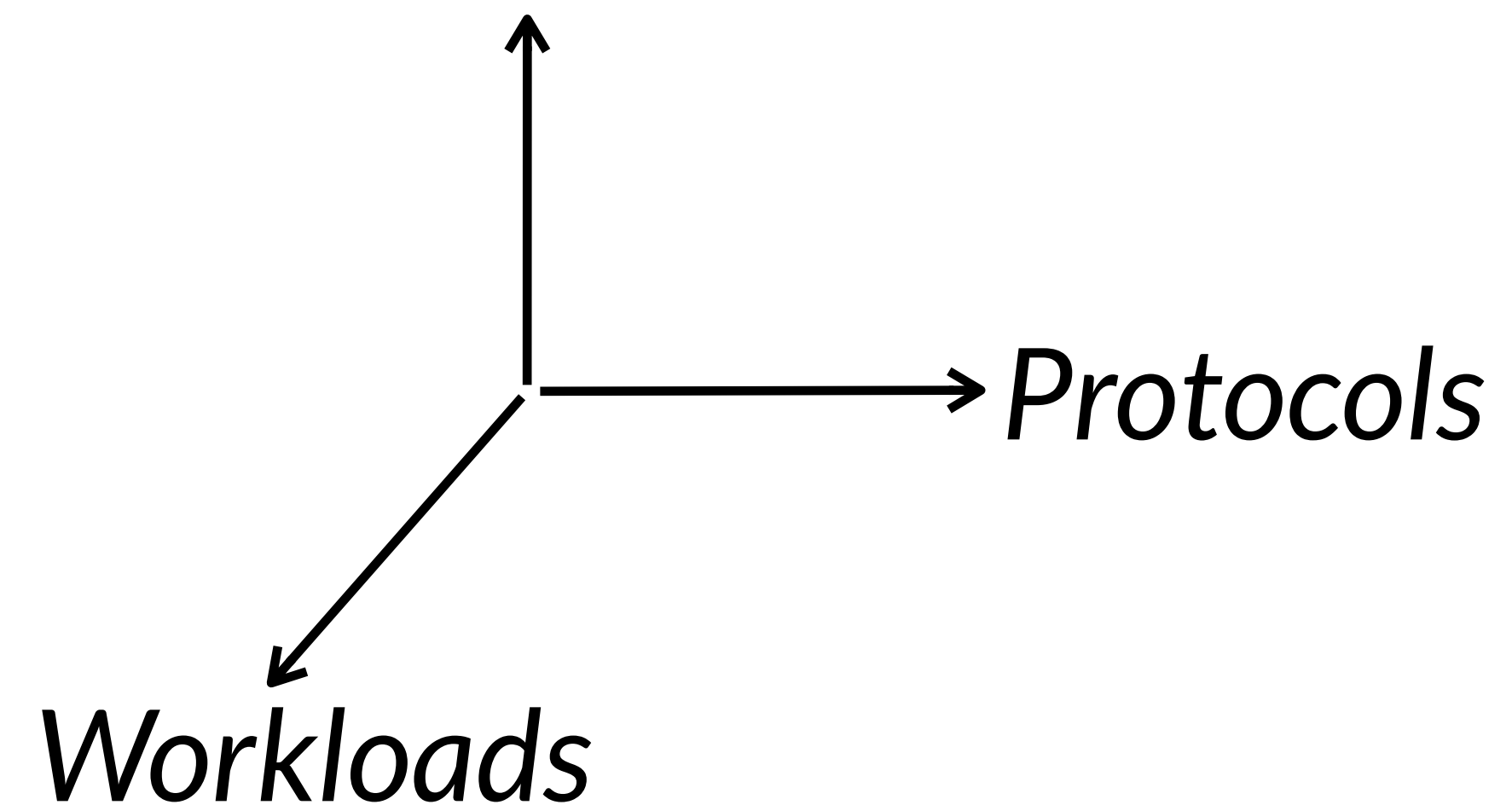⇨ One-sided vs. Two-sided, …

⇨ CX3, CX4, CX5, ROCE

OLTP workloads

⇨ TPC-C, TPC-E, TATP, Smallbank, …

Implements & Hardware

*Protocols*

*Workloads*

# This work: how to use RDMA for TXs

Focus on **OCC** in this work

⇨  Use **phase-by-phase** approach

**Well-tuned** RDMA execution framework

⇨  Representative RNICs (CX3 - CX5)

Representative OLTP workloads

⇨  TPC-C, TPC-E, and Smallbank

Implements & Hardware

**O**ptimistic **C**oncurrency **C**ontrol

Widely used in

Centralized

Silo[SOSP'13]   Foedus[SIGMOD'15]
...

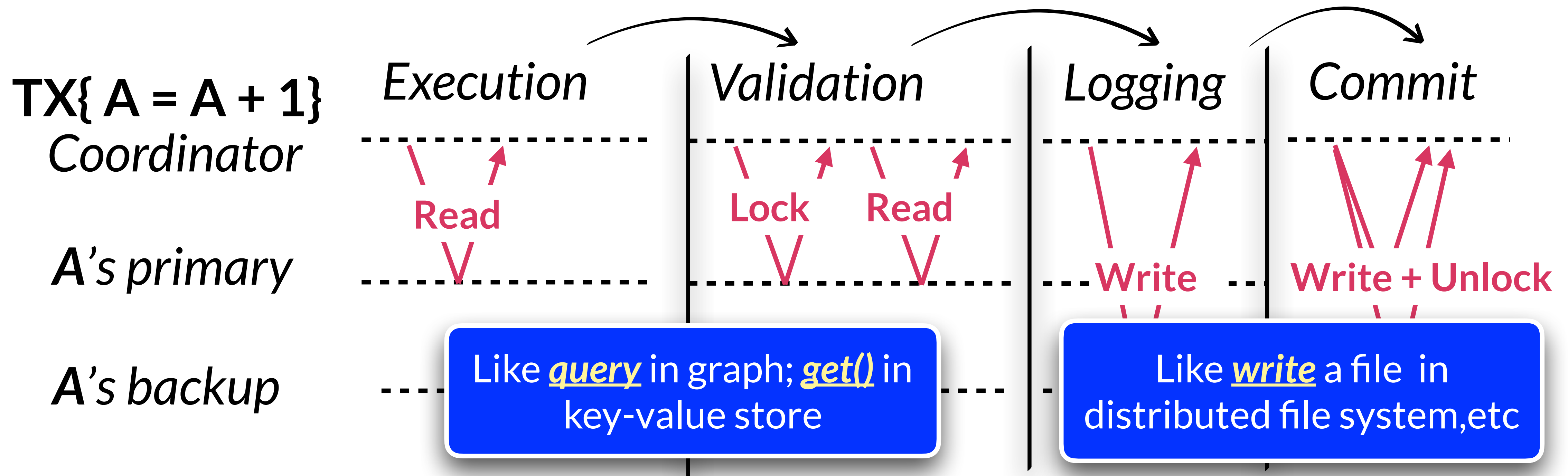Distributed

FaRM[SOSP'15]   TAPIR[SOSP'15]
...

*Workload*

# Phase-by-phase analysis is effective & useful

OCC uses **consecutive** phases

⇨ Better **phase** performance -> Better **overall** performance

**TX{ A = A + 1}**

| *Execution* | *Validation* | *Logging* | *Commit* |

*Coordinator*

Read        Lock    Read

*A's primary*                                    Write    Write + Unlock

*A's backup*

Like *query* in graph; *get()* in key-value store

Like *write* a file in distributed file system,etc

# Deconstructing TX with phase-by-phase analysis

OCC uses **consecutive** phases

  ▷  Better **phase** performance -> Better **overall** performance

| Execution | Validation | Logging | Commit |

🎓 **DrTM+H**    <span style="background:blue; color:white">_**No** single primitive wins **all** the time !_</span>
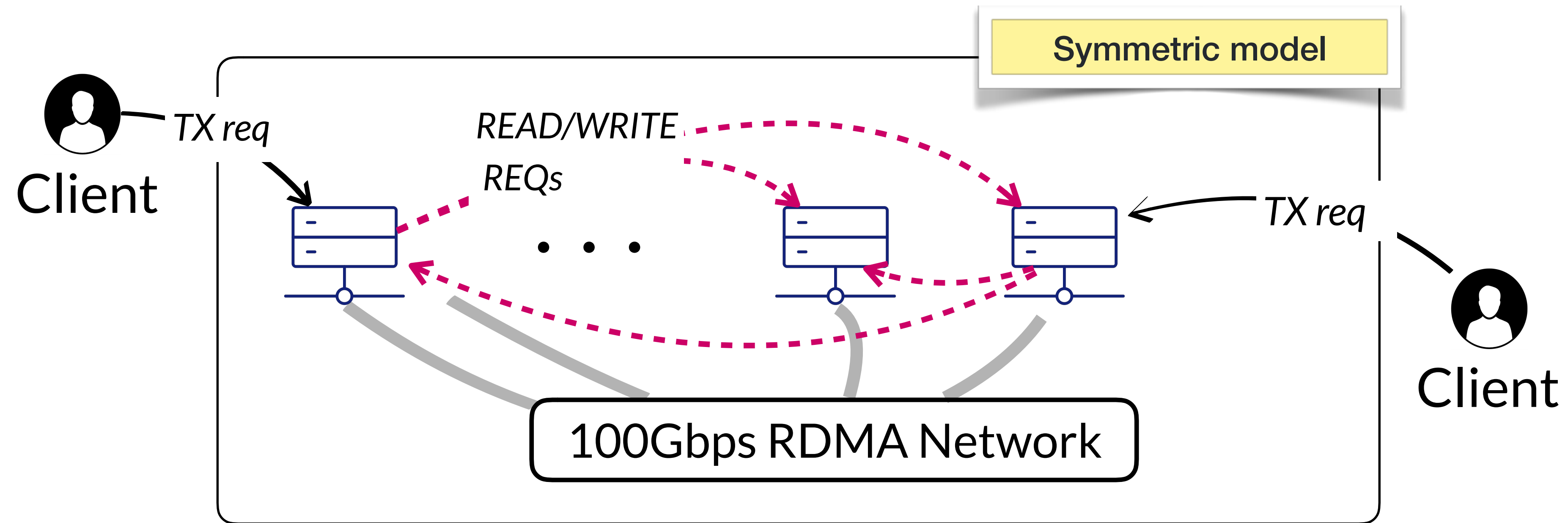
https://github.com/SJTU-IPADS/drtmh

# Outline

☑ RDMA primitive-level analysis

☑ Phase-by-phase analysis for TX

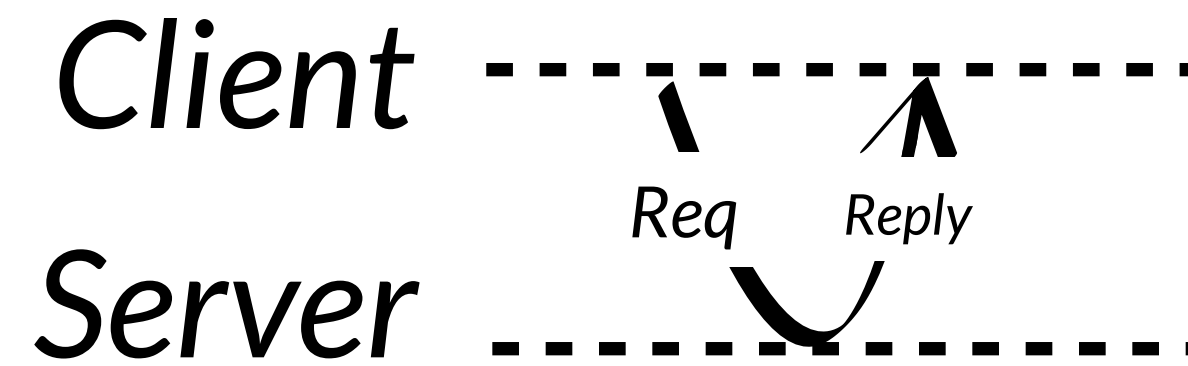☑ DrTM+H: Putting it all together

# System model & evaluation setup



**Evaluation setup**

| | CPU | RNIC | Link layer |
|---|---|---|---|
| **16 x** | **24 cores** | **2 * ConnectX-4 RNIC** | **Infiniband** |

# Primitive analysis

Client  ┈┈┈┈┈┈┈┈┈┈┈
*Req*   *Reply*
Server  ┈┈┈┈┈┈┈┈┈┈┈

## One-sided primitive

➡ Simple implementation
(Native verbs API)

➡ Optimized event loop
(Asyn

## Two-sided (RPC)

➡ FaSST RPC [OSDI'16]

➡ Fastest in our setting

**Throughput (m**

> READ/WRITE is *faster* w *known address*

> CAS is *slower*, but w *sufficient* performance
> *(48M per machine)*

120

90    103

48

80

30

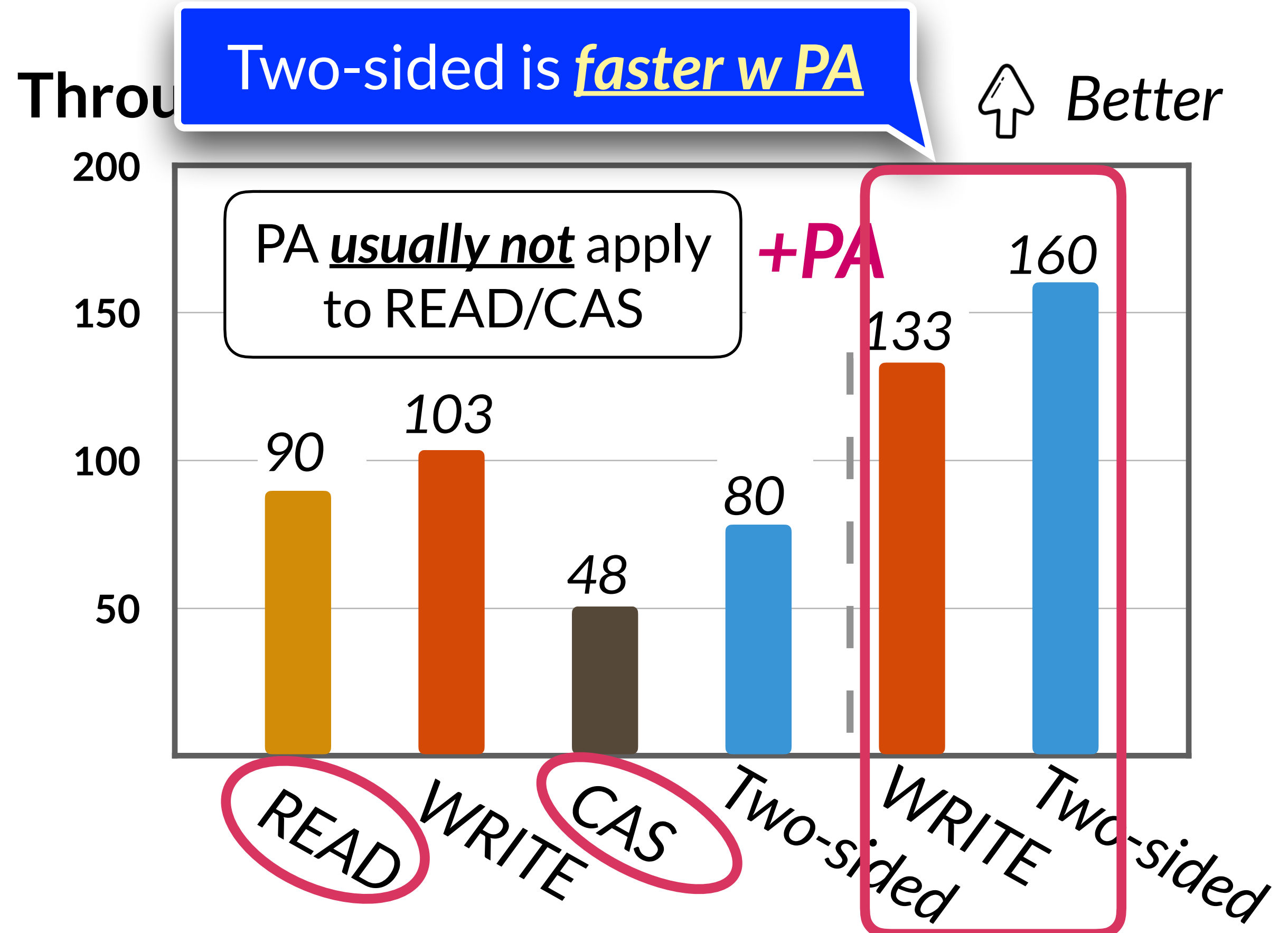READ   WRITE   CAS   Two-sided

# Passive ACK (PA)

**Opt**: when the reply is
**not on the critical path**
of the execution

One-sided primitive

⇨ **Unsignaled** requests

Two-sided primitive

⇨ **Batch** replies (passively)

*Client* - - - - - - - - - - - - - - - - -      - - - - - - - - - - - - - - -
                    *Req*   *Reply*        **+PA**        *Req*       *Reply*
*Server* - - - - - - - - - - - - - - - - -      - - - - - - - - - - - - - - -

**Throu**

Two-sided is *faster w PA*

⇧ *Better*

PA *usually not* apply
to READ/CAS

**+PA**

| | | | | | |
|---|---|---|---|---|---|
| 200 | | | | | |
| 150 | | | | 133 | 160 |
| 100 | 90 | 103 | | 80 | |
| 50 | | | 48 | | |

READ   WRITE   CAS   Two-sided   WRITE   Two-sided

12

# Towards phase-by-phase analysis

Transactional system

⇨ Built atop of our well-tuned execution framework (primitive analysis)

Workloads

⇨ TPC-C/no: new-order  (distributed)

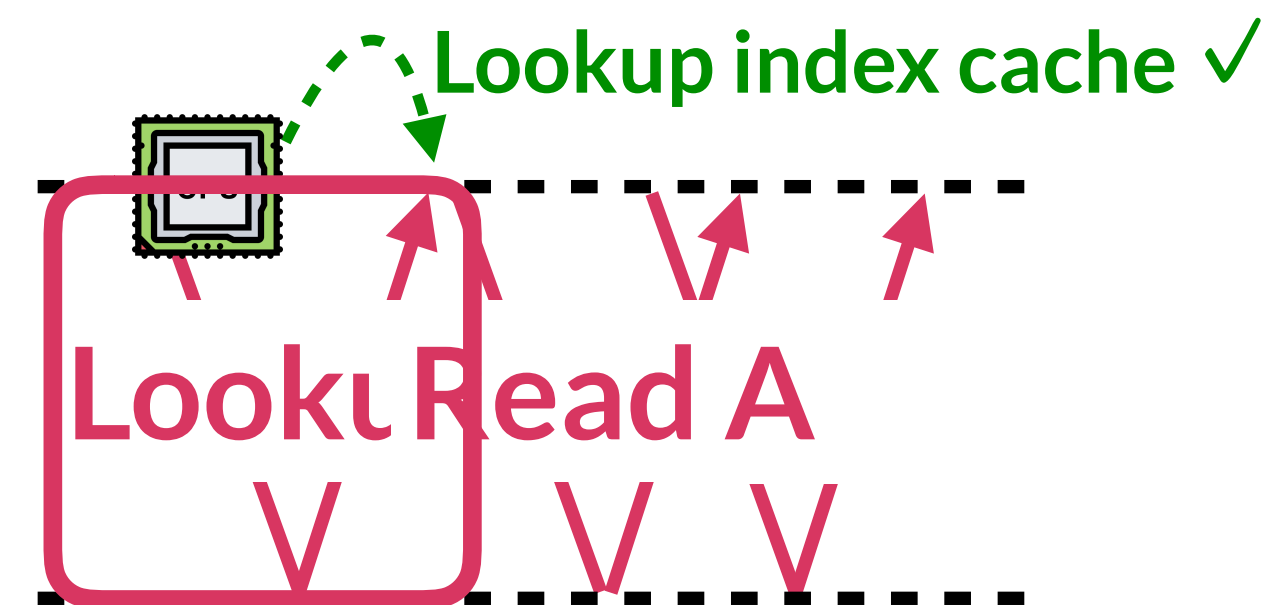⇨ Smallbank

⇨ TPC-E/cp: custom-position

# Execution = READs

**TX{A = A + 1}**

*Coordinator*

*A's store*

**One-sided (I) Cache**

Lookup index cache ✓
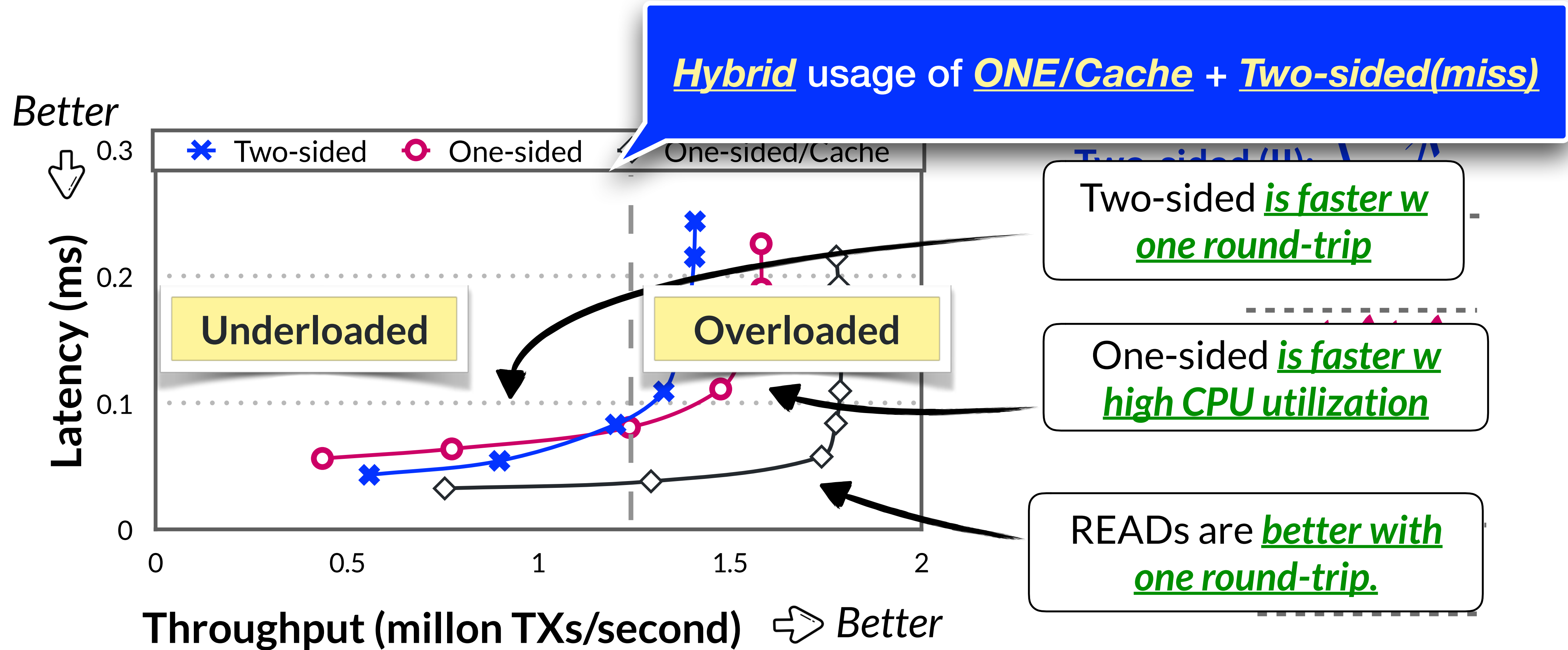
Looku Read A

**Two-sided (II)**

RPC
request

RPC
reply

CPU

Optimization for one-sided primitive

▷ **RDMA friendly store** (e.g. DrTM-KV) -> ~**One-round** lookup

▷ **Index cache**, cache **hot** items address -> **One-round** (lookup + read)

14

# Execution = READs

**Hybrid** usage of **ONE/Cache** + **Two-sided(miss)**

*Better*



Two-sided · One-sided · One-sided/Cache

**Underloaded**

**Overloaded**

Latency (ms)

0.3
0.2
0.1
0

0    0.5    1    1.5    2

**Throughput (millon TXs/second)** ⇨ *Better*

Two-sided *is faster w one round-trip*

One-sided *is faster w high CPU utilization*

READs are *better with one round-trip.*

# Validation = LOCKs + READs

**TX{A = A + 1}**
*Coordinator*

*A's store*

**One-sided (I)**

**Lookup CAS + Read**

*Lock(A)*     *Validate(A)*

**Two-sided (II)**

**RPC request**     **RPC reply**
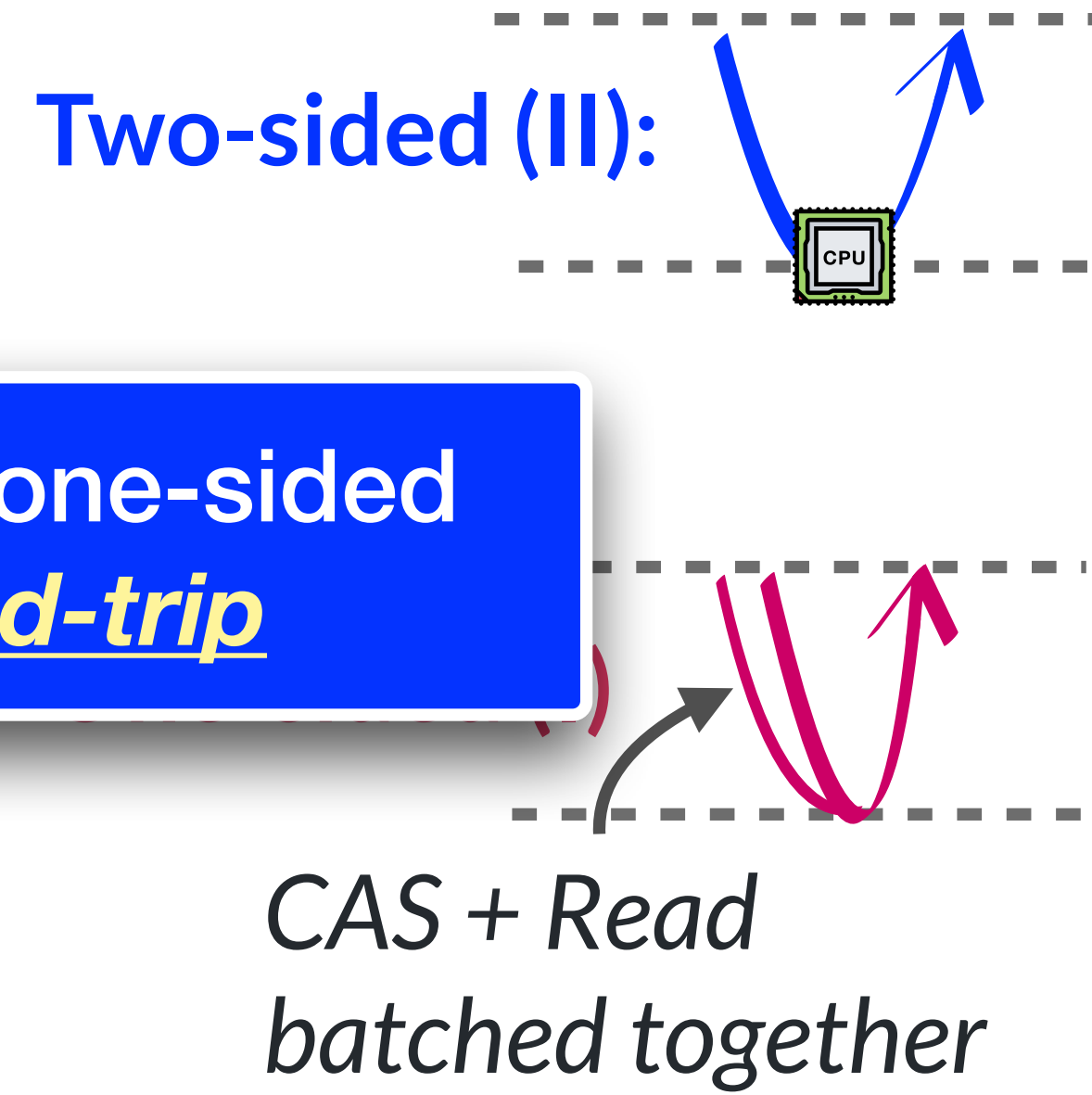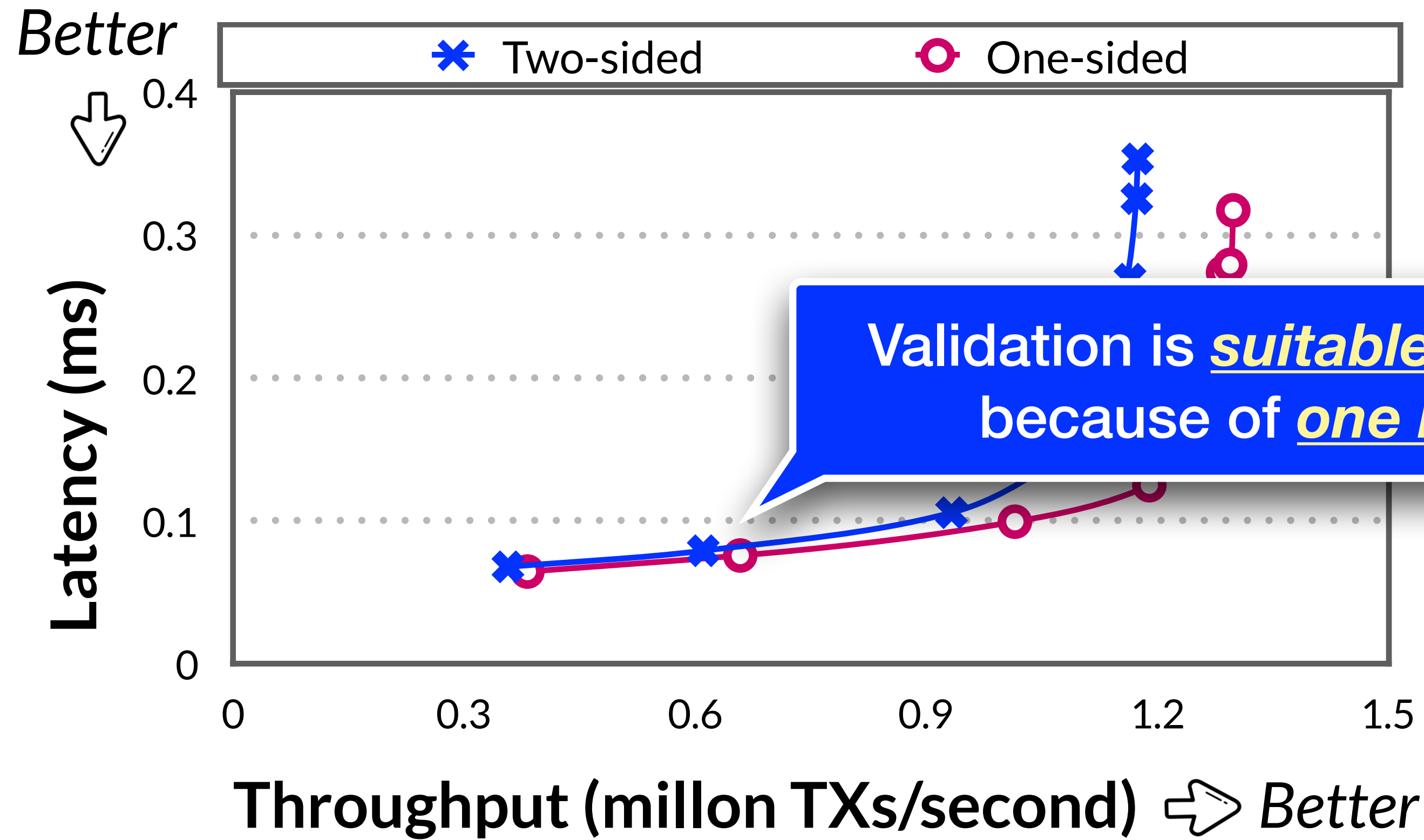
CPU

Optimization for one-sided primitive（for **one round-trip**）

▷ Address **known w the execution phase** -> no need for lookup

▷ Locked value cannot be changed -> doorbell **batch READs w CASs**
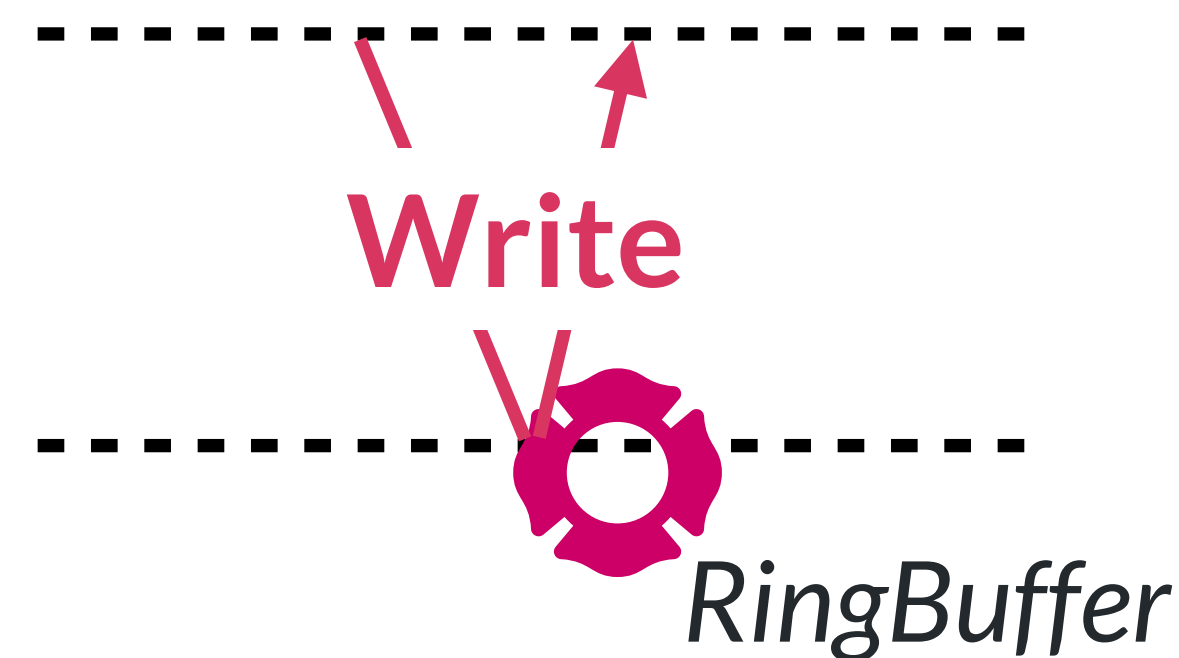
16

# Validation = LOCKs + READs

Exe | Val | Log | Commit



Better

Latency (ms)

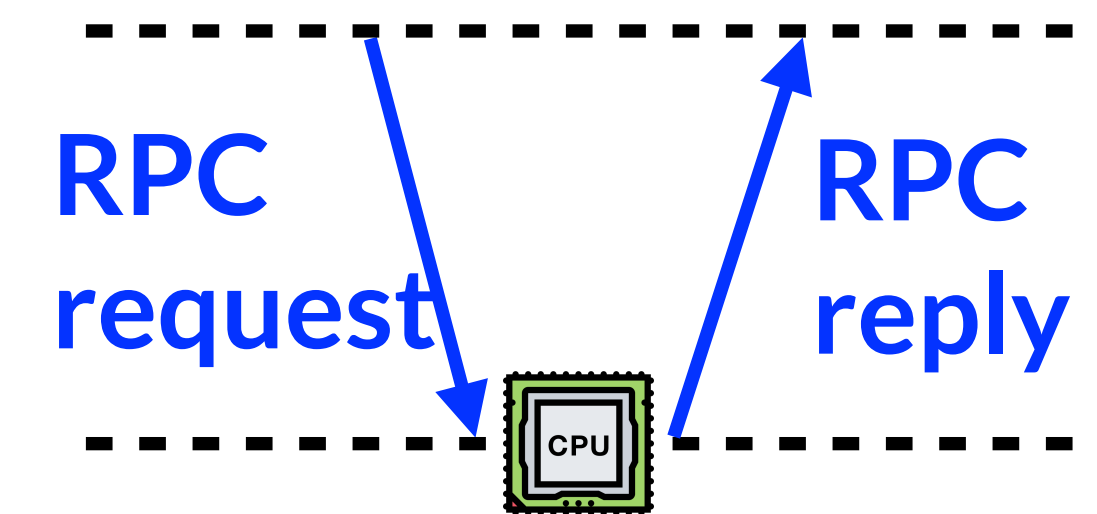Throughput (millon TXs/second) ⇨ Better

- ✶ Two-sided
- ○ One-sided

Two-sided (II):

Validation is *suitable* for one-sided because of *one round-trip*

CAS + Read batched together

# **Logging = WRITEs**

*Exe* | *Val* | ***Log*** | *Commit*

**One-sided (I)**

**Two-sided (II)**

TX{A = A + 1}

*Coordinator*

*A's backup*

**Write**

*RingBuffer*

**RPC request**

**RPC reply**

CPU
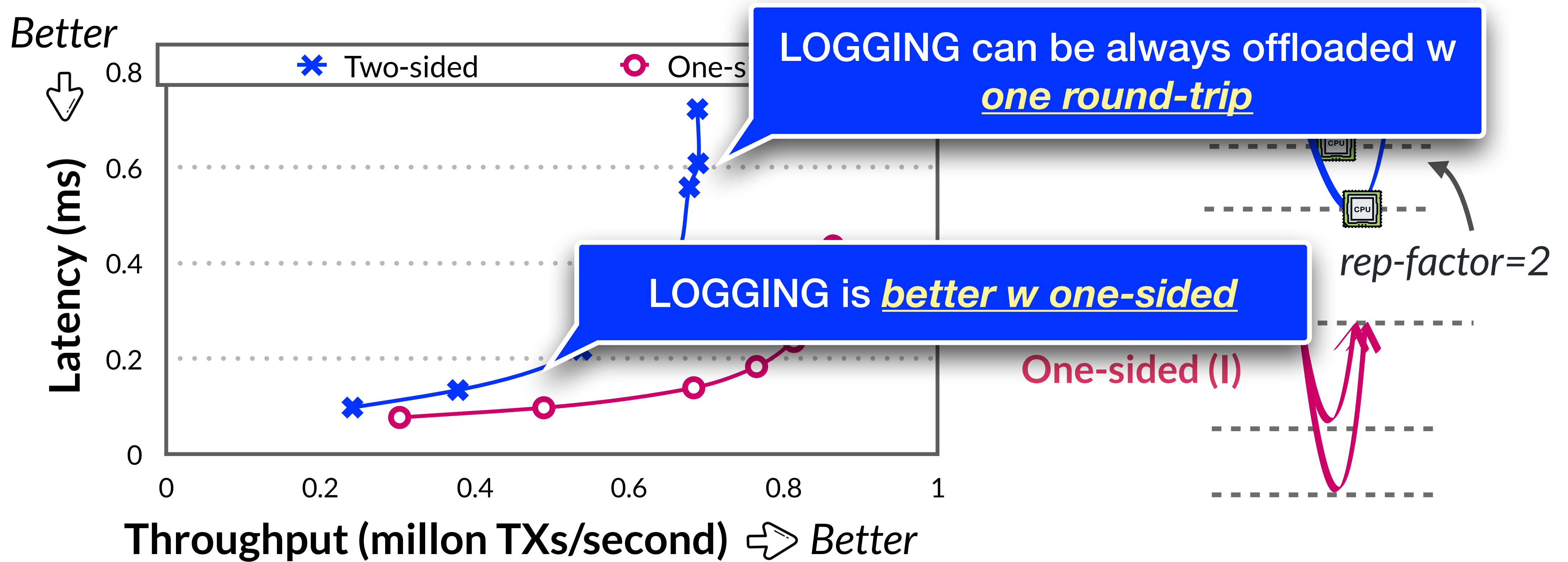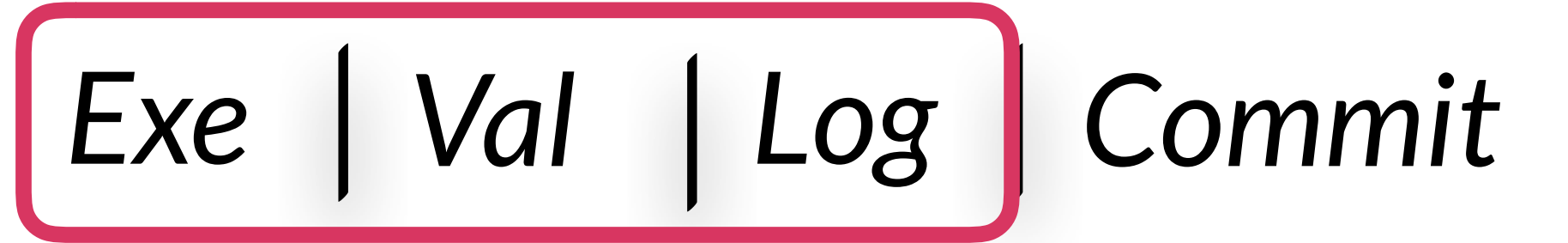
One round-trip for one-sided primitive

⇨ **Ring buffer** based log management [FaRM@NSDI'14]

⇨ **RNIC ack** -> logging succeed (Totally **bypassing** CPU)

18

# Logging = WRITEs

$Exe \mid Val \mid Log$ *Commit*



*Better*

**Latency (ms)**

0.8 — Two-sided — One-s[ided]

0.6

0.4

0.2

0

0    0.2    0.4    0.6    0.8    1

**Throughput (millon TXs/second)** *Better*

**LOGGING can be always offloaded w _one round-trip_**

**LOGGING is _better w one-sided_**

*rep-factor=2*

One-sided (I)

# Commit = WRITEs + UNLOCKs

*Exe* | *Val* | *Log* | ***Commit***

**One-sided (I)**

**Two-sided (II)**

TX{A = A + 1}

*Coordinator*

Lookup Write A

RPC request

RPC reply

CPU
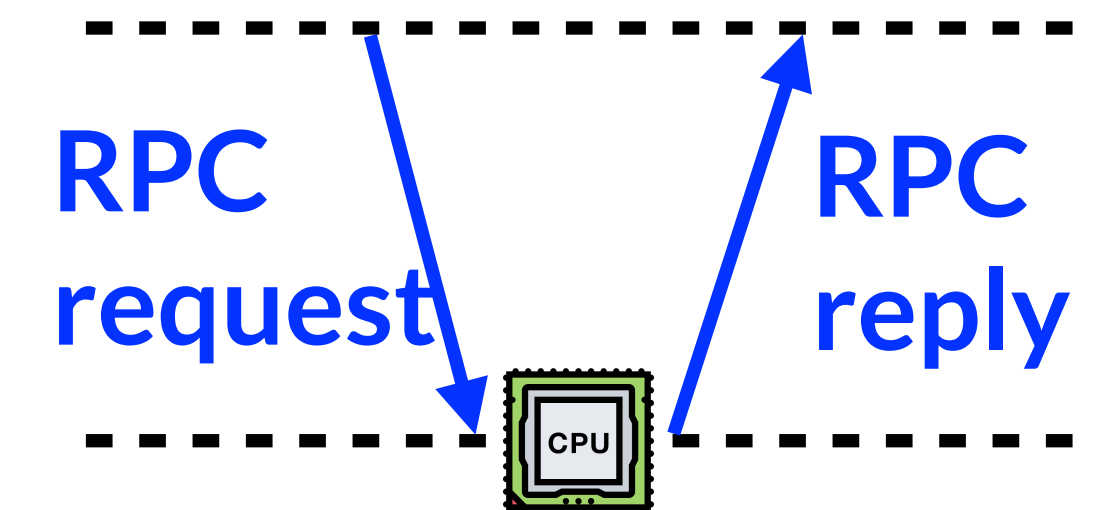
*A's store*

*Unlocks implemented as WRITEs*

One round-trip for one-sided primitive

⇨ Address **known w the execution phase** -> no need for lookup
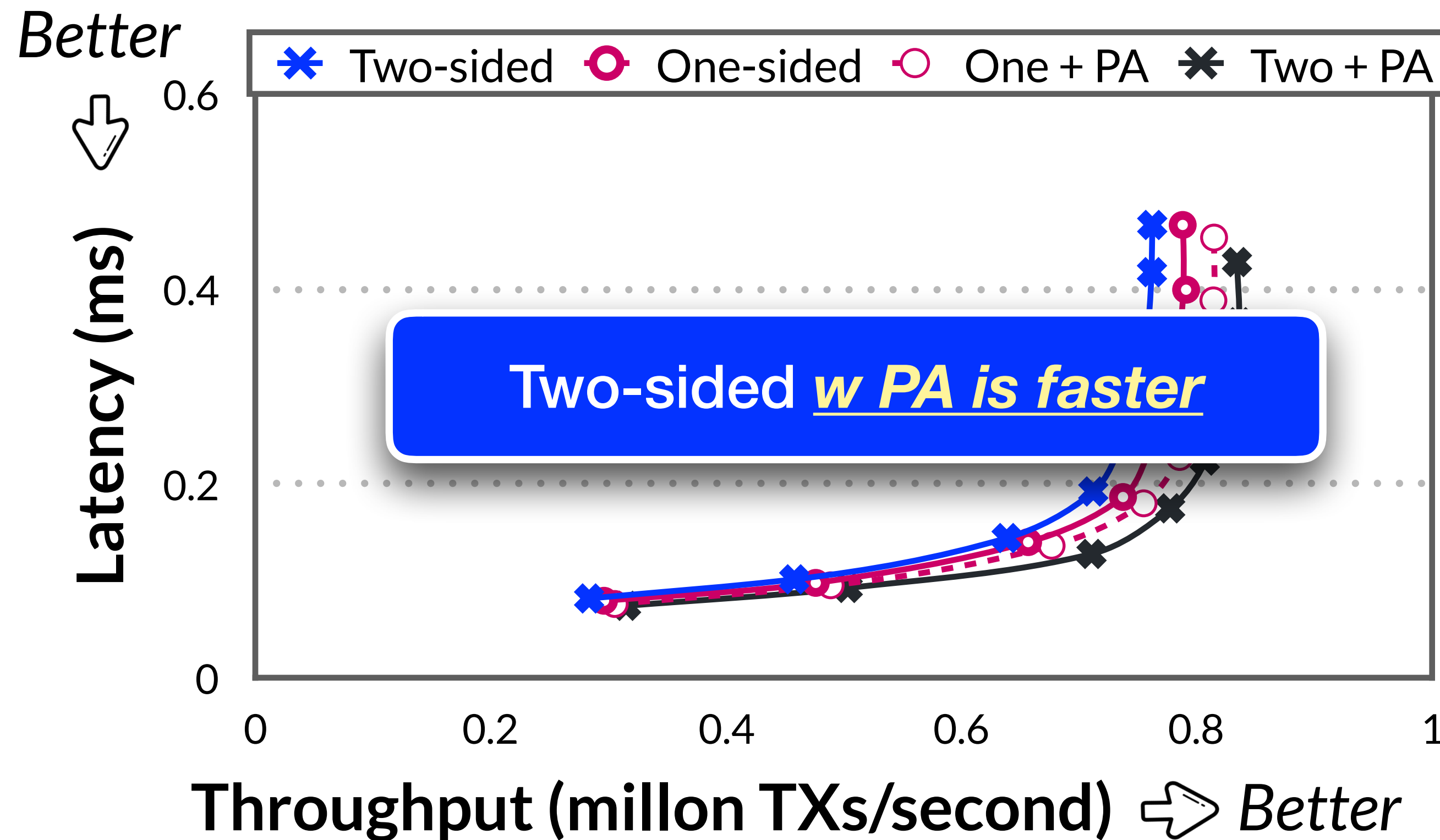
Adding passive ACK to **both primitives**

⇨ **Log succeed** indicates TX's **commit**
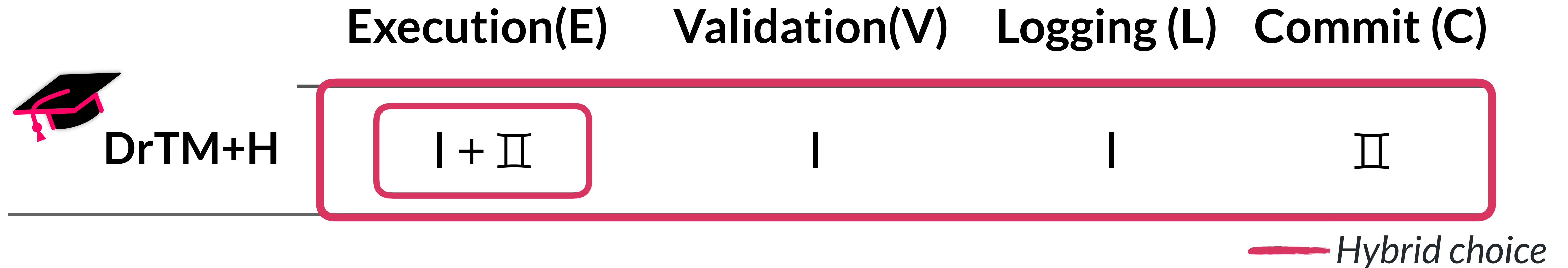
# Commit = WRITEs + UNLOCKs

*Exe* | *Val*   *Log*   *Commit*



**Better** ⬇

Latency (ms)

0.6
0.4
0.2
0

* Two-sided   ◯ One-sided   ◯ One + PA   ✳ Two + PA

**Two-sided *w PA is faster***

0    0.2    0.4    0.6    0.8    1

**Throughput (millon TXs/second)** ⇨ *Better*

⇨ Two-sided w PA has *higher peak throughput*

⇨ Commit RPC *costs is small*

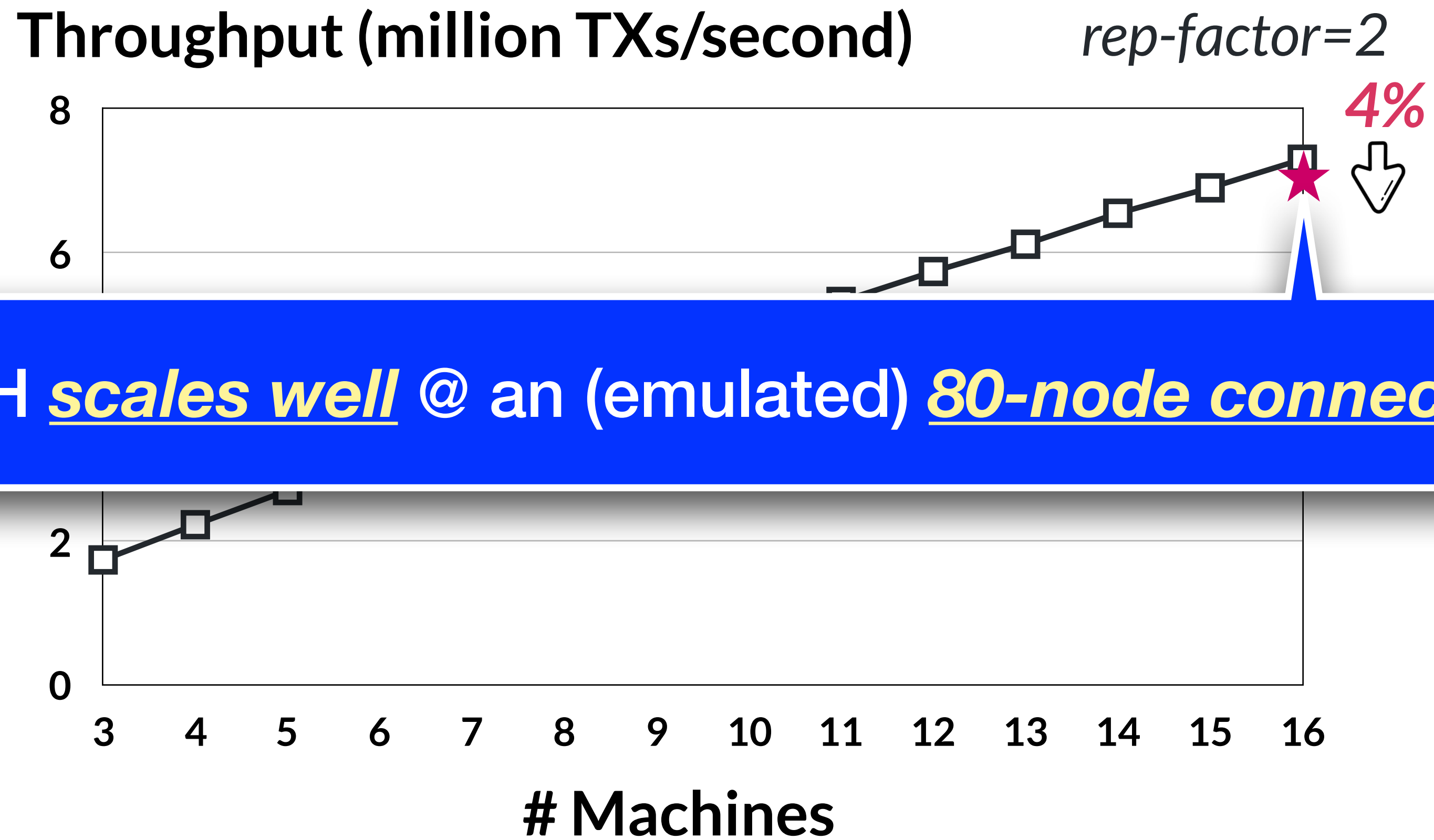⇨ Two-sided saves CPU *at sender*

# DrTM+H: Hybrid is better !

Hybrid system supports **serializability** & **high availability**

| | Execution(E) | Validation(V) | Logging (L) | Commit (C) |
|---|---|---|---|---|
| DrTM+H | Ⅰ + Ⅱ | Ⅰ | Ⅰ | Ⅱ |

—— *Hybrid choice*

## Specific optimizations

▷ **Passive ACK** to the commit phase ( & log cleaning message)

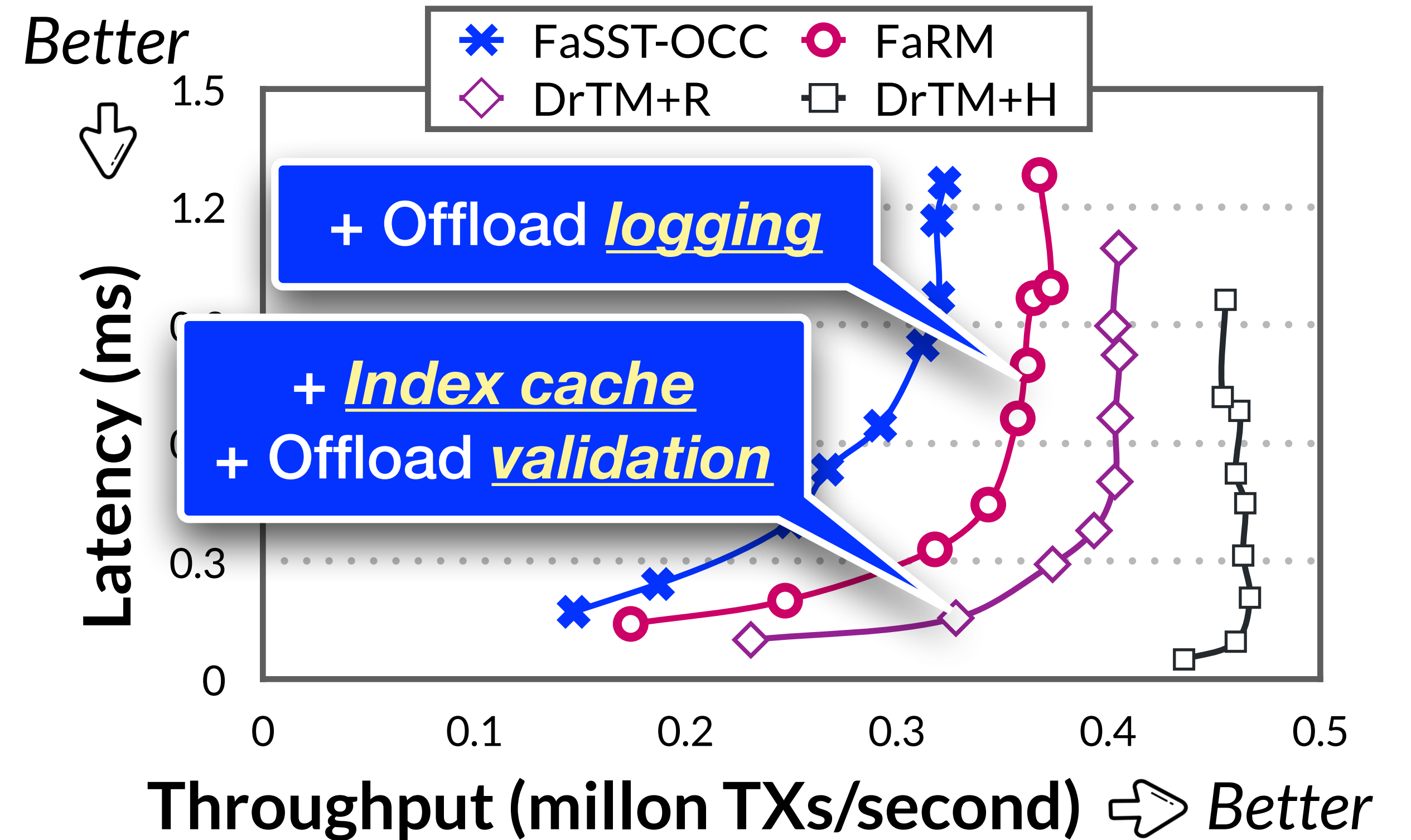▷ Speculative execution to send **outstanding requests (OR)** from one TX

# Performance & scalability on TPC-C/no

**Throughput (million TXs/second)**       *rep-factor=2*



*4%*

DrTM+H *scales well* @ an (emulated) *80-node connections*

**# Machines**

# End-to-end comparison against prior designs

In the **same platform**, the **same protocol**, but w **different choices**

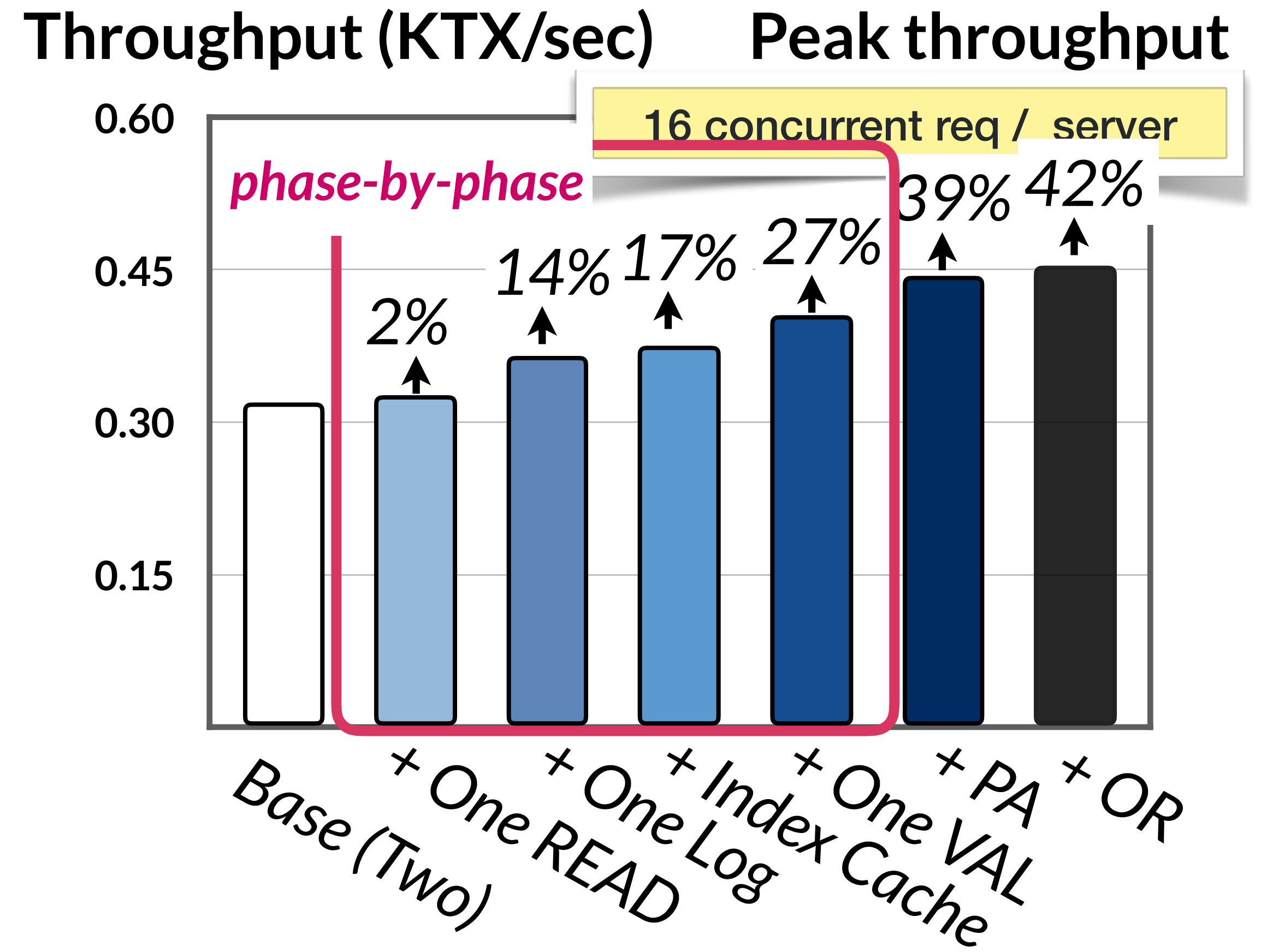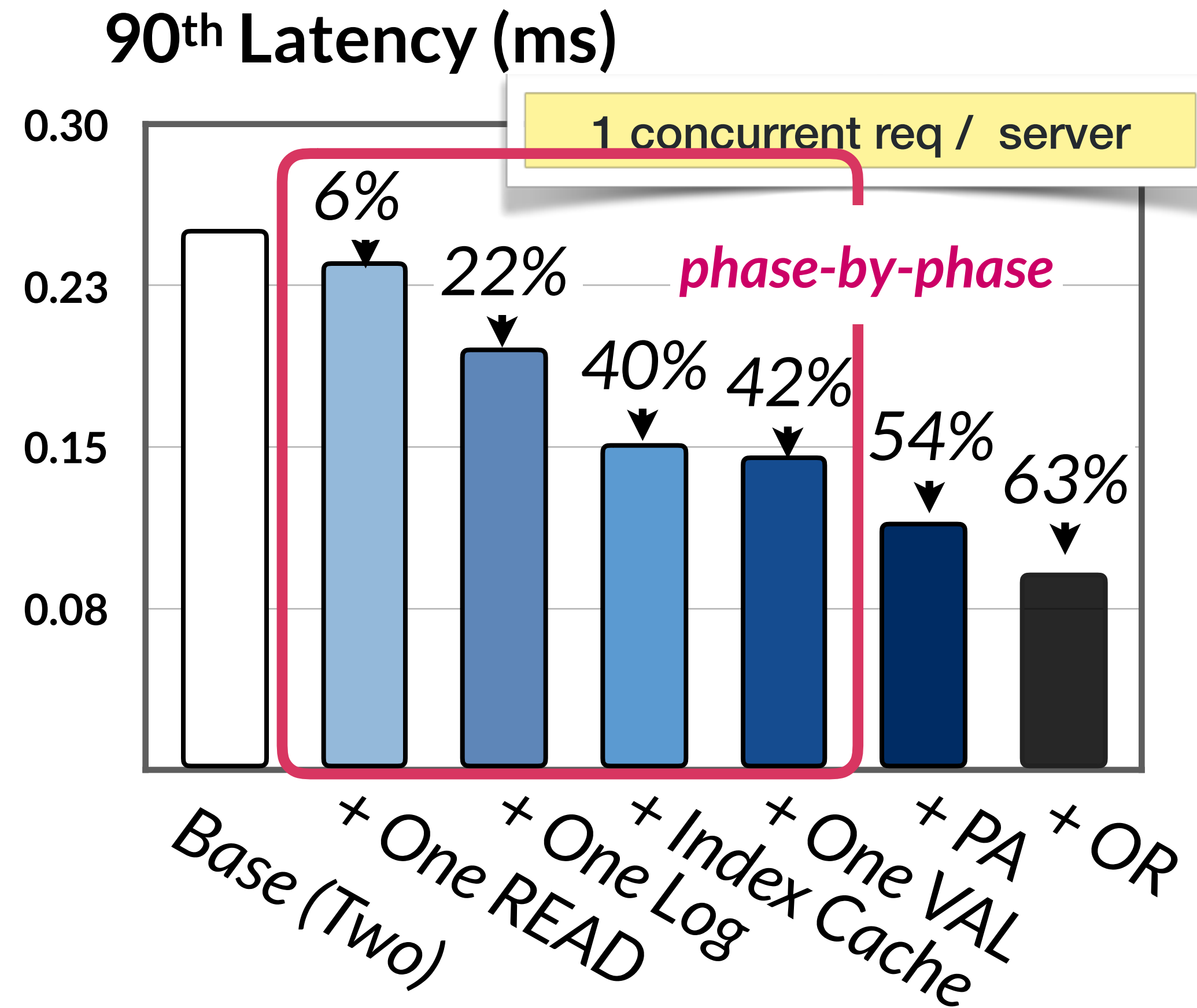| | E | V | L | C |
|---|---|---|---|---|
| **FaSST-OCC**[1] | ⏗ | ⏗ | ⏗ | ⏗ |
| **DrTM+R** | ❙[w cache] | ❙ | ❙ | ❙ |
| **FaRM** | ❙[w/o cache] | ❙+⏗ | ❙ | ⏗ |
| **DrTM+H** | ❙ + ⏗ | ❙ | ❙ | ⏗ |

*Better* ⬇

**Latency (ms)**

Legend: ✖ FaSST-OCC   ⬤ FaRM   ◇ DrTM+R   ☐ DrTM+H

+ Offload *logging*

+ *Index cache*
+ Offload *validation*

Throughput axis: 0, 0.1, 0.2, 0.3, 0.4, 0.5
Latency axis: 0, 0.3, 1.2, 1.5

**Throughput (millon TXs/second)** ⇨ *Better*

[1] FaSST uses a simplified OCC protocol compared to FaRM & DrTM+R.
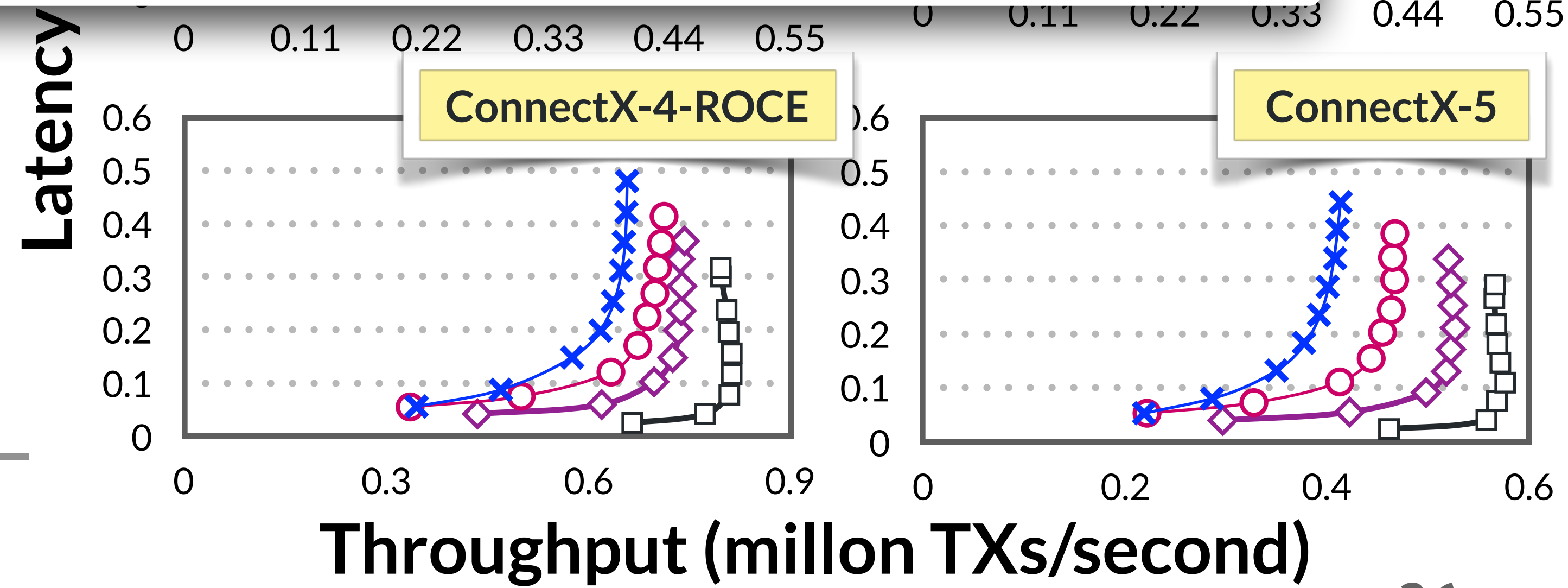
24

# Where do the performance gains come from?



**90th Latency (ms)**

1 concurrent req / server

phase-by-phase

6% 22% 40% 42% 54% 63%

Base (Two) + One READ + One Log + Index Cache + One VAL + PA + OR

0.30 0.23 0.15 0.08

**Throughput (KTX/sec)    Peak throughput**

16 concurrent req / server

phase-by-phase

2% 14% 17% 27% 39% 42%

Base (Two) + One READ + One Log + Index Cache + One VAL + PA + OR

0.60 0.45 0.30 0.15

# Not a hard conclusion !

May depends on **RNIC's characteristic** & **network setting**

| Legend | | | |
|---|---|---|---|
| ✕ FaSST-OCC | | ○ FaRM | |
| ◇ DrTM+R | | ⊟ DrTM+H | |

**RN...**

> The results start from the *primitive level analysis.*

| | | | | |
|---|---|---|---|---|
| **CX3[1]** | ⊞ | ⊞ | ⊞ | ⊞ |
| **CX4[2]** | ❘ + ⊞ | ❘ | ❘ | ⊞ |
| **CX4-ROCE[1]** | ❘ + ⊞ | ❘ | ❘ | ⊞ |
| **CX5[1][3]** | ❘ + ⊞ | ❘ | ❘ | ⊞ |

**ConnectX-3**

**ConnectX-4**

**ConnectX-4-ROCE**

**ConnectX-5**

**Latency**

**Throughput (millon TXs/second)**

[1]1-way replication used due to cluster limitation
[2] Main results in this talk
[3]1-RNIC per machine, others uses 2

26

# Evaluation summary

**Offloading** w one-sided **improves the performance**

> ⇨ Especially **w/o adding more round-trips**

> ⇨ Less affected by **CPU load** at the server

One-sided primitive has **good scalability** on modern RNIC

> ⇨ Especially when RNIC is **not the bottleneck** of the application

> ⇨ Although one-sided primitive is restricted by hardware limitation

# More: check our paper!
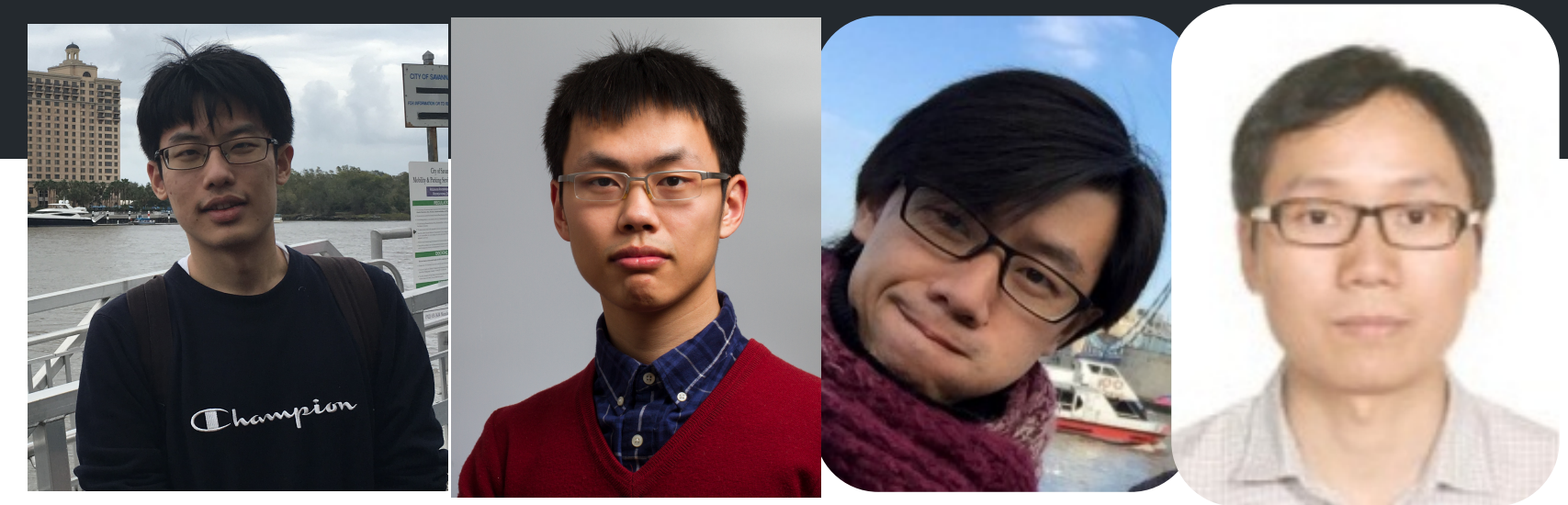
☑ Optimized execution framework

☑ Results of large scale

   Modern RNIC has good scalability for one-sided primitive

☑ Read-only Transactions

   A hybrid scheme also wins

☑ TPC-E, Smallbank
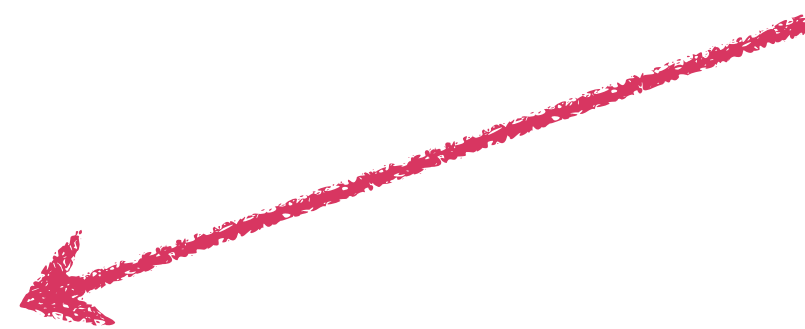
# Conclusions

The **first systematic** study on

⇨ How to use RDMA for OCC TXs

**Thanks & QA**

**No single** primitive is **better**!

⇨ Depends on workload pattern & primitive analysis

Execution framework & DrTM+H are available @

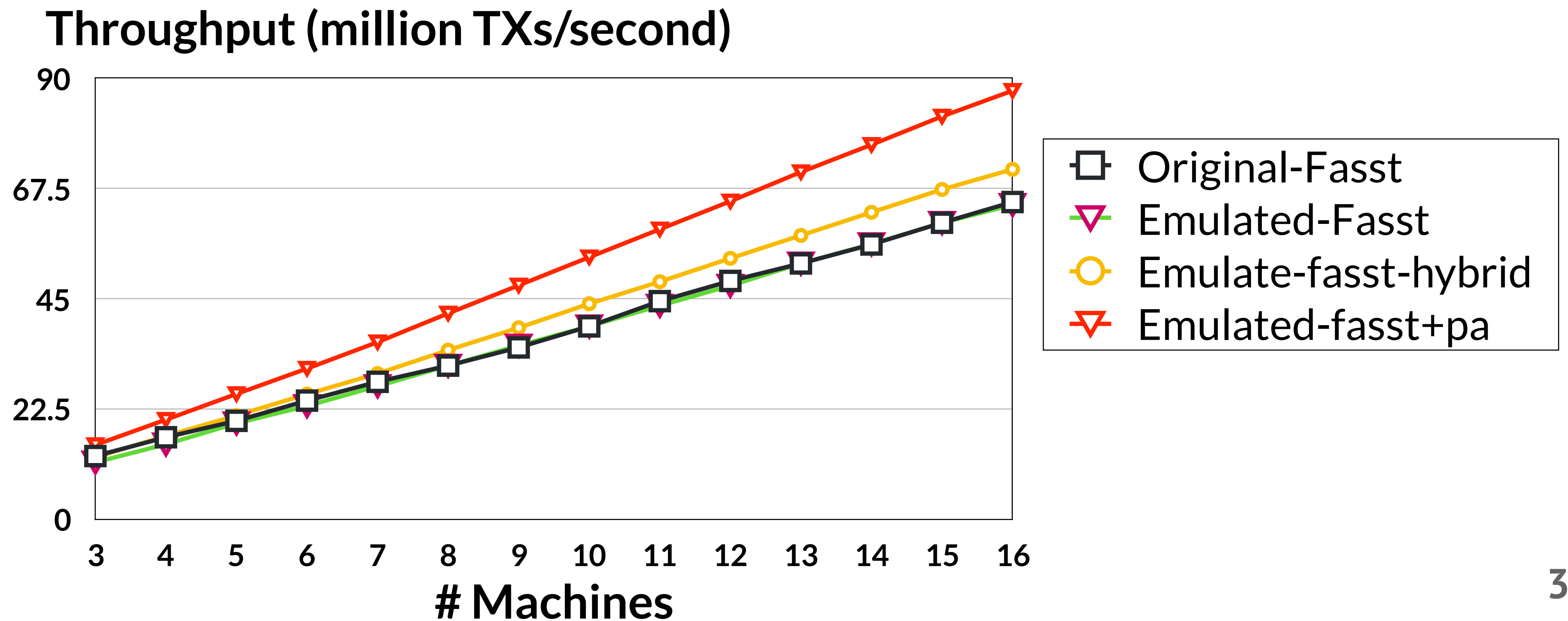https://github.com/SJTU-IPADS/drtmh

# Backups

# Improved overall systems

FaSST's simplified OCC protocol

Adding hybrid-schema for logging

**Throughput (million TXs/second)**

# Smallbank workloads

| | E | V | L | C |
|---|---|---|---|---|
| **CX3[1]** | ⟂ | ⟂ | ⟂ | ⟂ |
| **CX4** | \| + ⟂ | \| | \| | ⟂ |
| **CX4-ROCE[1]** | \| + ⟂ | \| | \| | ⟂ |
| **CX5[1][2]** | \| + ⟂ | \| | \| | ⟂ |

Legend:
- ✕ FaSST-OCC
- ◇ DrTM+R
- ○ FaRM
- ☐ DrTM+H



ConnectX-3

ConnectX-4

ConnectX-4-ROCE

ConnectX-5

Latency (ms)

Throughput (millon TXs/second)

[1]1-way replication used due to cluster limitation

[2]1-RNIC per machine, others uses 2

# RDMA based execution framework

Applied & based RDMA optimizations
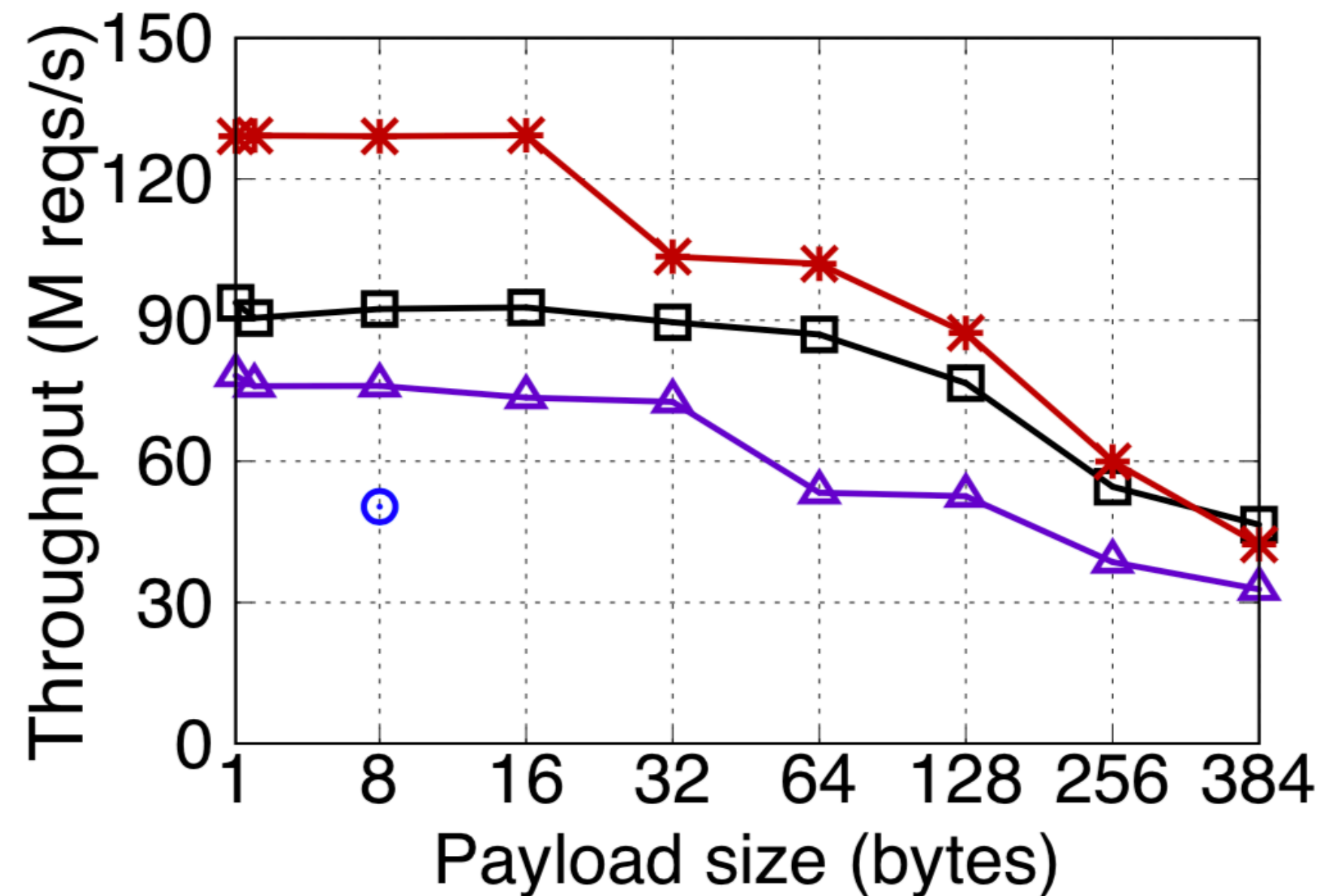
⇨ FaRM [NSDI'14,SOSP'15]

⇨ Herd [NSDI'14]

⇨ RDMA guideline [ATC'16]

⇨ FaSST [OSDI'16]
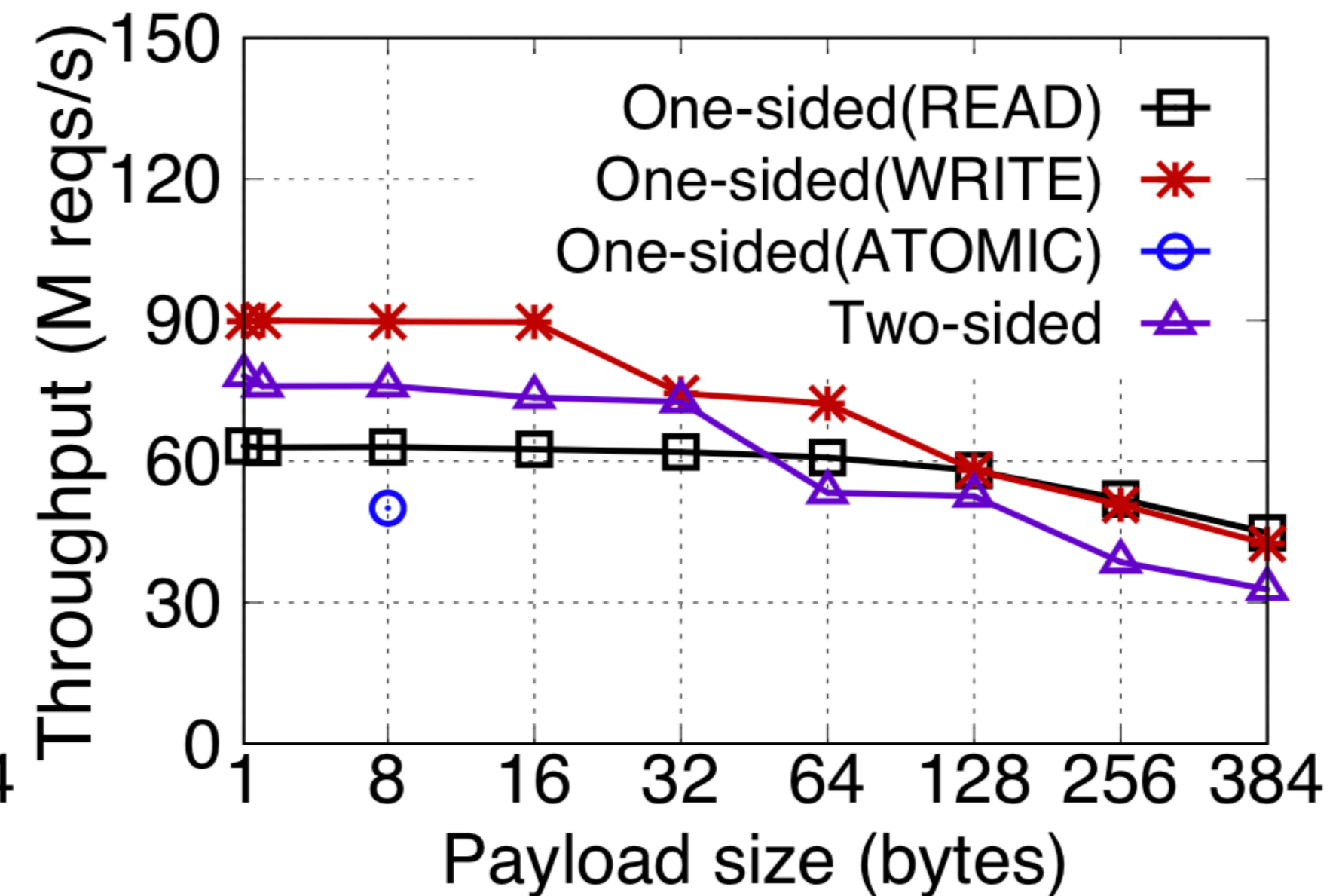
Others

⇨ LITE [SOSP'2017] -> Further improve one-sided's scalability

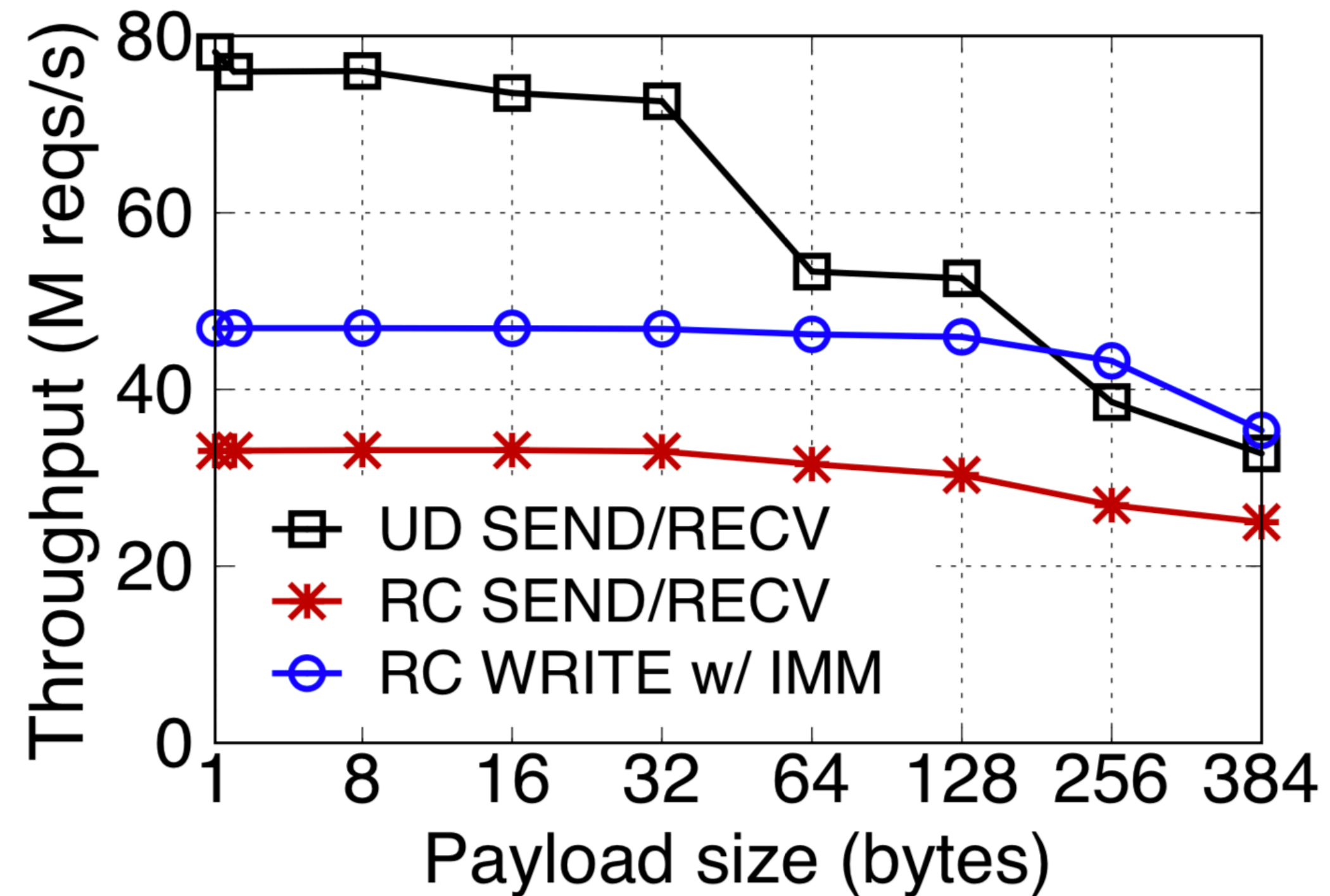# Results using large connections



**16 node**

**Emulate 80-node**

# Comparison of two-sided implementations

FaSST RPC uses UD SEND/RECV

# RDMA enabled application

Load balance framework


Distributed TXs


Graph processing systems


Distributed file system