# Deploying User-space TCP at Cloud Scale with LUNA

Lingjun Zhu, Yifan Shen, Erci Xu , Bo Shi, Ting Fu, *Shu Ma*, Shuguang Chen, Zhongyu Wang, Haonan Wu, Xingyu Liao, Zhendan Yang, Zhongqing Chen, Wei Lin, Yijun Hou, Rong Liu, Chao Shi, Jiaji Zhu, and Jiesheng Wu
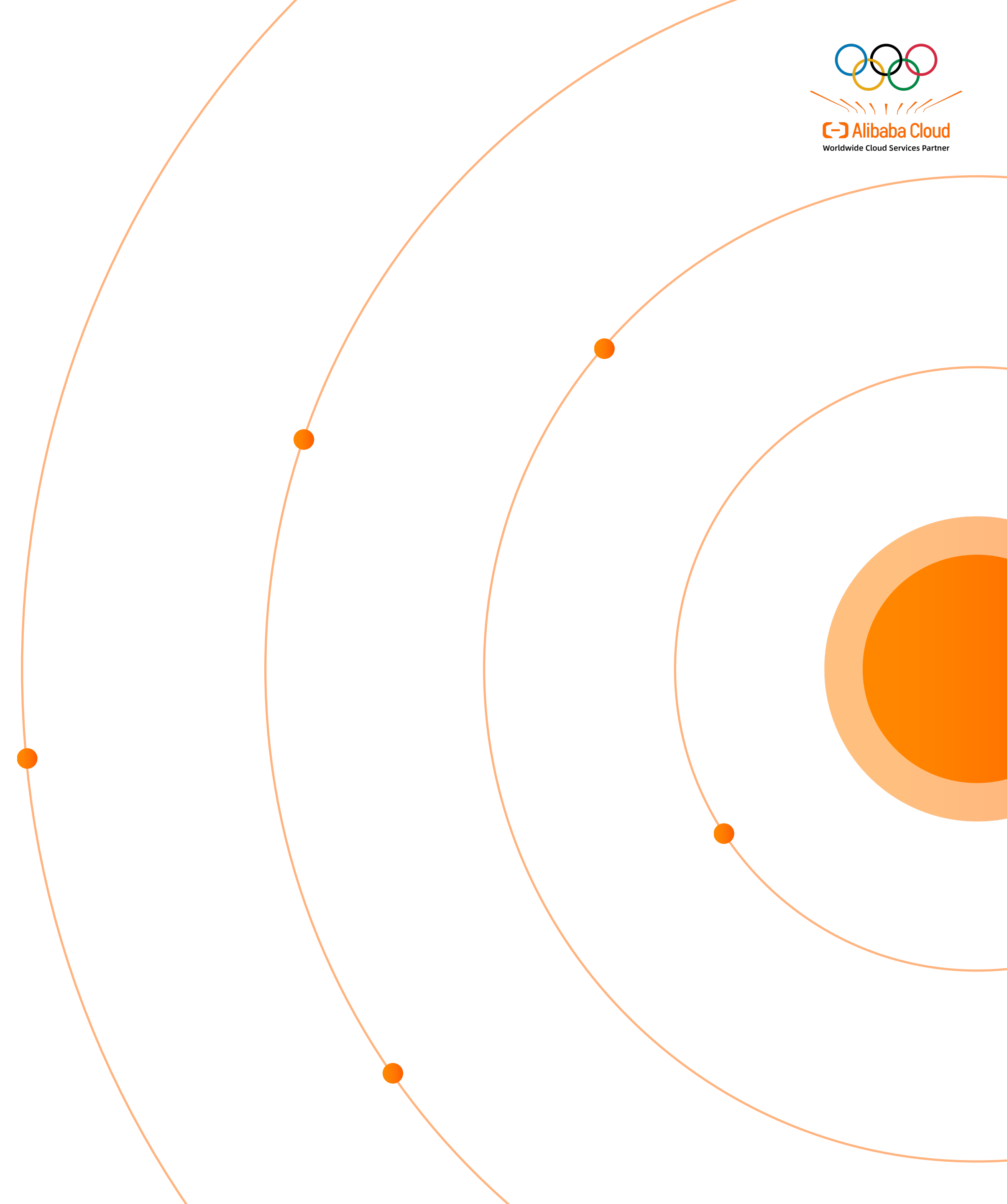
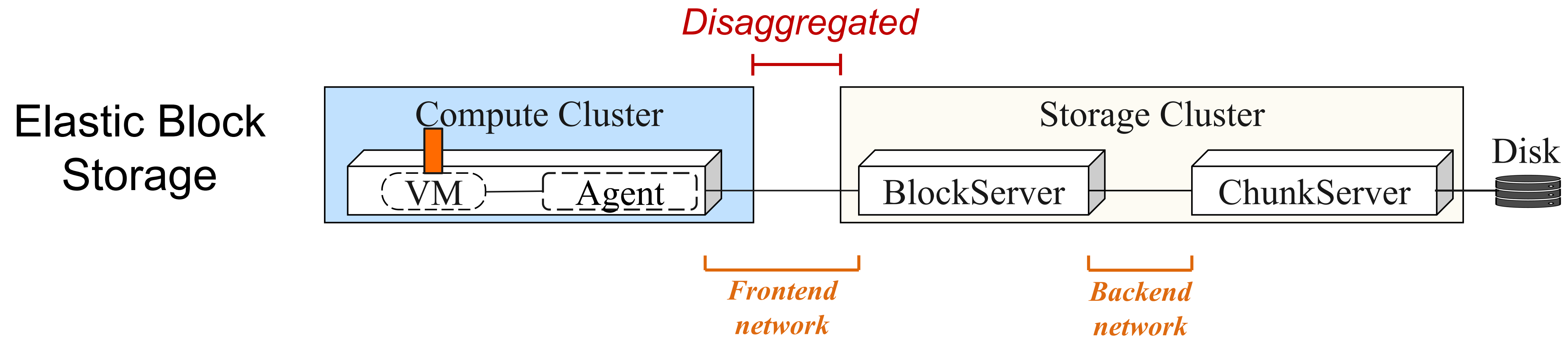**Alibaba Group**

# Background & Motivation
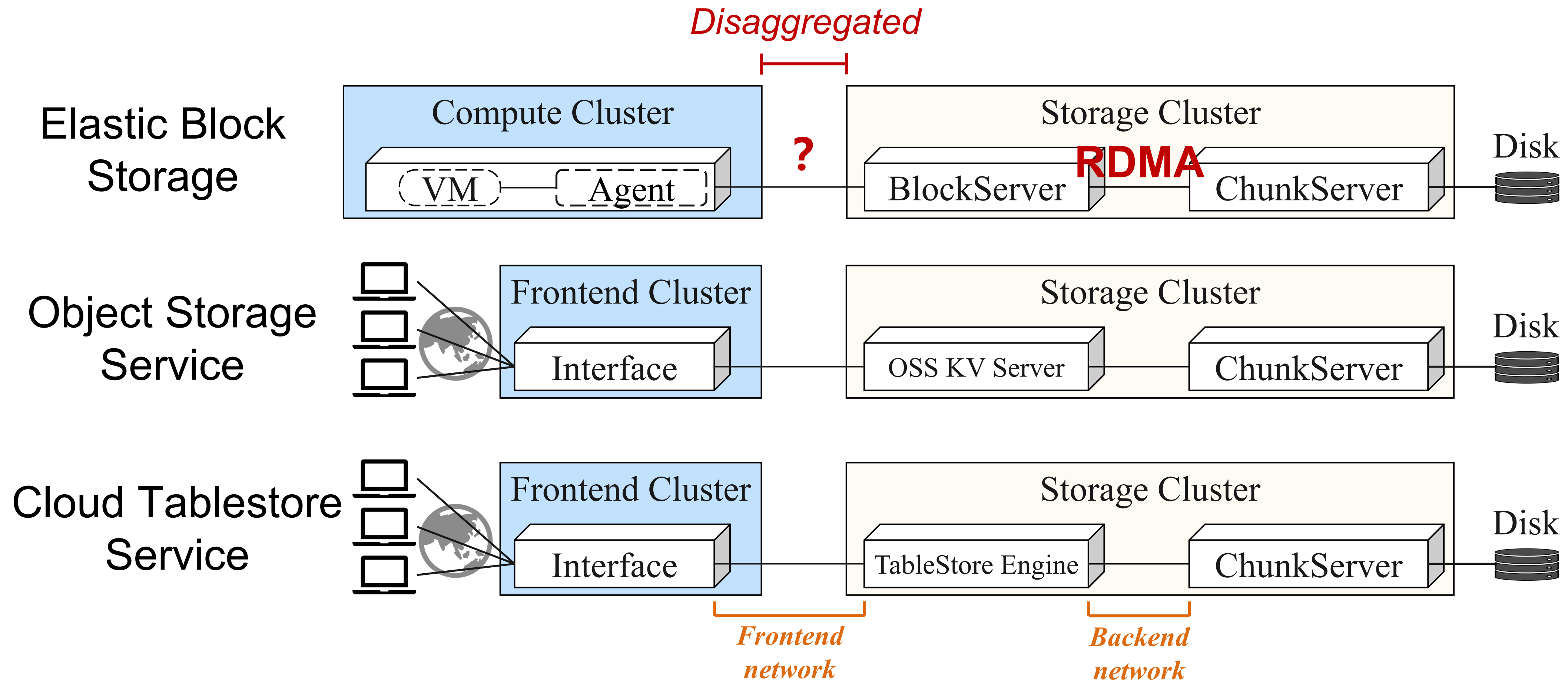
Design

Evaluation

Discussion & Lessons

Conclusion

# Alibaba Cloud Storage Network

# Alibaba Cloud Storage Network



*Disaggregated*

**Elastic Block Storage**

Compute Cluster | ? | Storage Cluster
VM — Agent | | BlockServer — **RDMA** — ChunkServer — Disk

**Object Storage Service**

Frontend Cluster | | Storage Cluster
Interface | | OSS KV Server — ChunkServer — Disk

**Cloud Tablestore Service**

Frontend Cluster | | Storage Cluster
Interface | | TableStore Engine — ChunkServer — Disk

*Frontend network*     *Backend network*
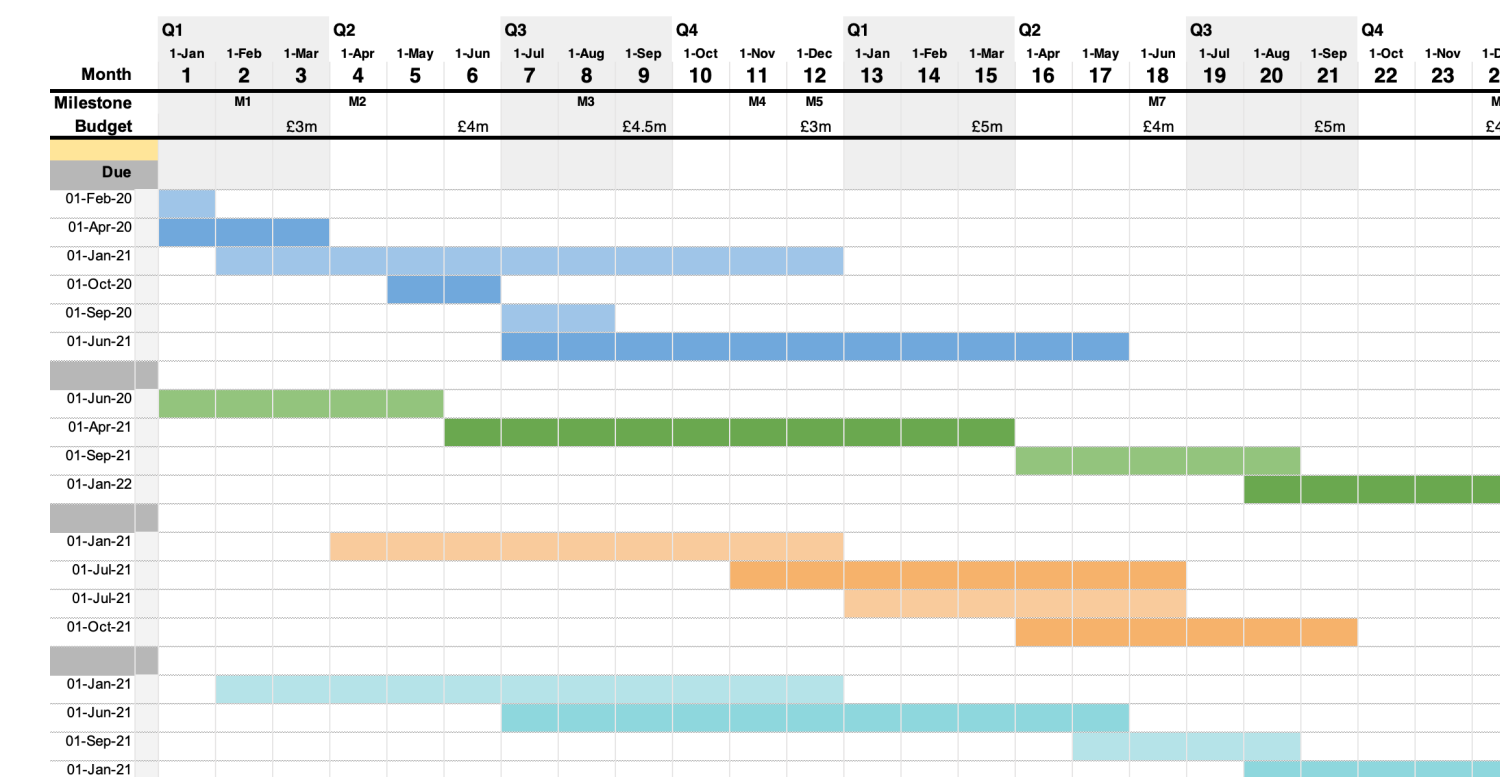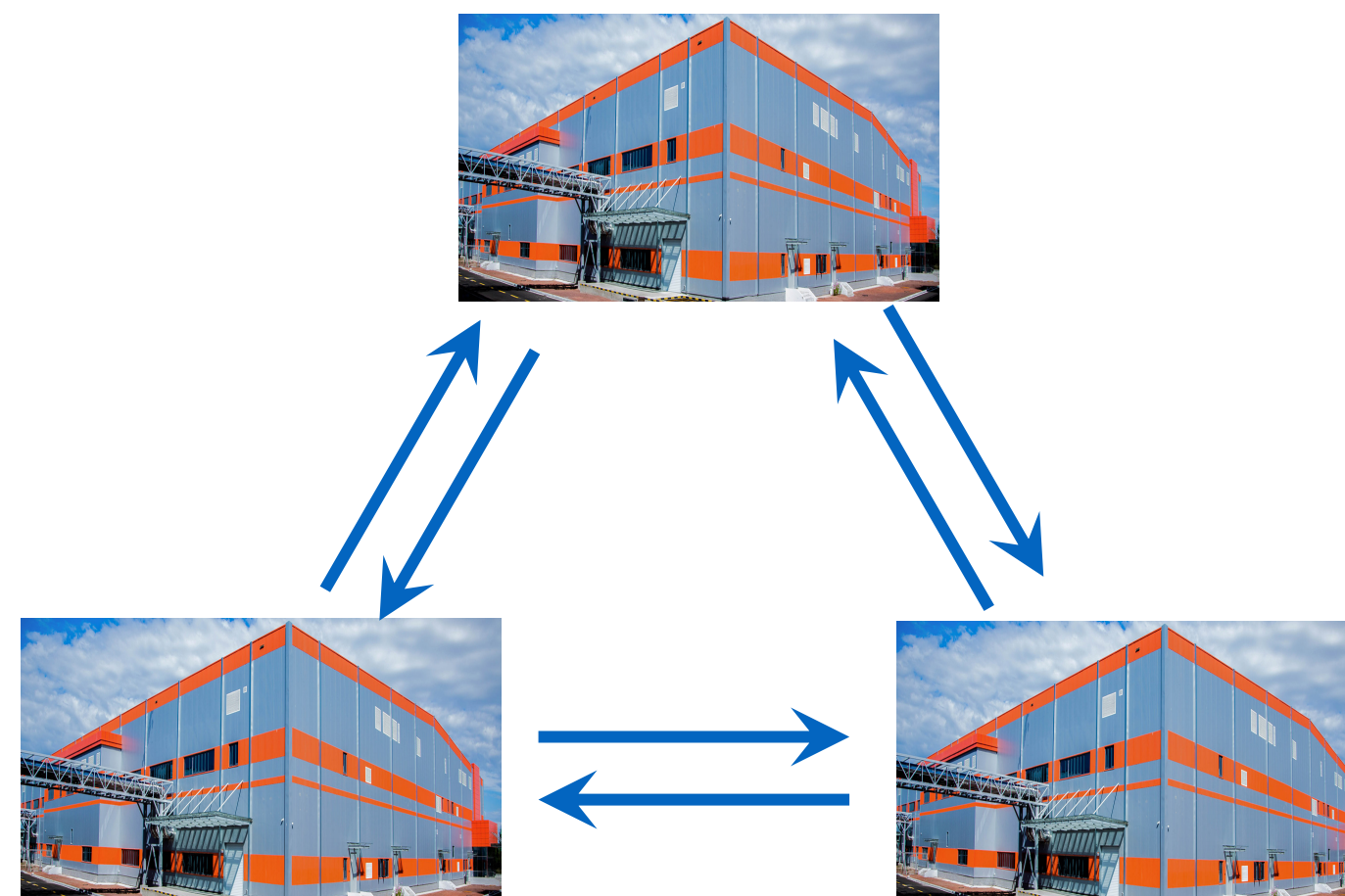
[1] NSDI 2021, When cloud storage meets RDMA
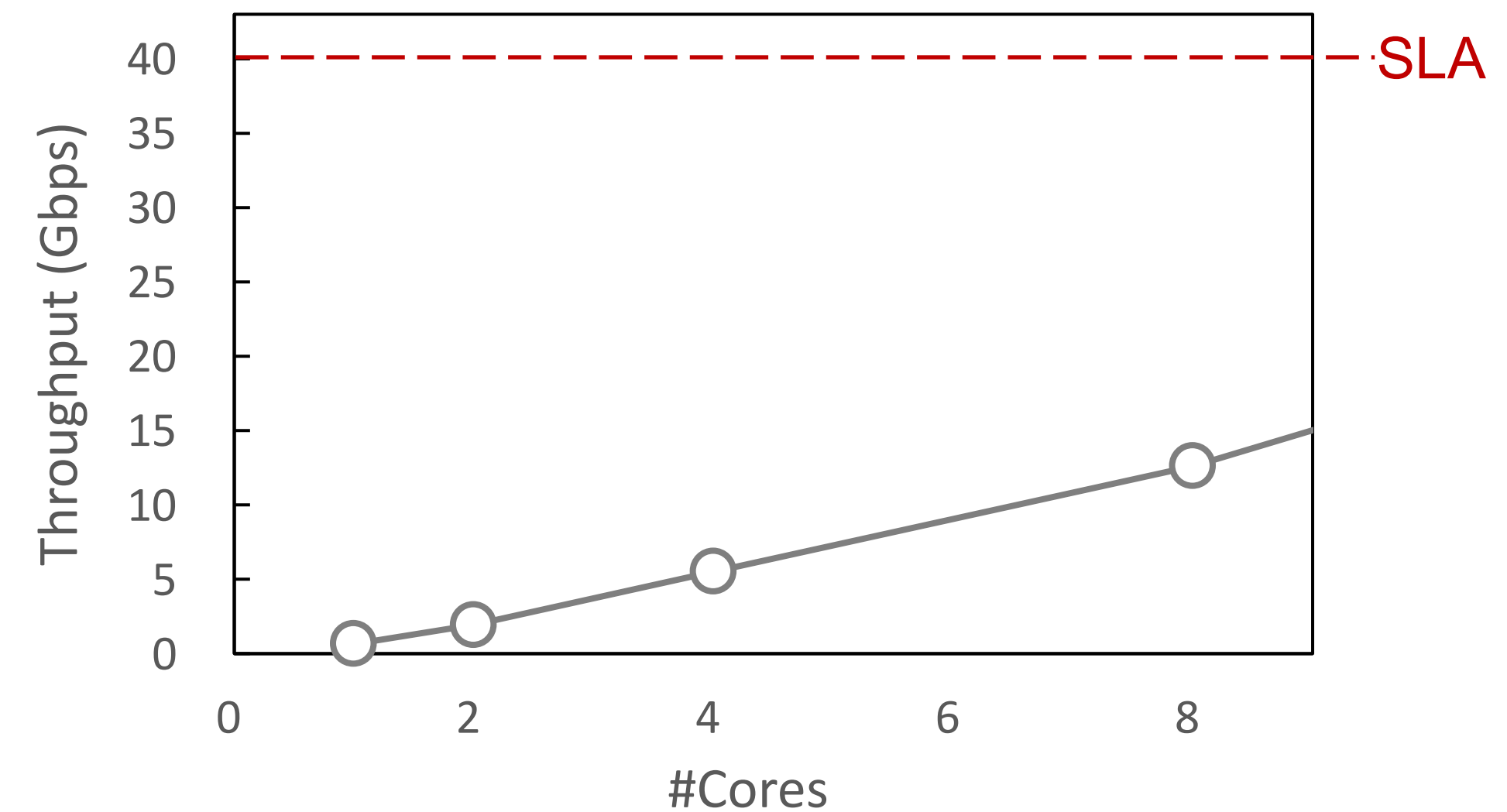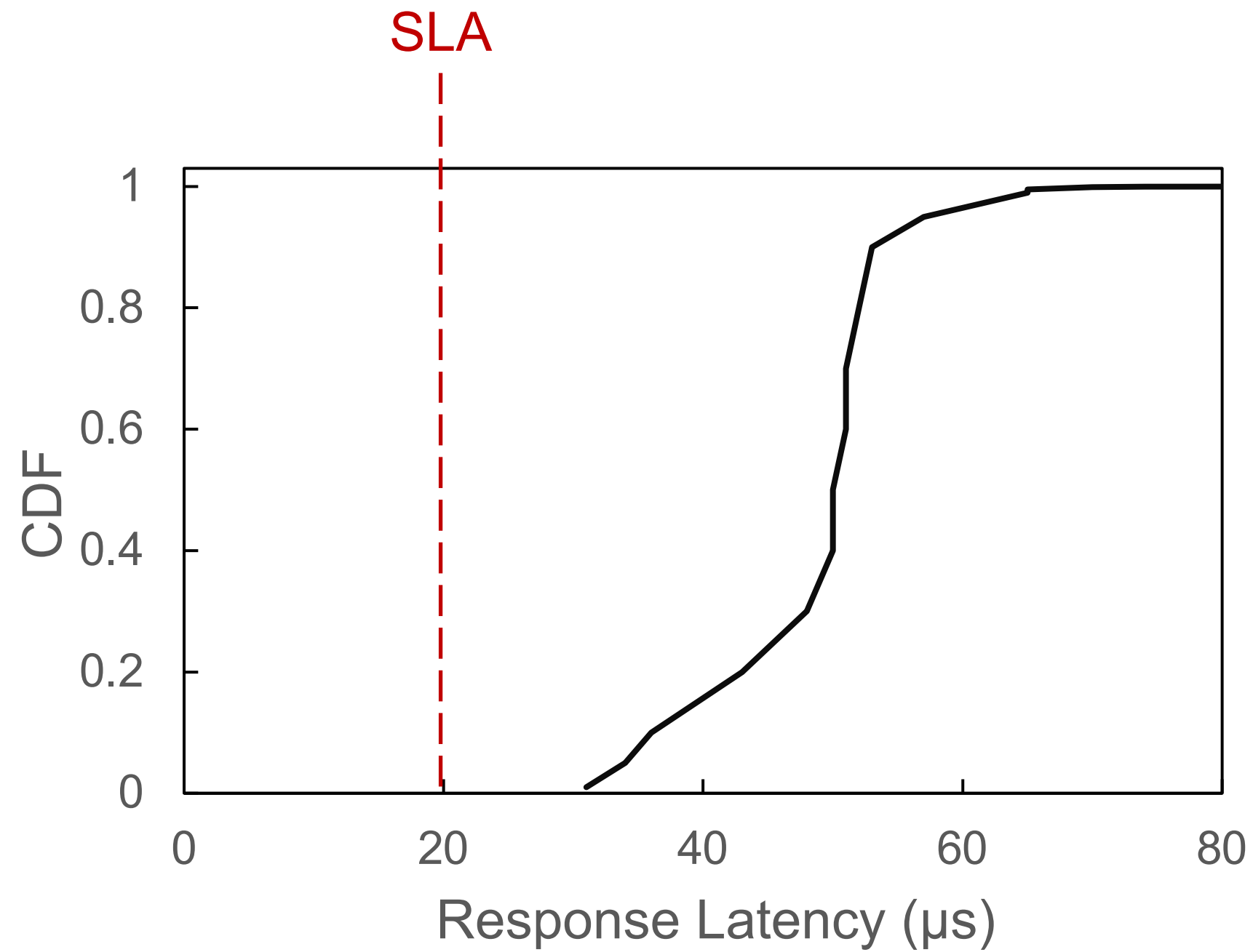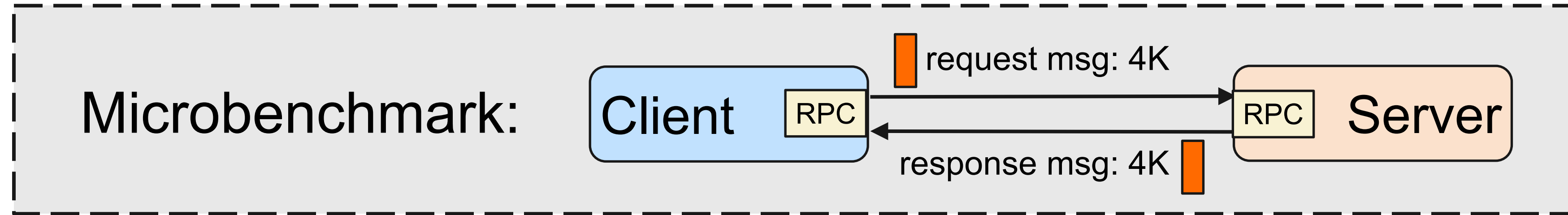
4

# Requirements of Frontend Networks



➢ **Diverse software and hardware**

➢ **Large amount of connections for inter-DC communications**

➢ **Engineering effort**

💡 **Kernel TCP can satisfy them all, but ...**

# Revisiting Kernel TCP
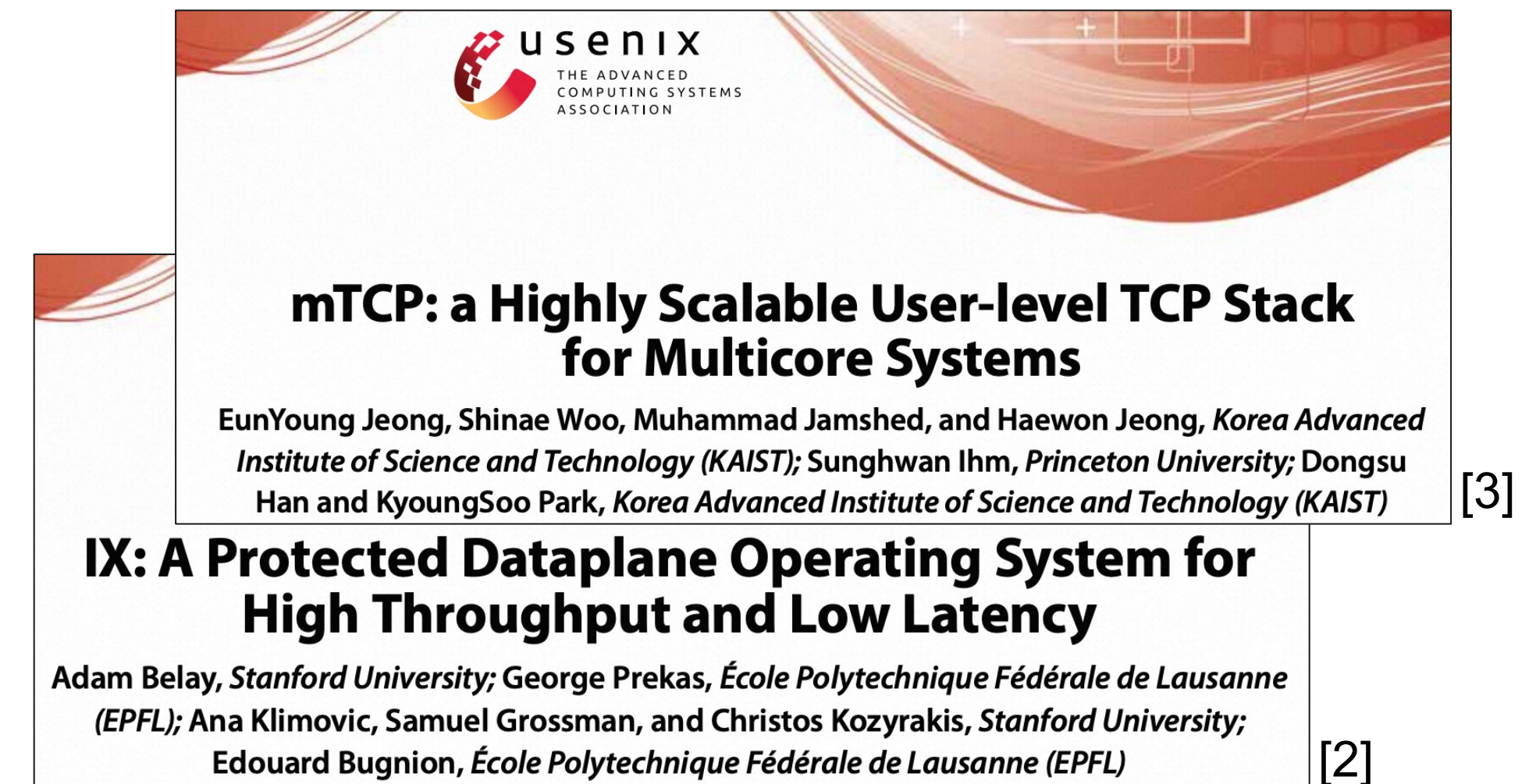


Microbenchmark:

Client — RPC → request msg: 4K → RPC — Server
Server — RPC → response msg: 4K → RPC — Client

CDF vs Response Latency (µs) — SLA at 20

Throughput (Gbps) vs #Cores — SLA at 40

**No longer suitable for cloud storage**

# Beyond Kernel TCP

## RDMA



When Cloud Storage Meets RDMA
Yixiao Gao, *Nanjing University and Alibaba Group;* Qiang Li, Lingbo Tang, Yongqing Xi, Pengcheng Zhang, Wenwen Peng, Bo Li, Yaohui Wu, Shaozong Liu, Lei Yan, Fei Feng, Yan Zhuang, Fan Liu, Pan Liu, Xingkui Liu, Zhongjie Wu, Junping Wu, and Zheng Cao, *Alibaba Group;* Chen Tian, *Nanjing University;* Jinbo Wu, Jiaji Zhu, Haiyong Wang, Dennis Cai, and Jiesheng Wu, *Alibaba Group* [1]

❌ Legacy devices
❌ Large amount of connections

## Existing User-space TCP

mTCP: a Highly Scalable User-level TCP Stack for Multicore Systems
EunYoung Jeong, Shinae Woo, Muhammad Jamshed, and Haewon Jeong, *Korea Advanced Institute of Science and Technology (KAIST);* Sunghwan Ihm, *Princeton University;* Dongsu Han and KyoungSoo Park, *Korea Advanced Institute of Science and Technology (KAIST)* [3]

IX: A Protected Dataplane Operating System for High Throughput and Low Latency
Adam Belay, *Stanford University;* George Prekas, *École Polytechnique Fédérale de Lausanne (EPFL);* Ana Klimovic, Samuel Grossman, and Christos Kozyrakis, *Stanford University;* Edouard Bugnion, *École Polytechnique Fédérale de Lausanne (EPFL)* [2]

❌ Performance issues
❌ Compatibility

### *Solution: building our own user-space TCP - Luna*

[1] NSDI'21, When Cloud Storage Meets RDMA
[2] NSDI'14, mTCP: A Highly Scalable User-level TCP Stack for Multicore Systems
[3] OSDI'14, IX: A Protected Dataplane Operating System for High Throughput and Low Latency
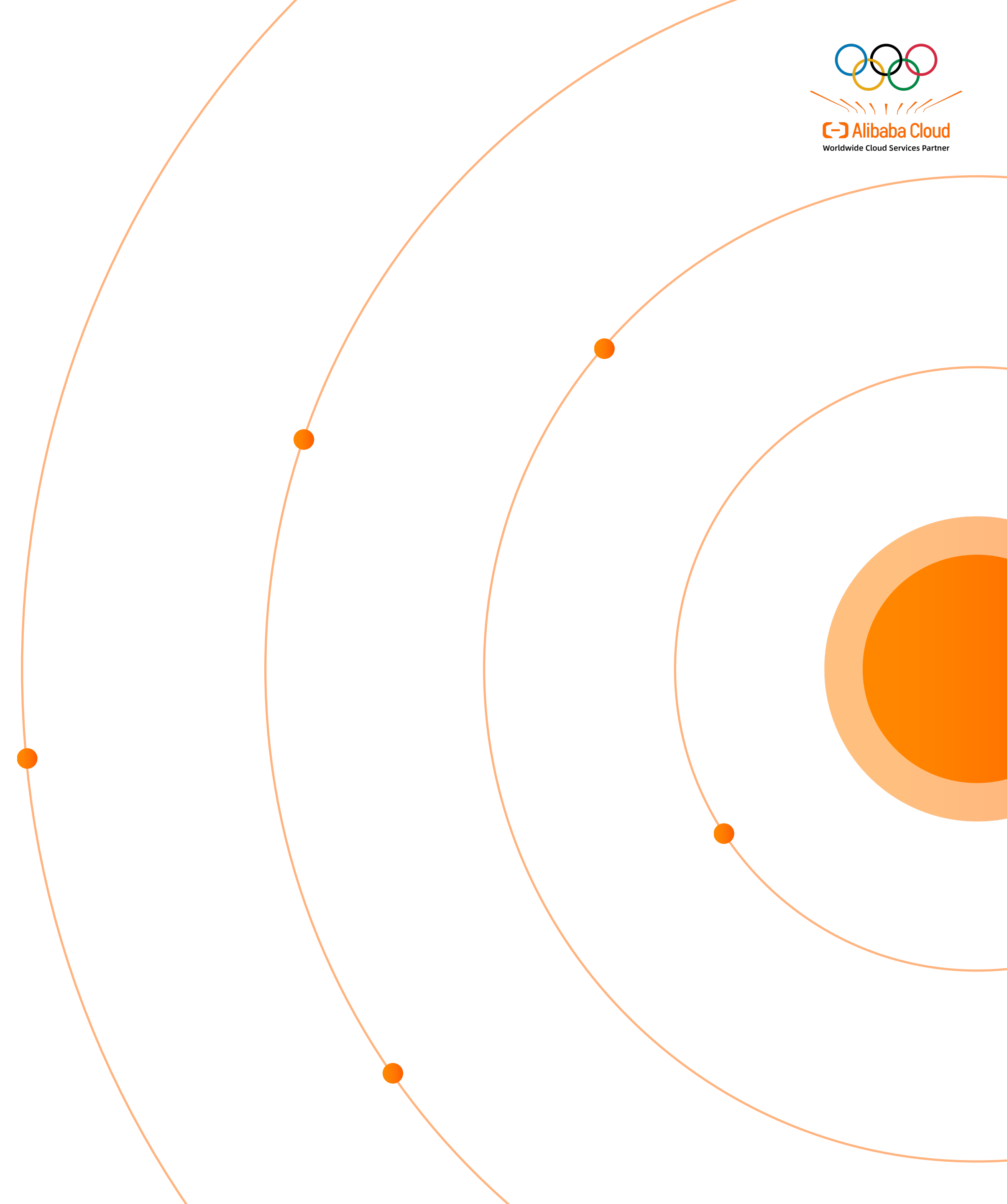
# Luna Architecture

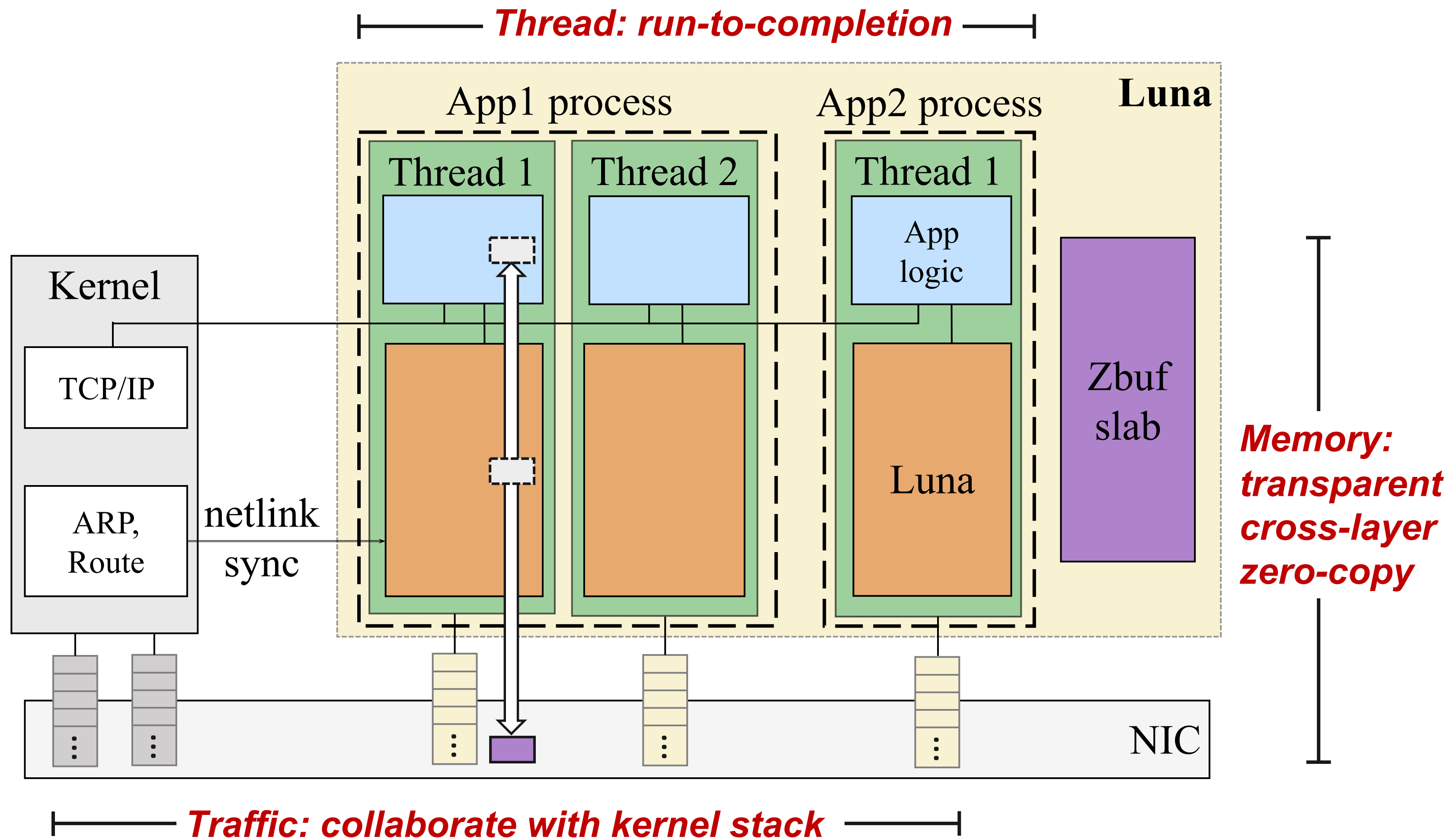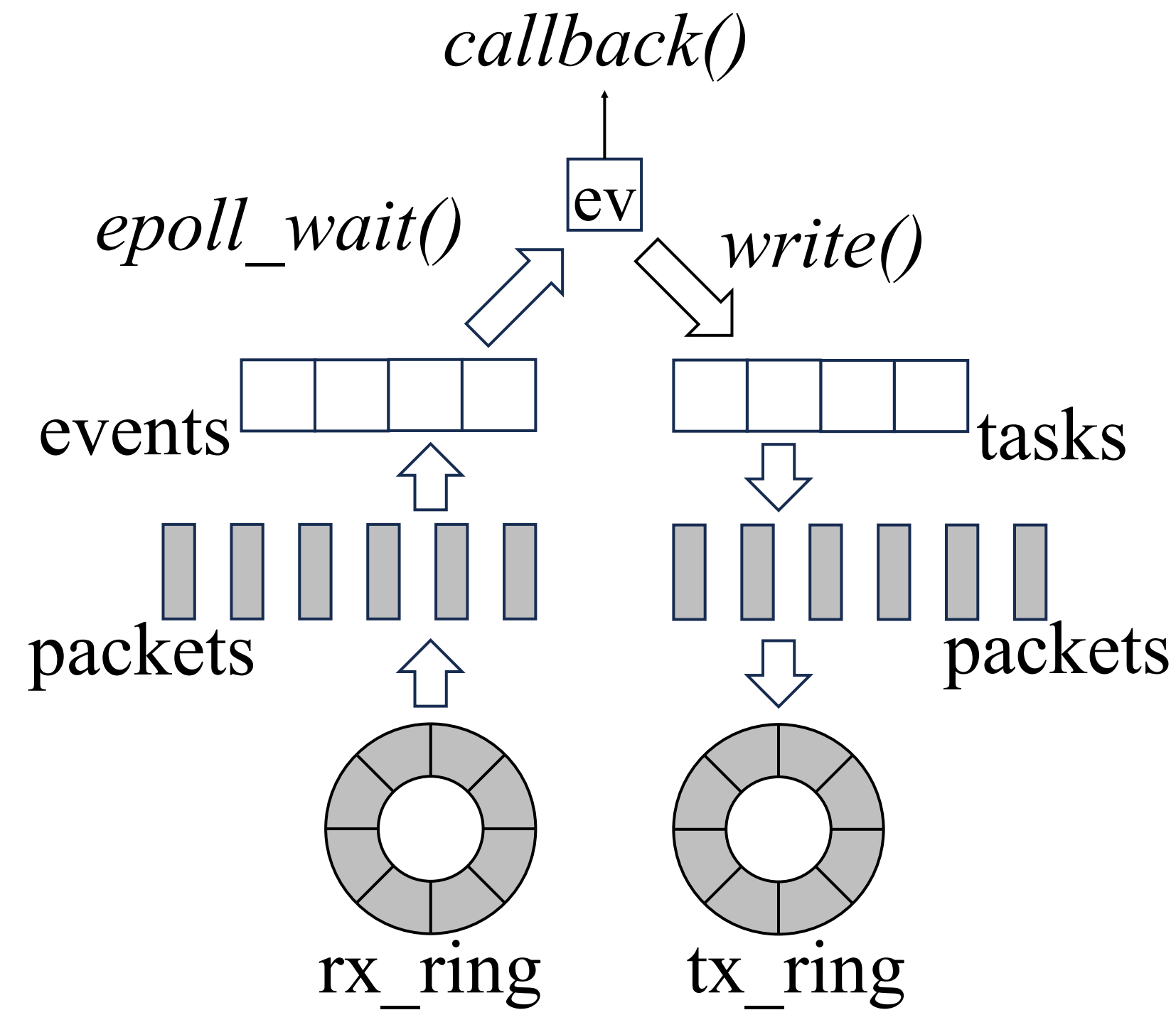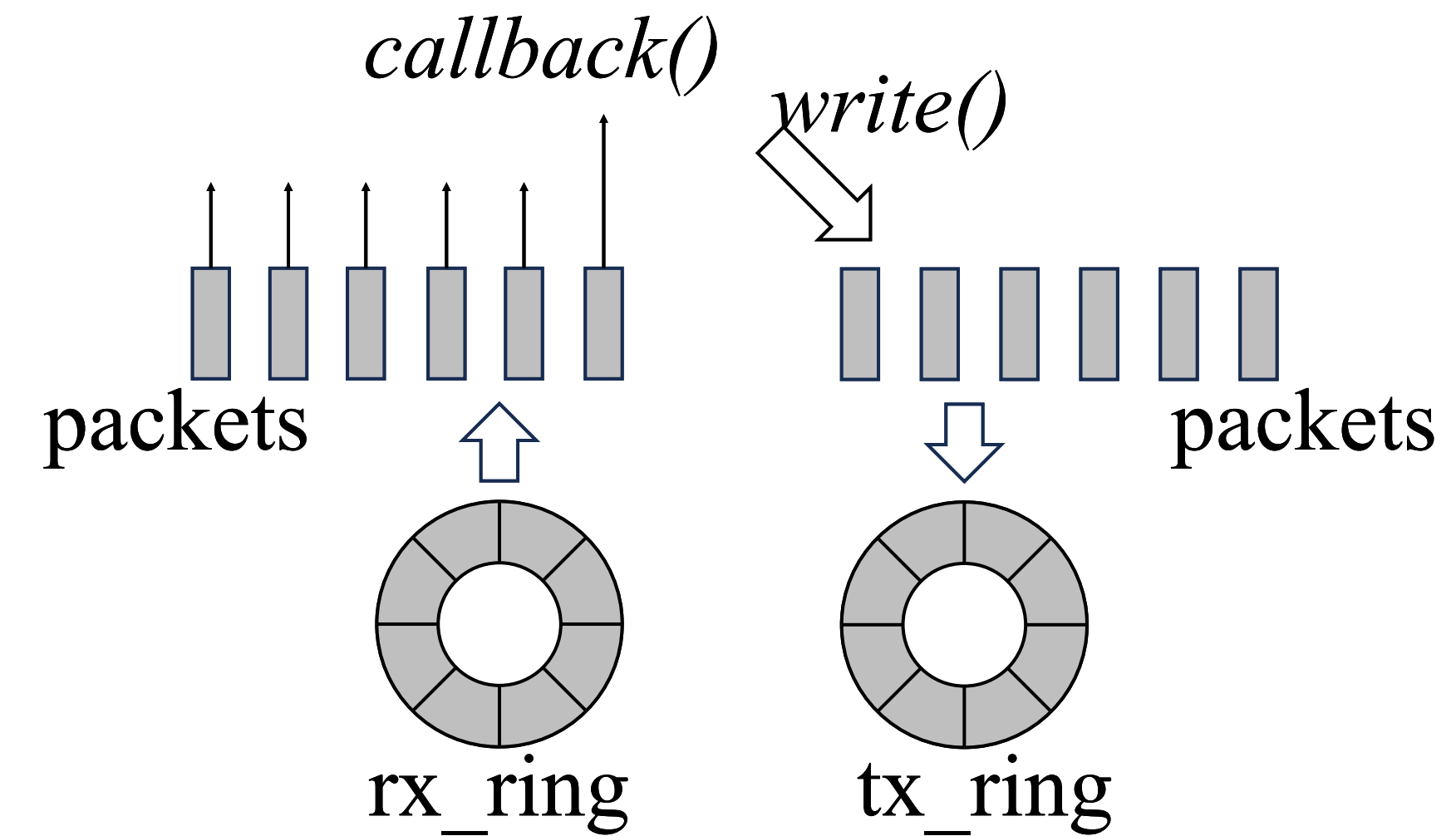# Thread Model: Run-to-completion



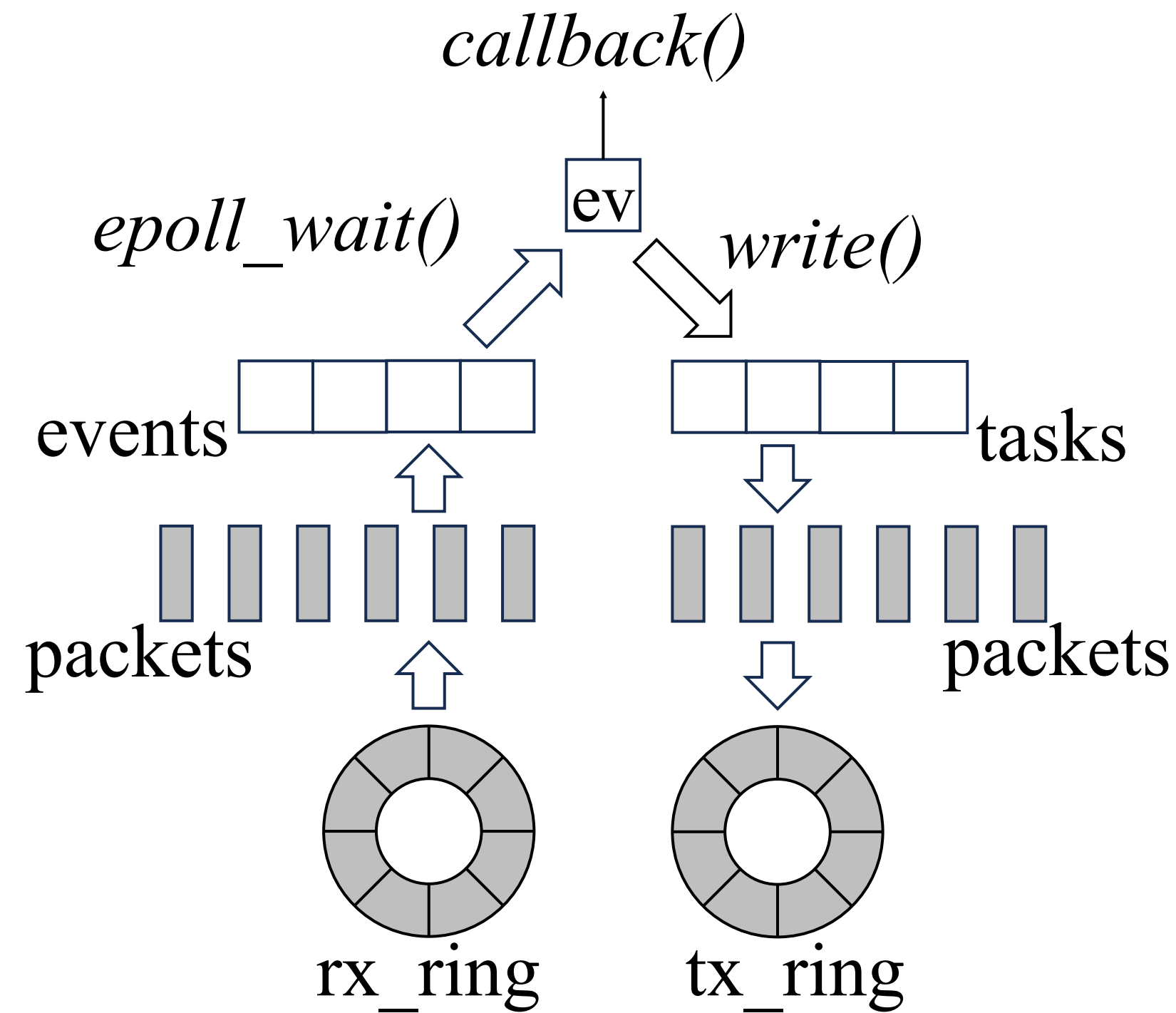Batch-r2c
*(similar to mTCP[1], IX[2])*

Inline-r2c

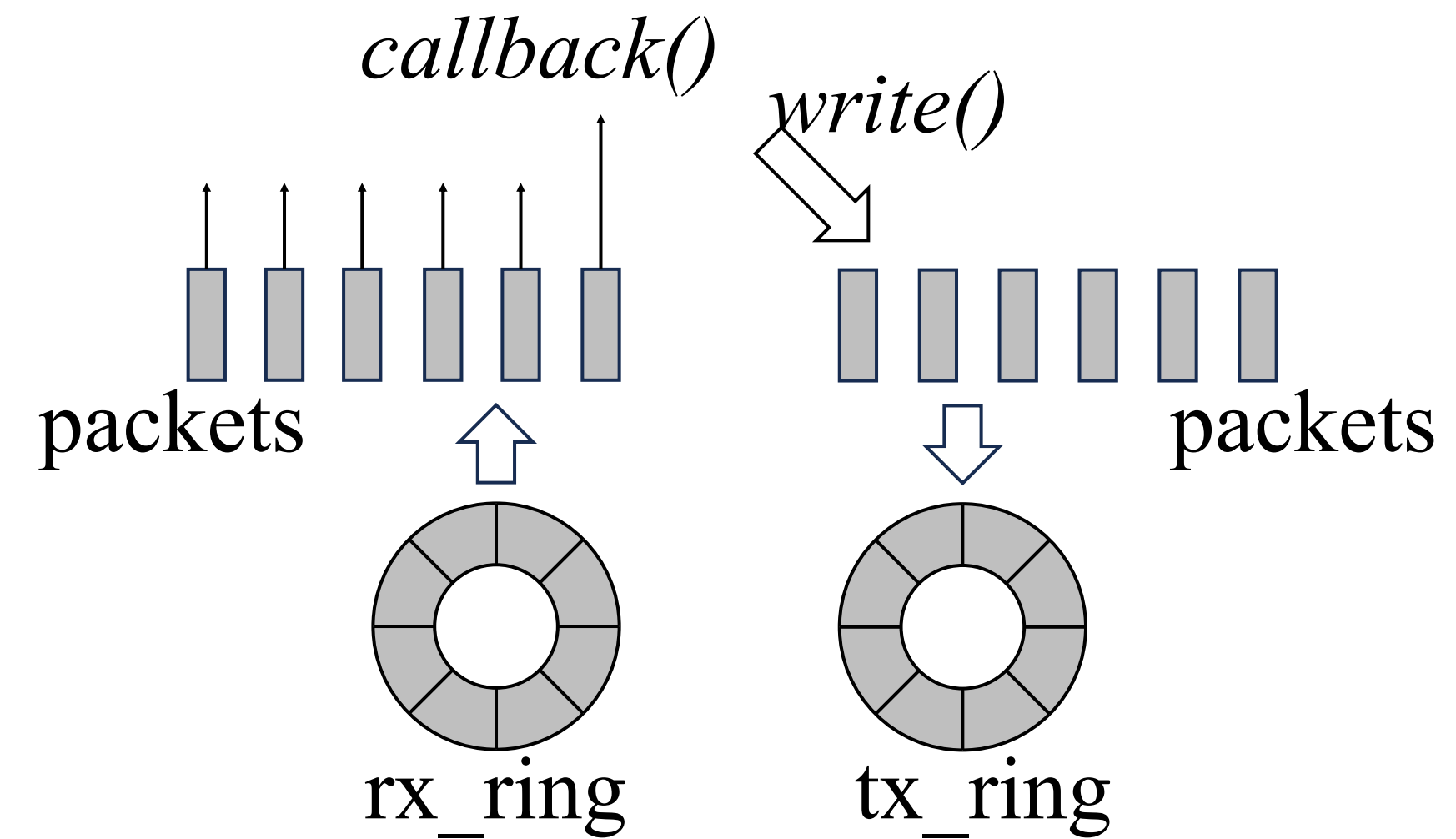[1] NSDI'14, mTCP: A Highly Scalable User-level TCP Stack for Multicore Systems
[2] OSDI'14, IX: A Protected Dataplane Operating System for High Throughput and Low Latency
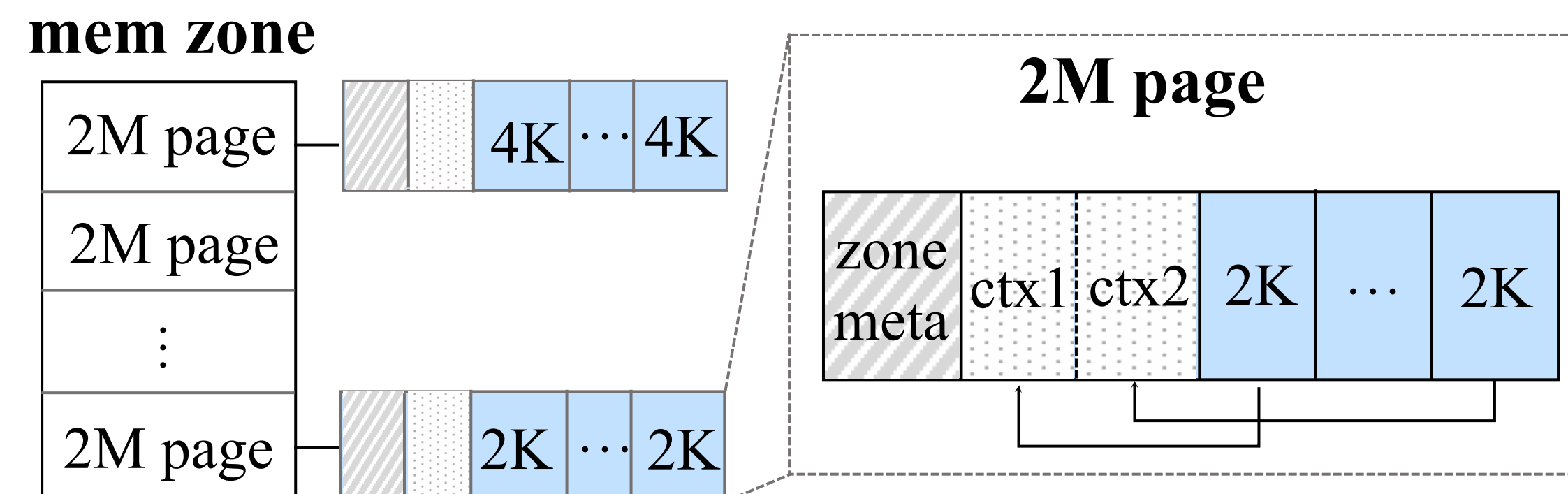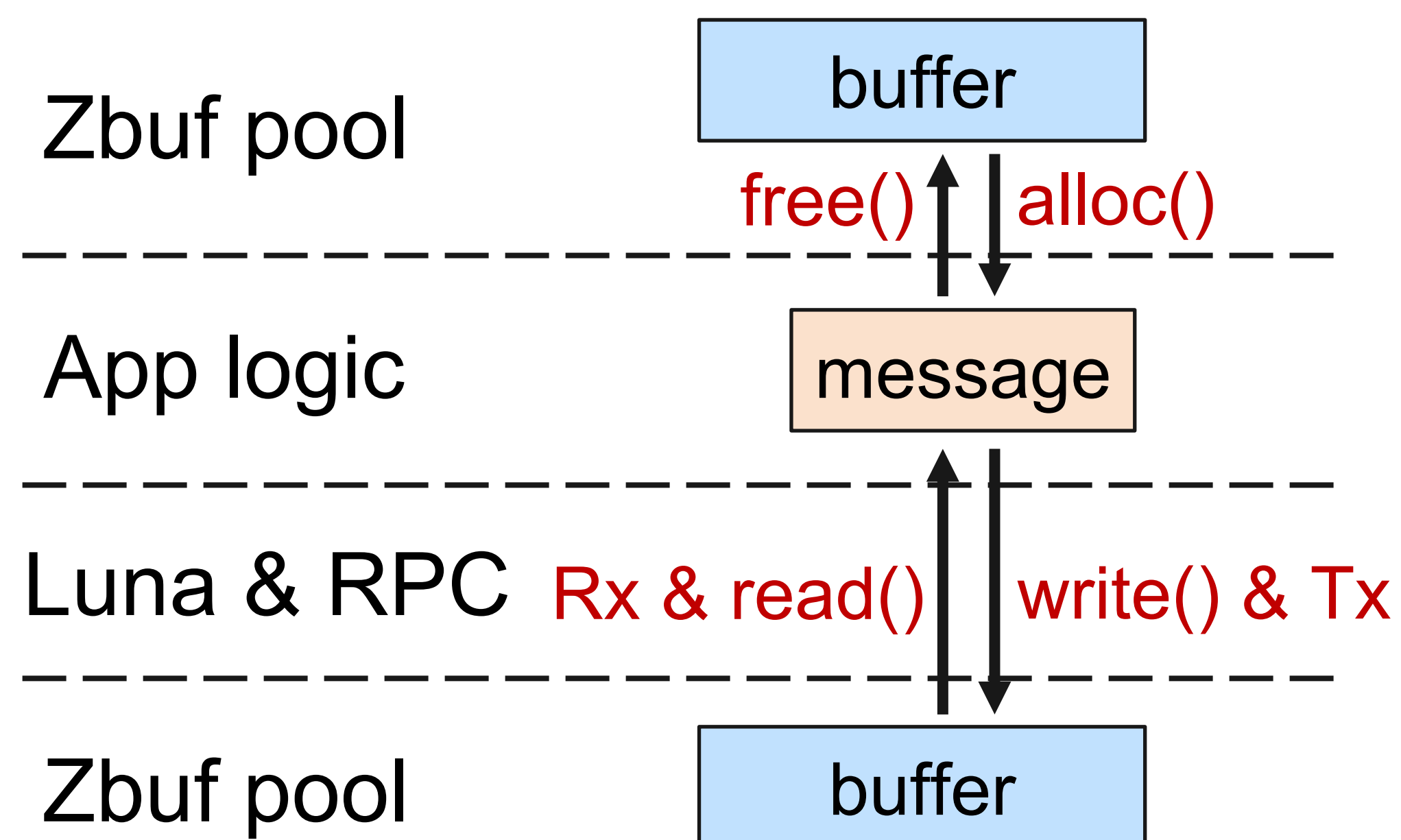
# Thread Model: Run-to-completion



*callback()*

*epoll_wait()*  ev  *write()*

events          tasks

packets         packets

rx_ring   tx_ring

## Batch-r2c

✓  Pros: Compatibility
✗  Cons: Event framework overhead

*callback()*   *write()*

packets         packets

rx_ring   tx_ring

## Inline-r2c

✓  Pros: Performance
✗  Cons: Programing model change

# Memory Model: Cross-layer Zero-copy

Zbuf pool    **buffer**

free() ↑ ↓ alloc()

App logic    message

Luna & RPC    Rx & read() ↑ ↓ write() & Tx

Zbuf pool    **buffer**

**mem zone**

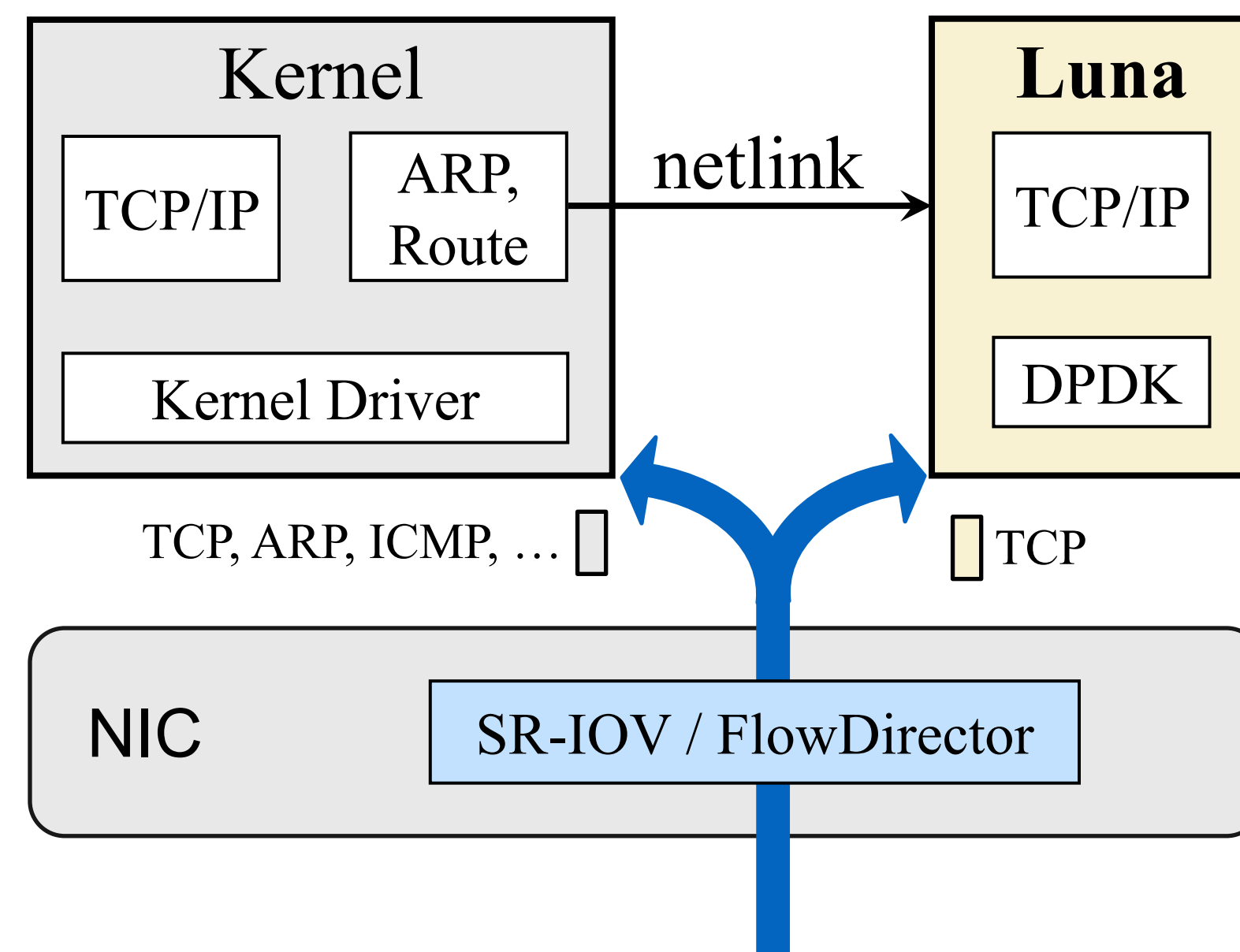| 2M page | — | 4K | ⋯ | 4K |
| 2M page | | | | |
| ⋮ | | | | |
| 2M page | — | 2K | ⋯ | 2K |

**2M page**

| zone meta | ctx1 | ctx2 | 2K | ⋯ | 2K |

Transparent buffer management —— *just like heap* 😊

# Traffic Model: Collaborate with Kernel Stack



**Other user-space TCP**

✗ Exclusive with kernel stack Apps

✗ Complex control plane logic and configuration

*Luna*

✓ Cohost with other Apps
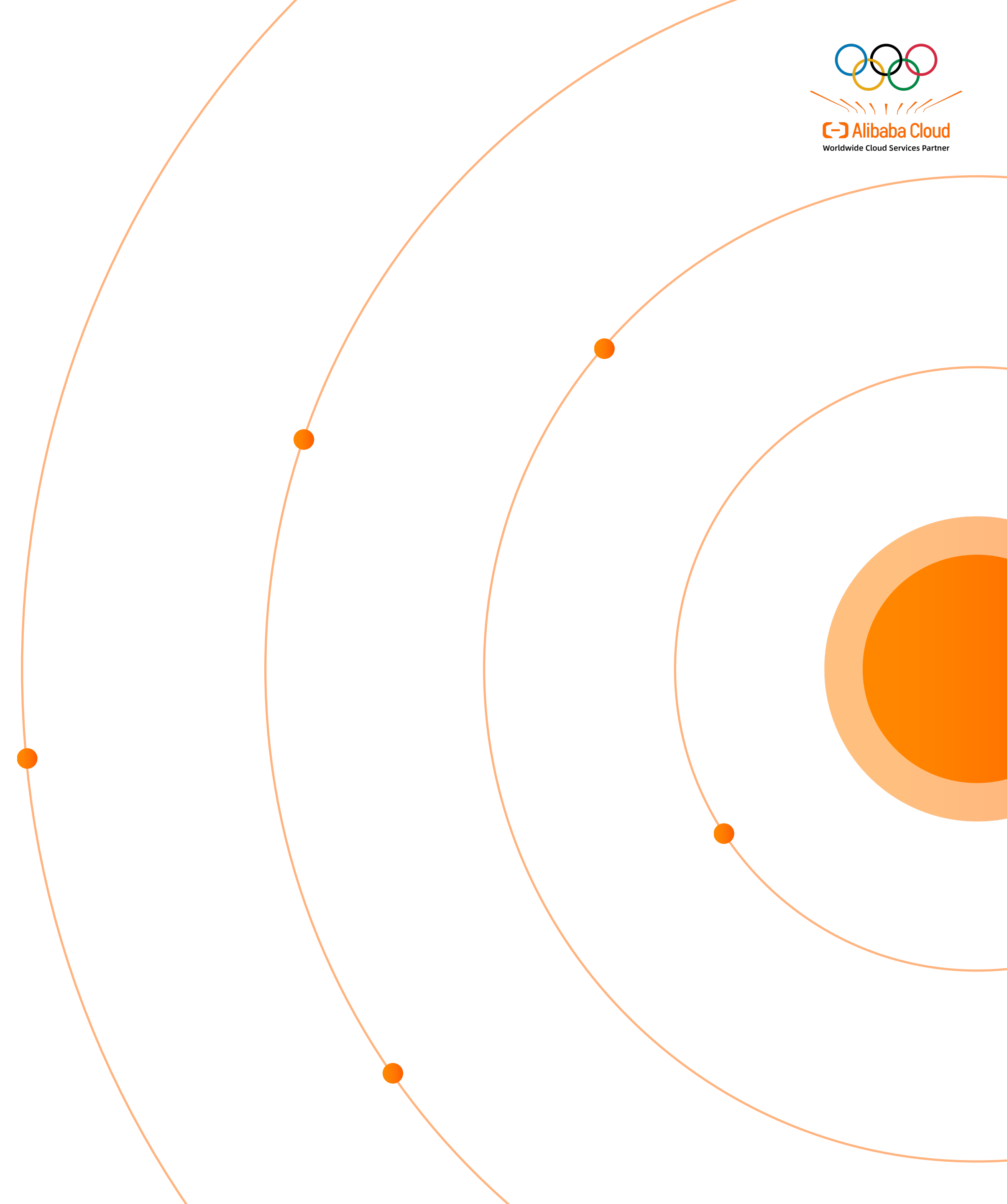
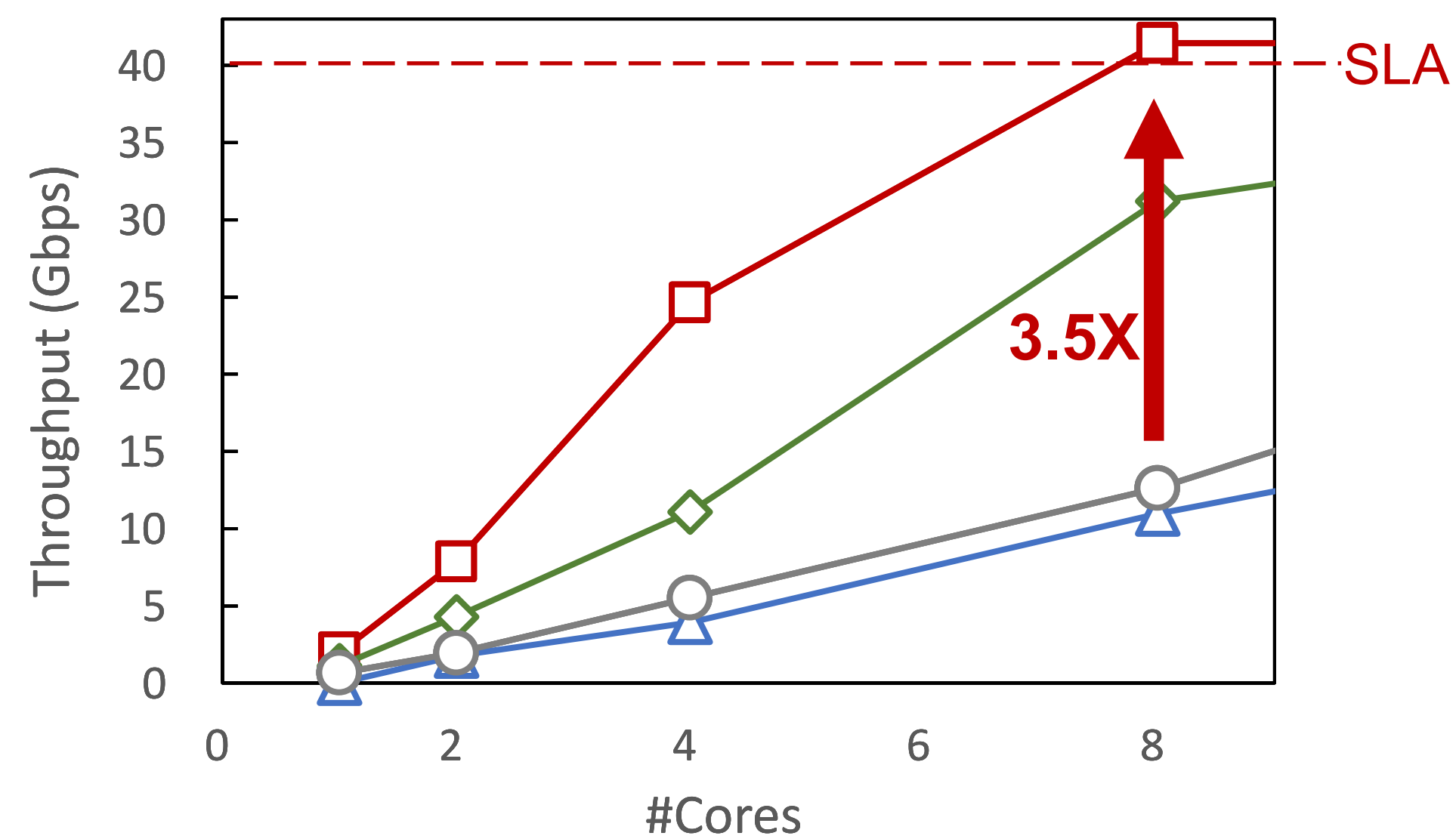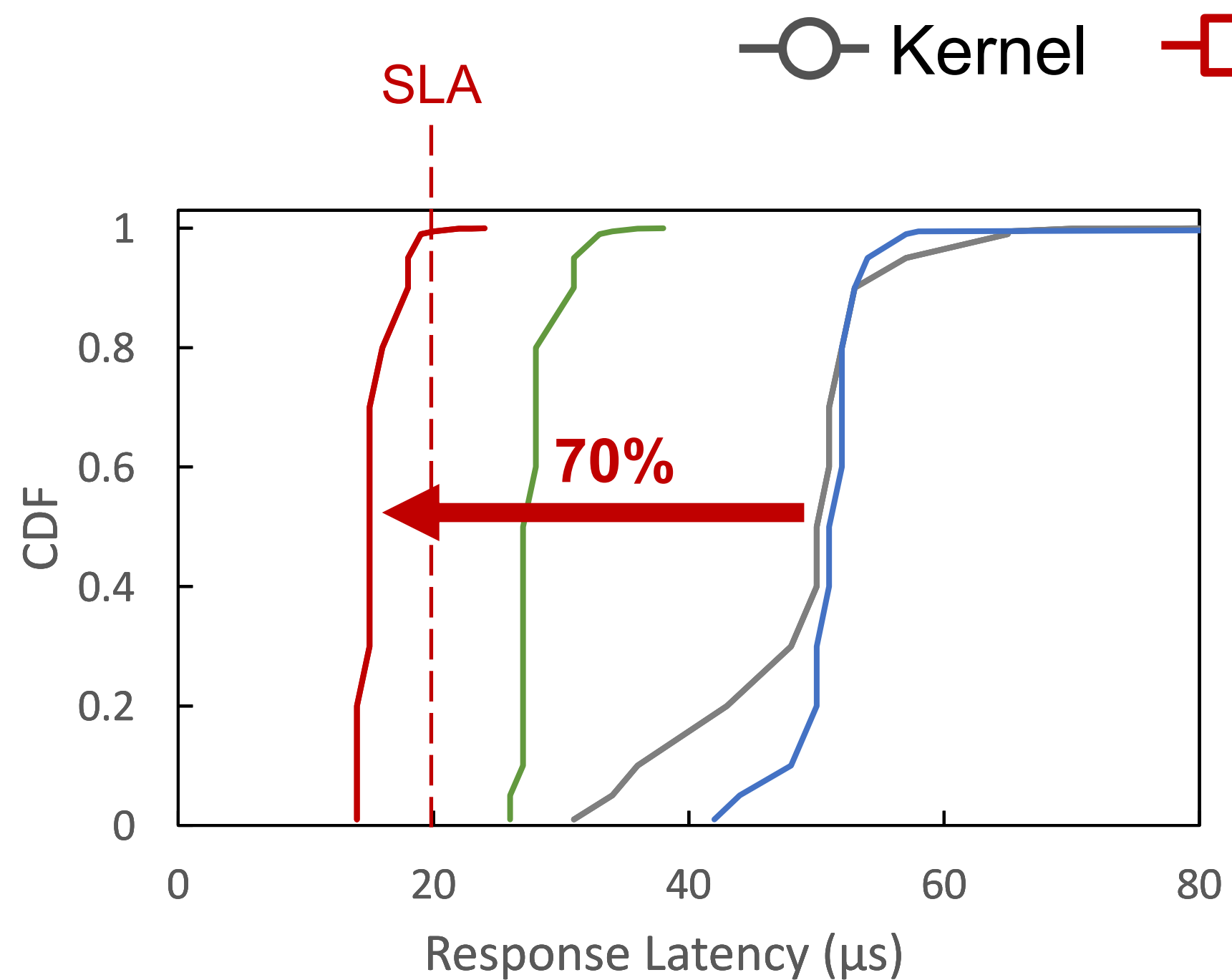✓ Get rid of complex control plane implementation

# Evaluation - Microbenchamrk

○— Kernel    □— Luna    ◇— mTCP    △— VPP

Elastic Block Storage

Table Store Service
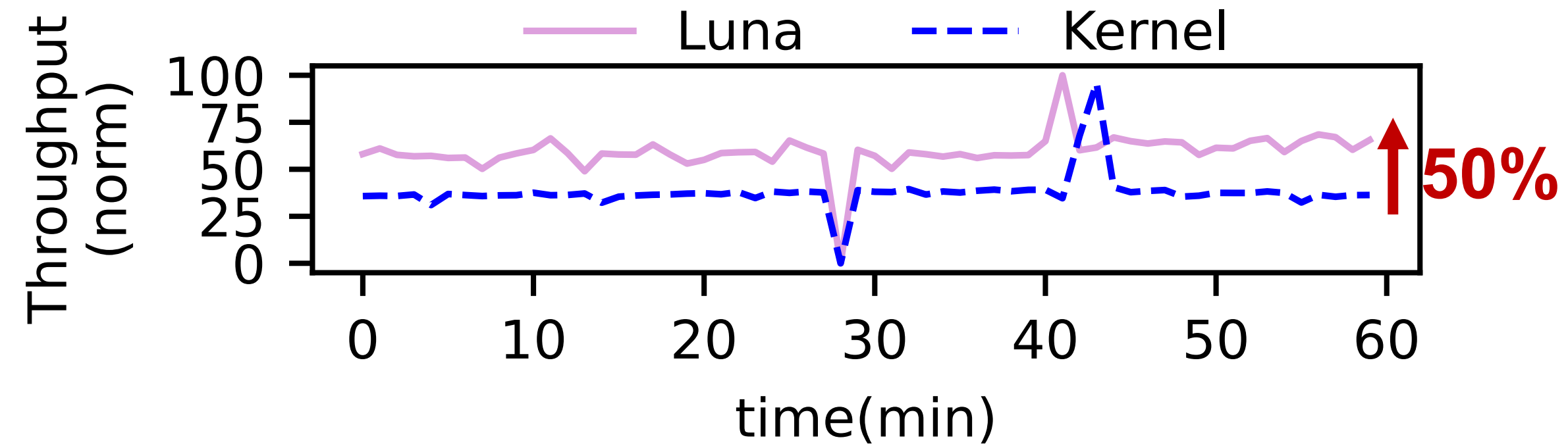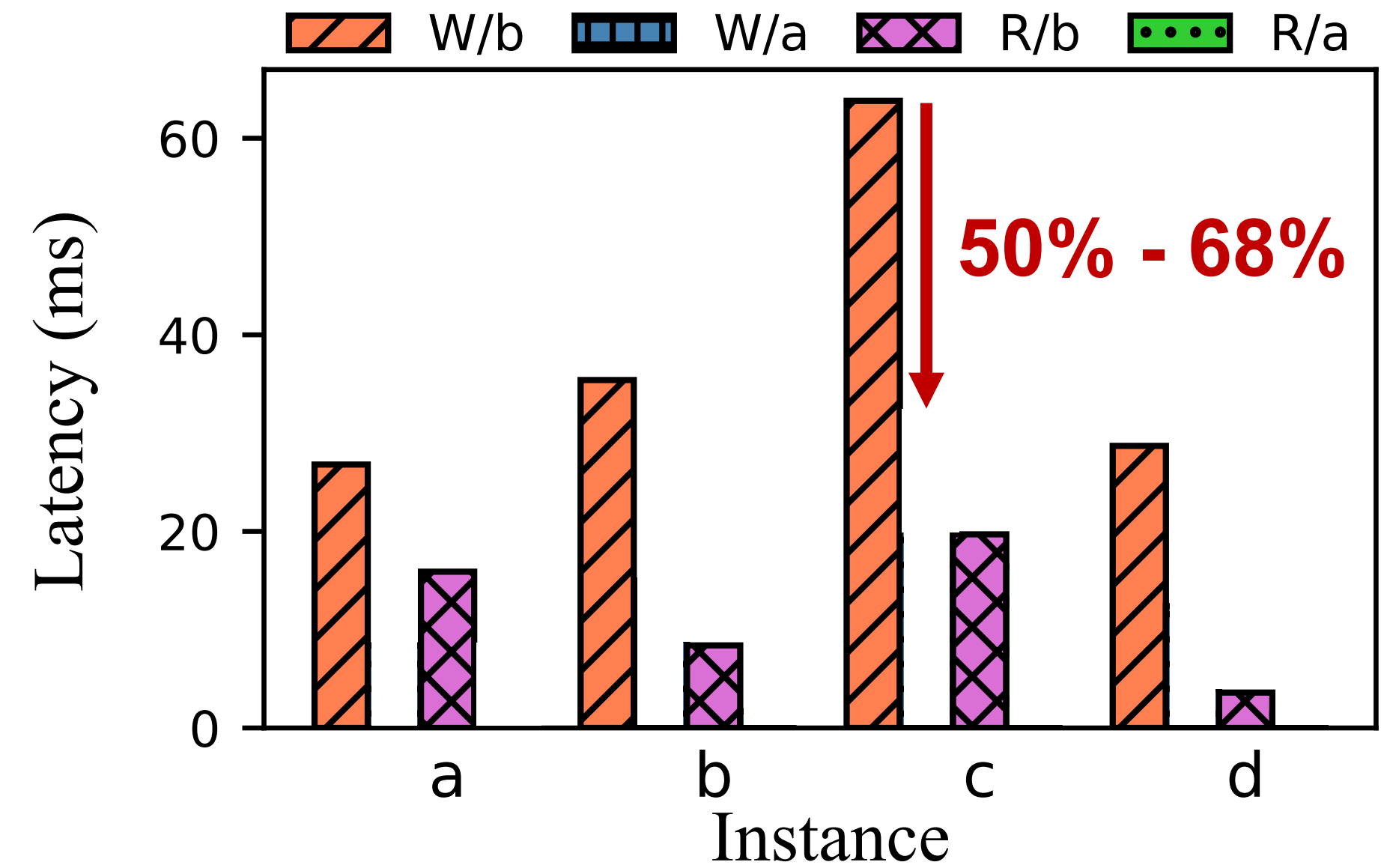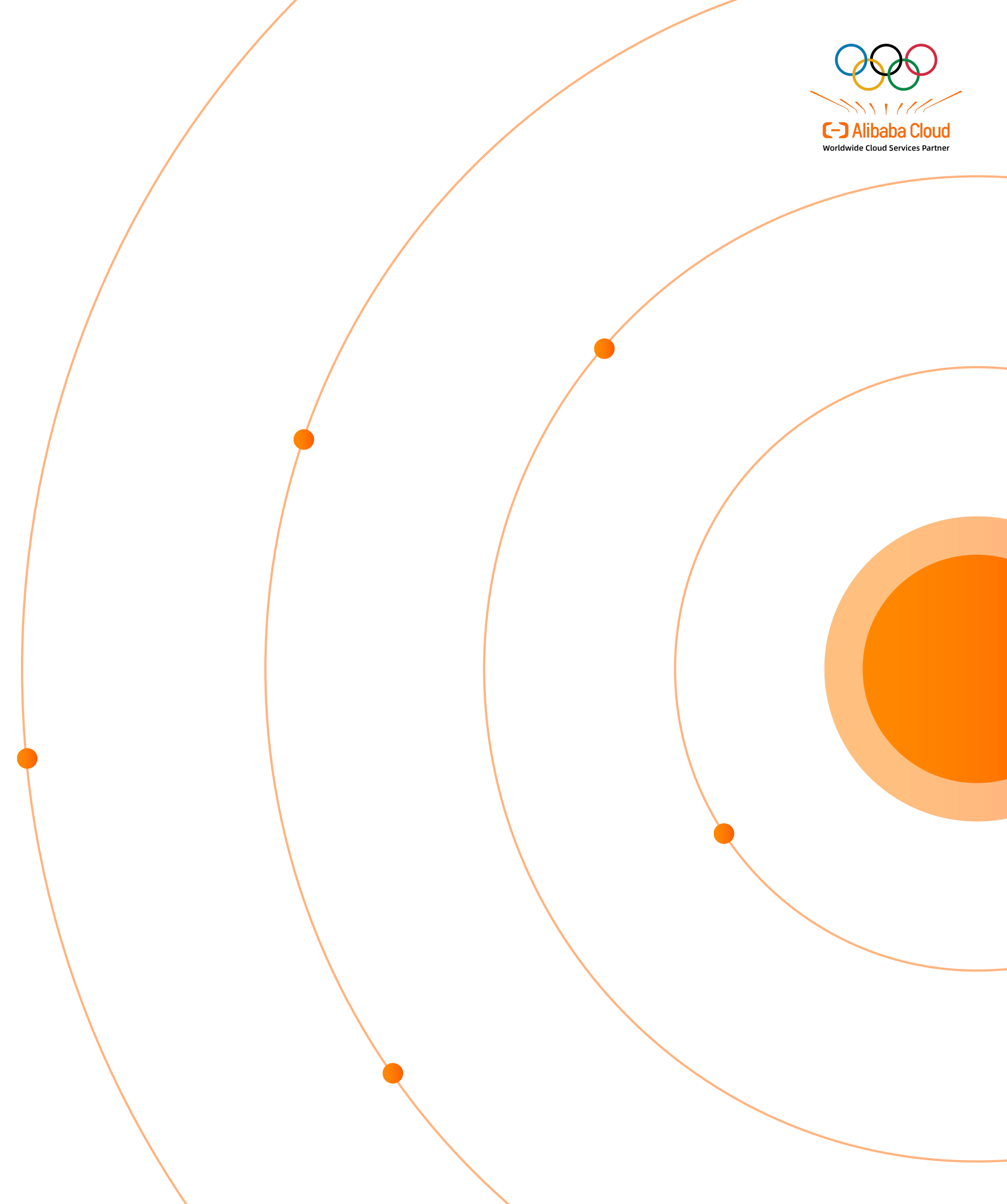
Background & Motivation

Design

Evaluation

**Discussion & Lessons**

Conclusion

# Design Choices: LibOS with R2C

LibOS

## Microkernel

APP₁ APP₂

Snap

## Sidecar

APP₁ APP₂

## Dispatch

APP₁ APP₂

## r2c

APP₁ APP₂

*Luna*

APP    Network Stack    process    thread

➤ Performance is the most critical
➤ No frequent upgrades

# Design Choices: Customized TCP

New transport
protocols

pFabric,

pHost,

…

High-performance,
but complex

Customized TCP

Simplified fast path
NewReno
Fast recovery

Simple, but works

# Network Evolution
## *(in Alibaba Cloud Storage)*

UDP-based,
multipath,
hardware offload,
DPU co-designed,
…

iov-oriented API,
inline-r2c,
Tx zero-copy

*Solar*

new protocol to a larger scale

Socket-like API,
batch-r2c,
copy on Tx

LUNA-customized
RPC

LUNA

default protocol

Kernel TCP

[1] SigComm 2022, From luna to solar

2018          2019          2020          20

# Conclusion

➢ Luna, a user-space network designed for cloud storage service

➢ Network architecture of Alibaba Cloud Storage

➢ Three design pieces

- Thread model: run-to-completion

- Memory model: cross-layer zero-copy

- Traffic model: collaborate with kernel stack

➢ Discussion

- LibOS vs. Microkernel

- TCP and tailoring

- Evolution

# Thanks for Listening!
## Q&A

*Email: mashu.ms@alibaba-inc.com*