



湖南大學  
HUNAN UNIVERSITY



# UCM: Fast and Maintainable User-space RDMA Connection Setup

Huijun Shen<sup>1</sup>, Jian Yang<sup>2</sup>, Zelong Yue<sup>2</sup>, Xingyu Guo<sup>1</sup>, Xijin Yin<sup>1</sup>, Lang An<sup>2</sup>, Yulin Chen<sup>2</sup>, Jie Ding<sup>2</sup>,  
Hongyu Wu<sup>2</sup>, Yong Zhang<sup>2</sup>, Jianxi Ye<sup>2</sup>, Guo Chen<sup>1</sup>

<sup>1</sup>Hunan University <sup>2</sup>ByteDance

# RDMA is widely deployed

---

- RDMA in production-level applications:
  - Cloud storage, Recommendation system, LLM inference/training...



...

- RDMA connection setup approaches
  - Socket APIs : out-of-band connection
  - CM\* APIs : in-band connection

\* The **RDMA CM** is a **communication manager** used to setup reliable, connected and unreliable datagram data transfers, and provides standard APIs defined by librdmacm library. [https://man7.org/linux/man-pages/man7/rdma\\_cm.7.html](https://man7.org/linux/man-pages/man7/rdma_cm.7.html)

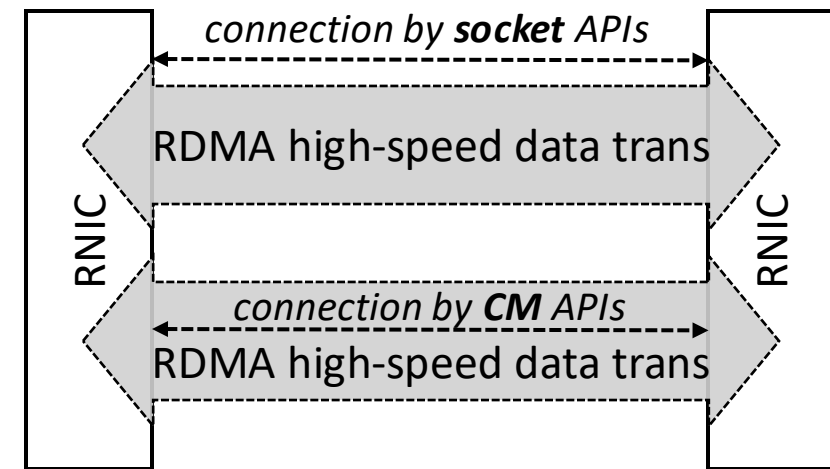
# RDMA is widely deployed

- RDMA in production-level applications:
  - Cloud storage, Recommendation system, LLM inference/training...



...

- RDMA connection setup approaches
  - Socket APIs : out-of-band connection
  - CM\* APIs : in-band connection



\* The **RDMA CM** is a **communication manager** used to setup reliable, connected and unreliable datagram data transfers, and provides standard APIs defined by librdmacm library. [https://man7.org/linux/man-pages/man7/rdma\\_cm.7.html](https://man7.org/linux/man-pages/man7/rdma_cm.7.html)

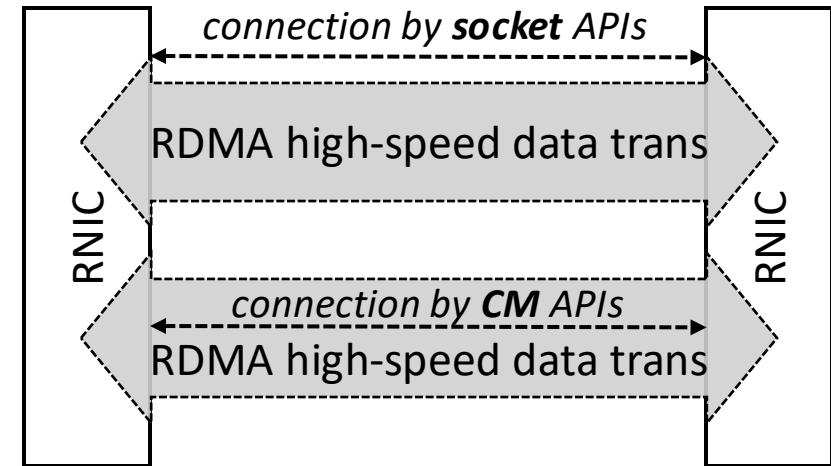
# RDMA is widely deployed

- RDMA in production-level applications:
  - Cloud storage, Recommendation system, LLM inference/training...



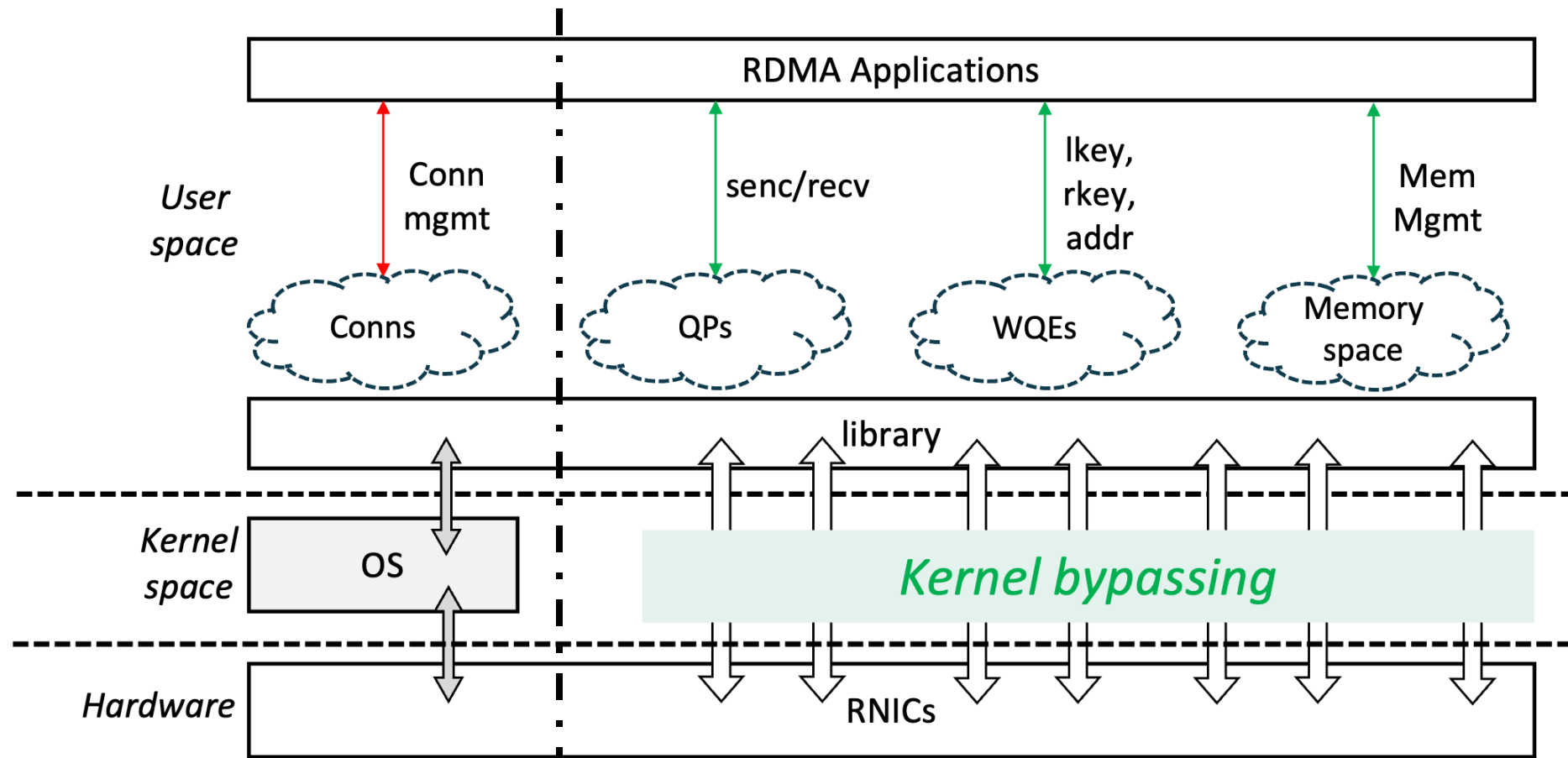
...

- RDMA connection setup approaches
  - Socket APIs : out-of-band connection
  - **CM\* APIs : in-band connection**



\* The **RDMA CM** is a **communication manager** used to setup reliable, connected and unreliable datagram data transfers, and provides standard APIs defined by librdmacm library. [https://man7.org/linux/man-pages/man7/rdma\\_cm.7.html](https://man7.org/linux/man-pages/man7/rdma_cm.7.html)

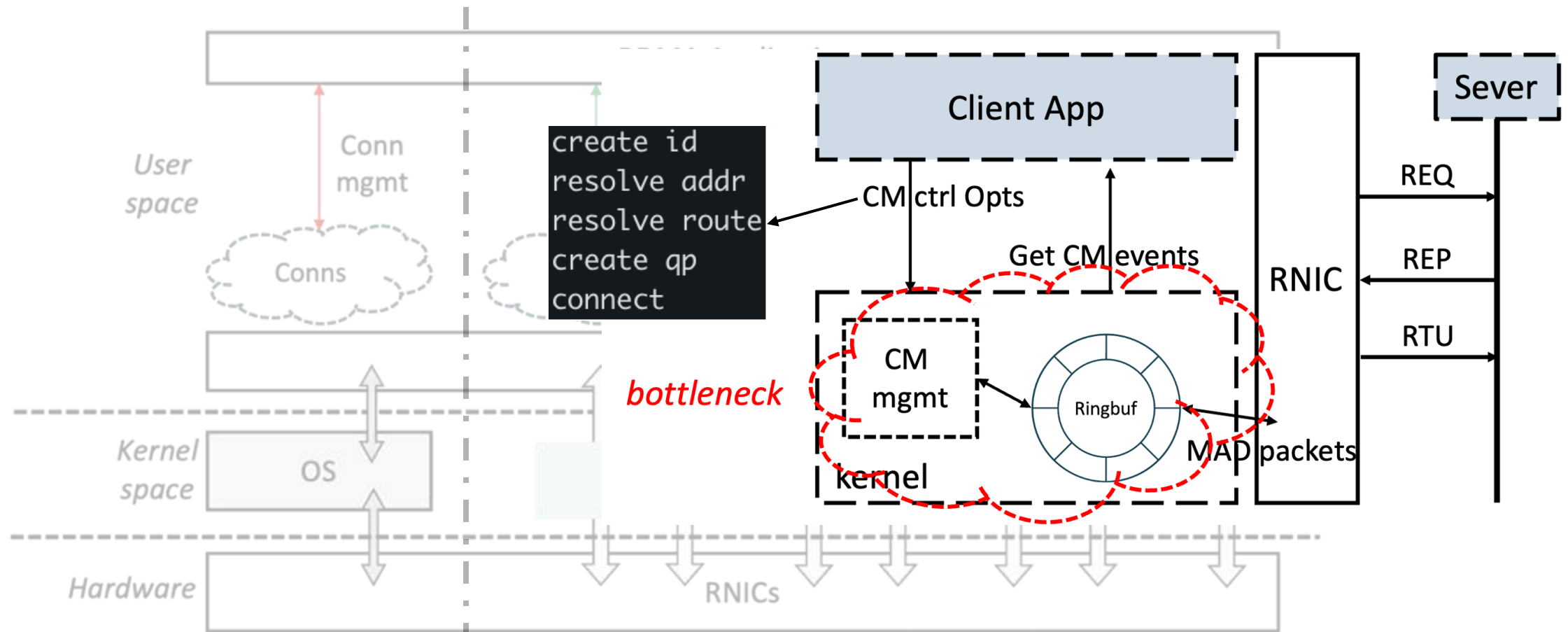
# RDMA fast path and slow path



➤ Slow control path: Connection mgmt, *etc.*

➤ Fast data path: Data trans , Memory mgmt *etc.*

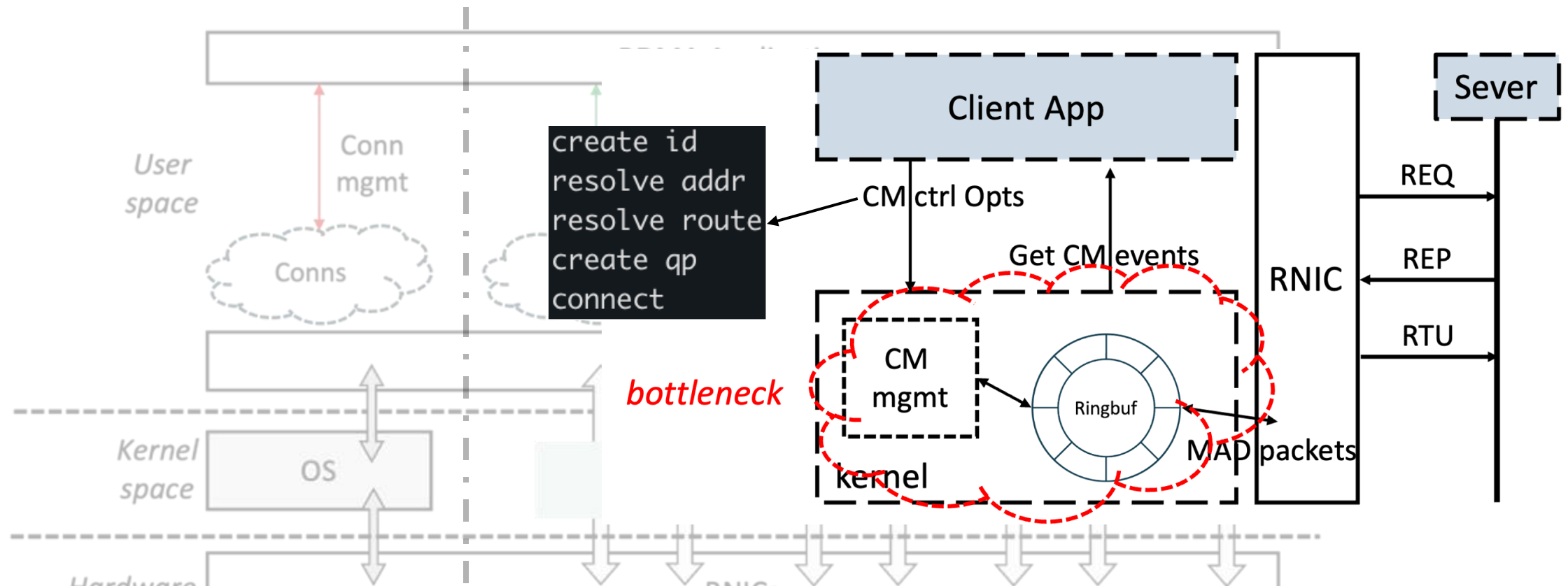
# RDMA control path is slow



➤ Slow control path: Connection mgmt, *etc.*

➤ Fast data path: Data trans, Memory mgmt *etc.*

# RDMA control path is slow



**RDMA conn setup process is inefficient and difficult to monitor.**

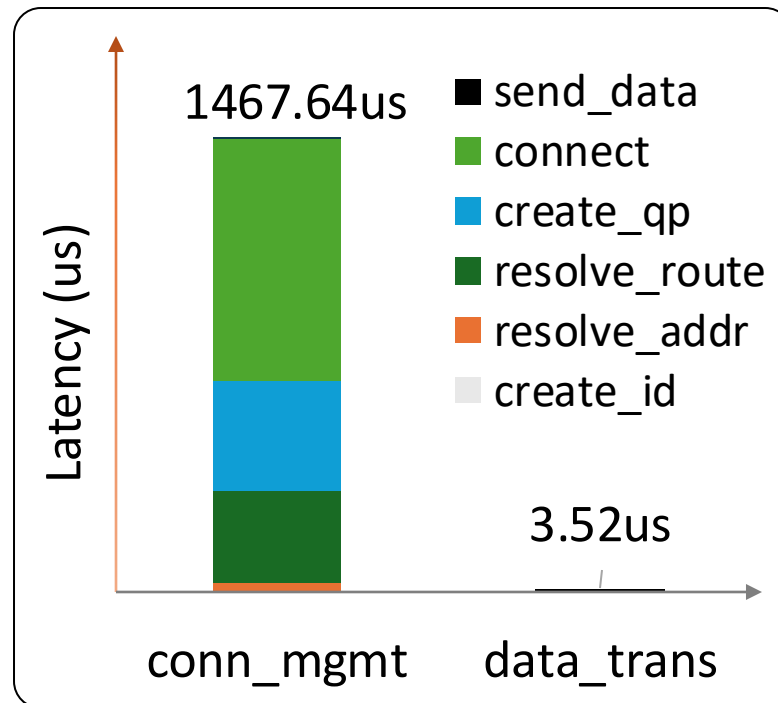
➤ Slow control path. Connection mgmt, etc. ➤ Fast data path. Data trans, memory mgmt etc.

Connection mgmt, etc.

# Problems of RDMA connection setup

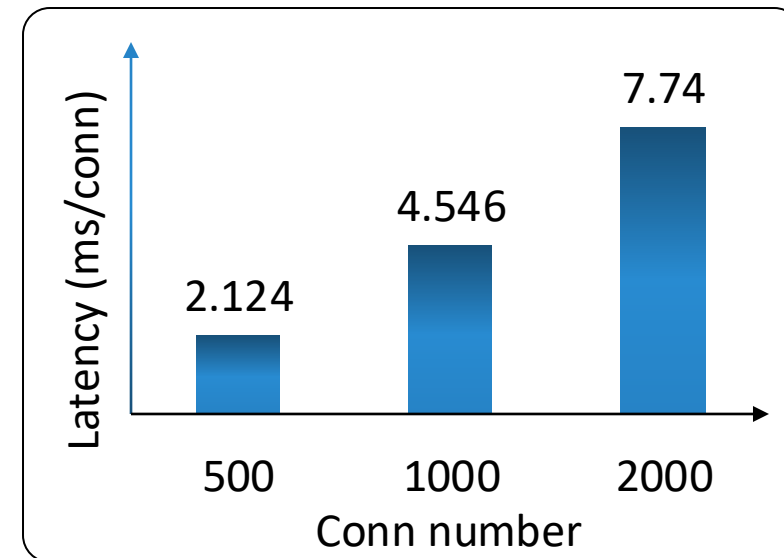
## ■ RDMA connection management (CM) is **Inefficient**.

Overheads in one RDMA transmission



➤ **High cost:** 1-2ms/conn.

Overheads in large-scale connections

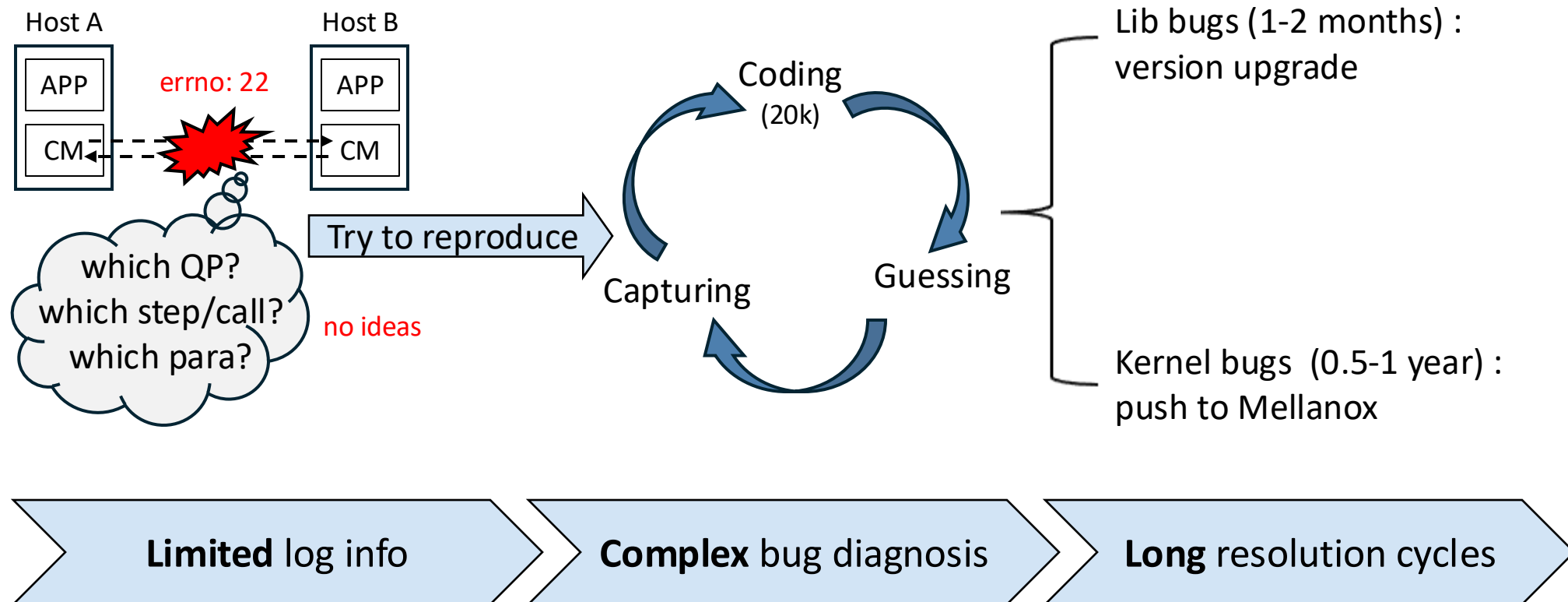


➤ **Bad scalability:** Connection setup efficiency further decreases as scale increases.



# Problems of RDMA connection setup

## ■ Production Deployment Practices for RDMA Connections



# Goals

---

## ■ Production Deployment Practices for RDMA Connections

Host A

Host B

Lib bugs (1-2 months) :

***Our goal is to develop an new User-space  
RDMA CM setup approach***

which para.

kernel bugs (0.5-1 year) .  
push to Mellanox

Limited log info

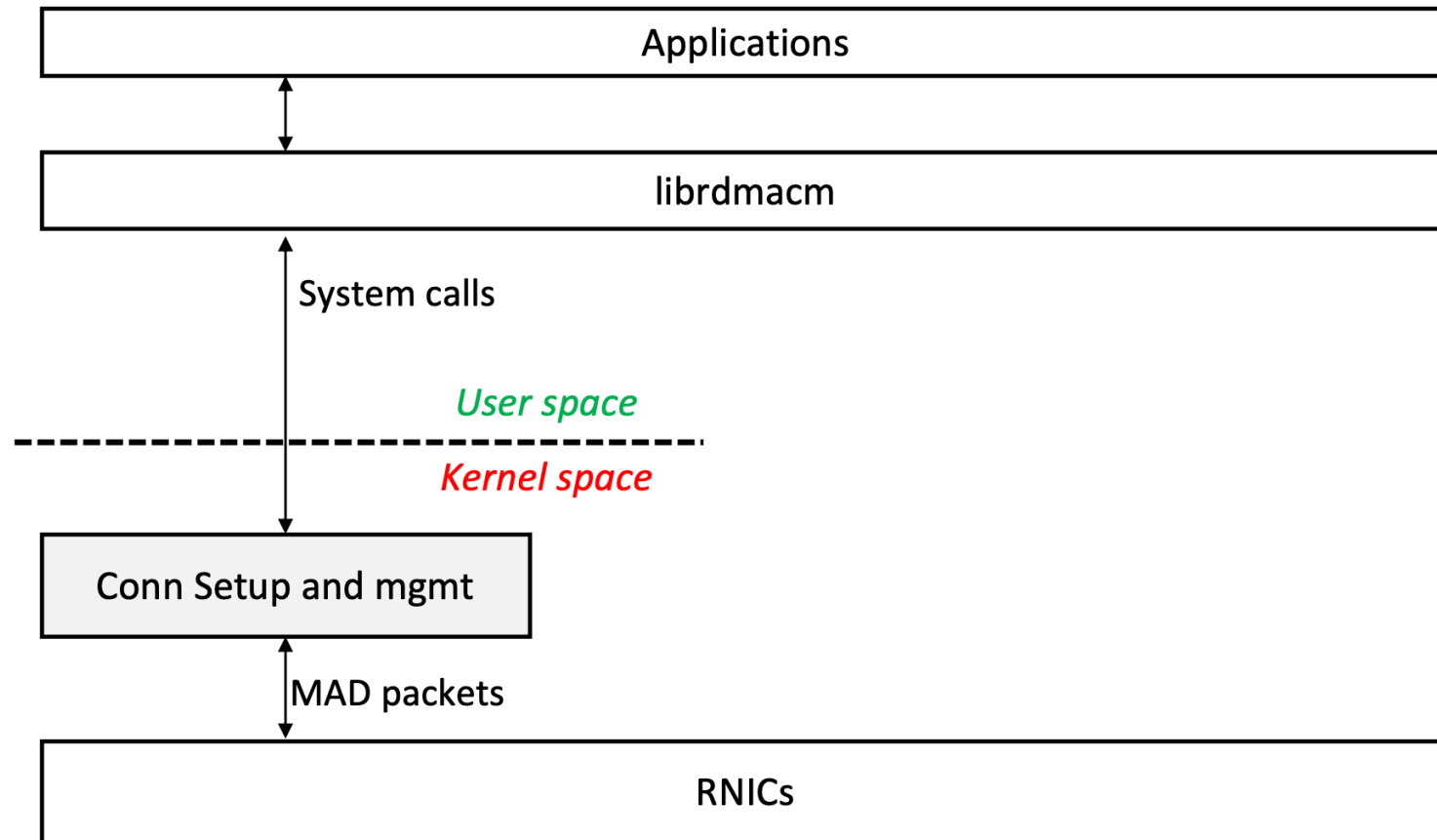
Complex bug diagnosis

Long resolution cycles

# From Kernel to User-space

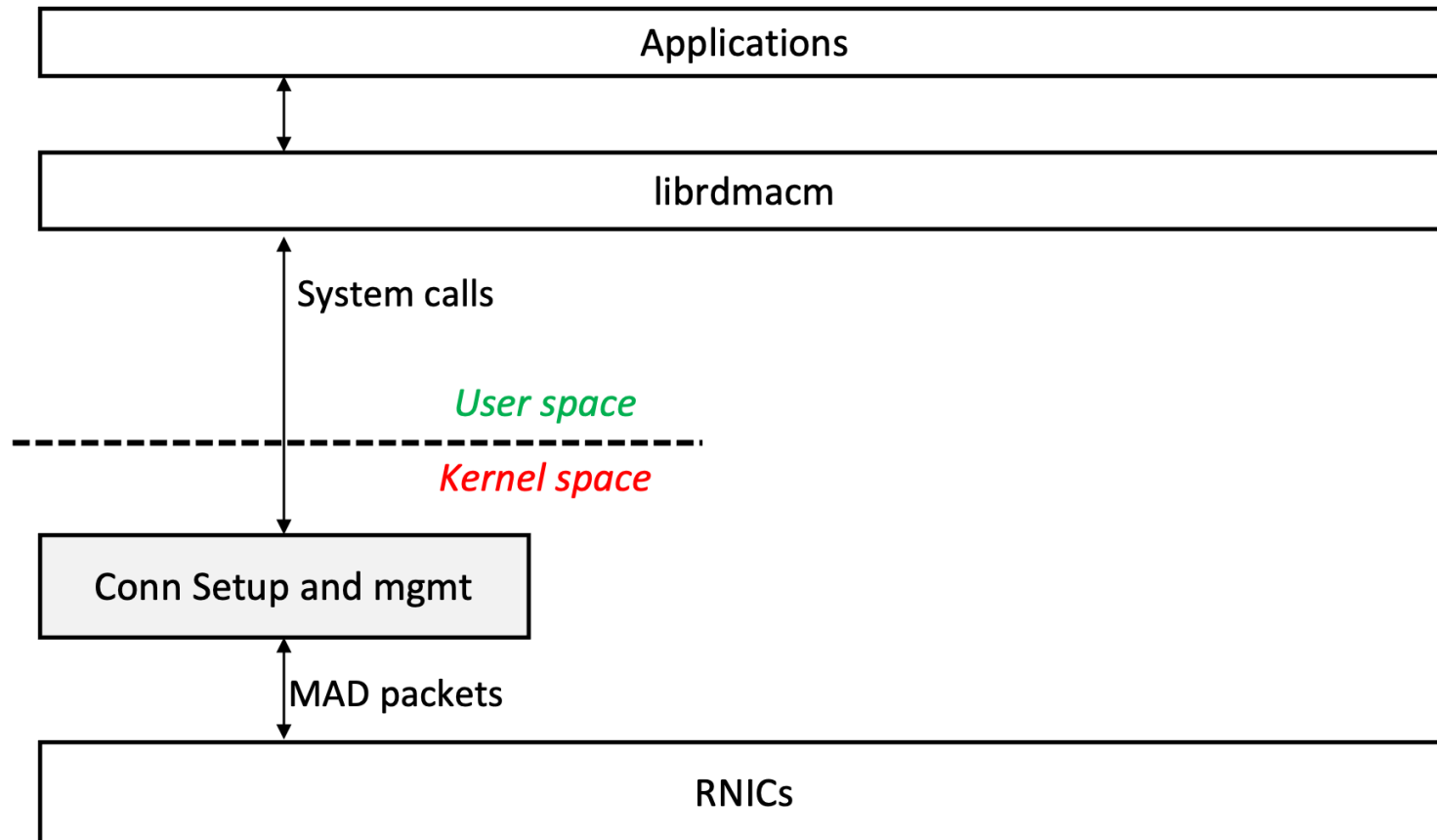
---

## ➤ Original RDMA connection management in Linux kernel



# From Kernel to User-space

## ➤ Original RDMA connection management in Linux kernel

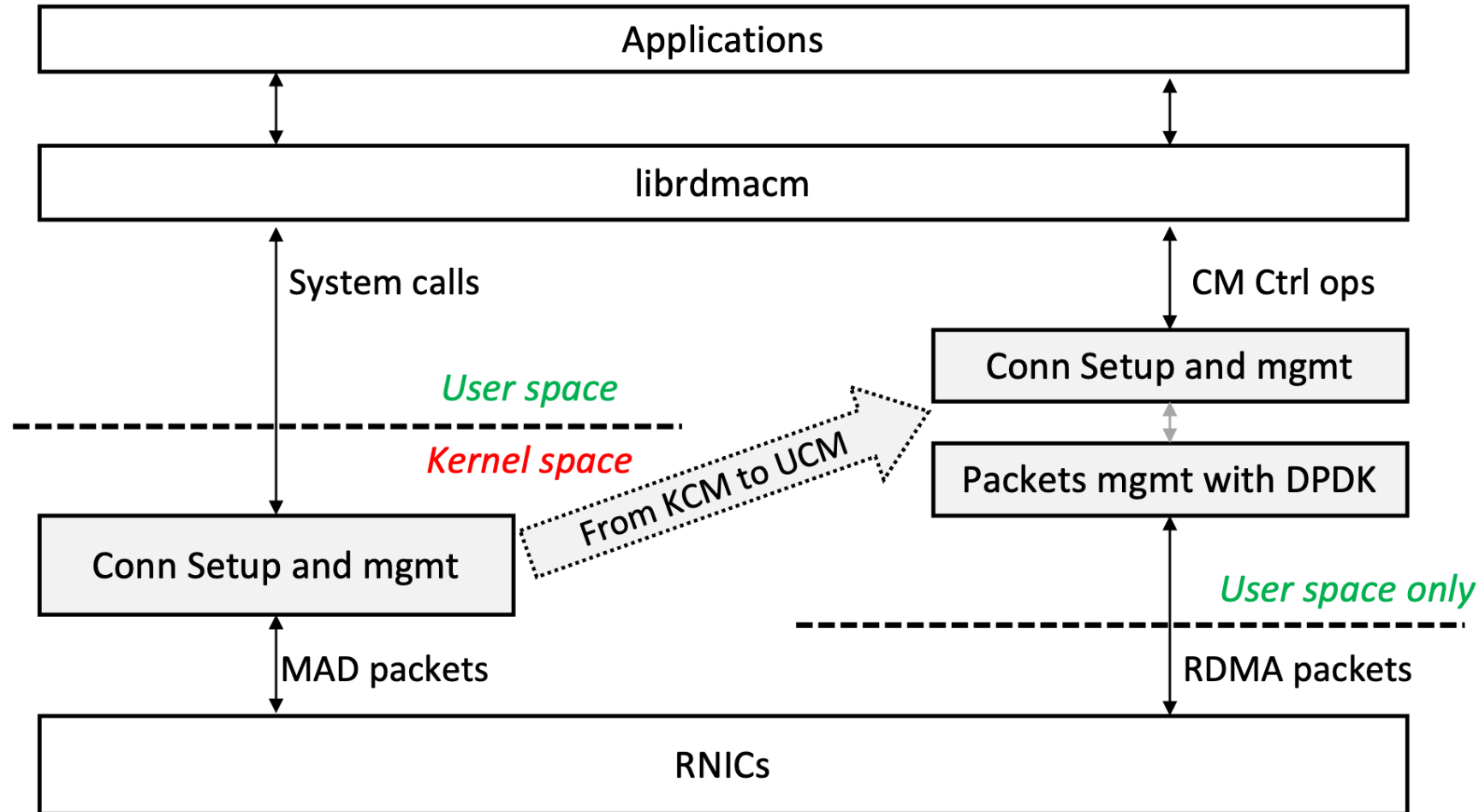


### Insights:

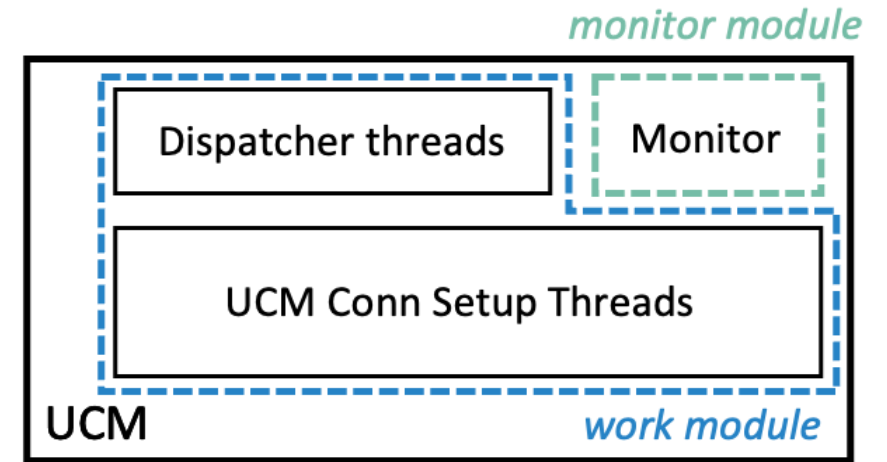
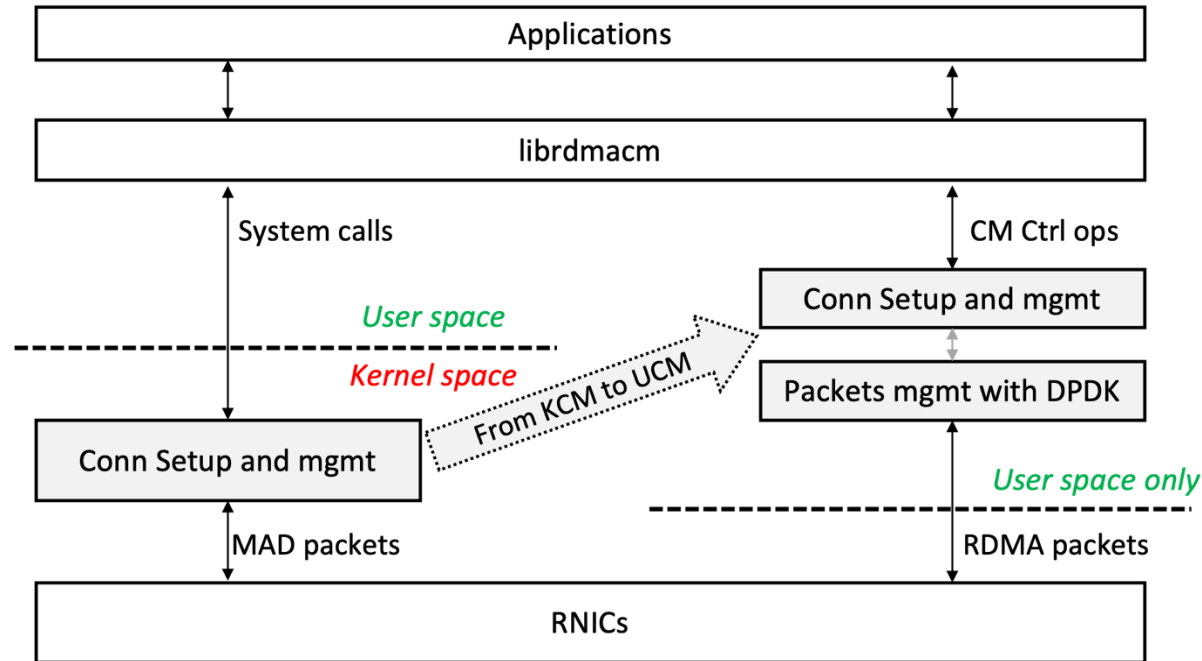
- Choosing CM conn
  - unified APIs
  - in-band path
  - path detection
- Recent User-Space Tech
  - DPDK
  - LibOS

# UCM: User-space RDMA Connection Management

- Our idea to bypass kernel for better performance



# UCM overview

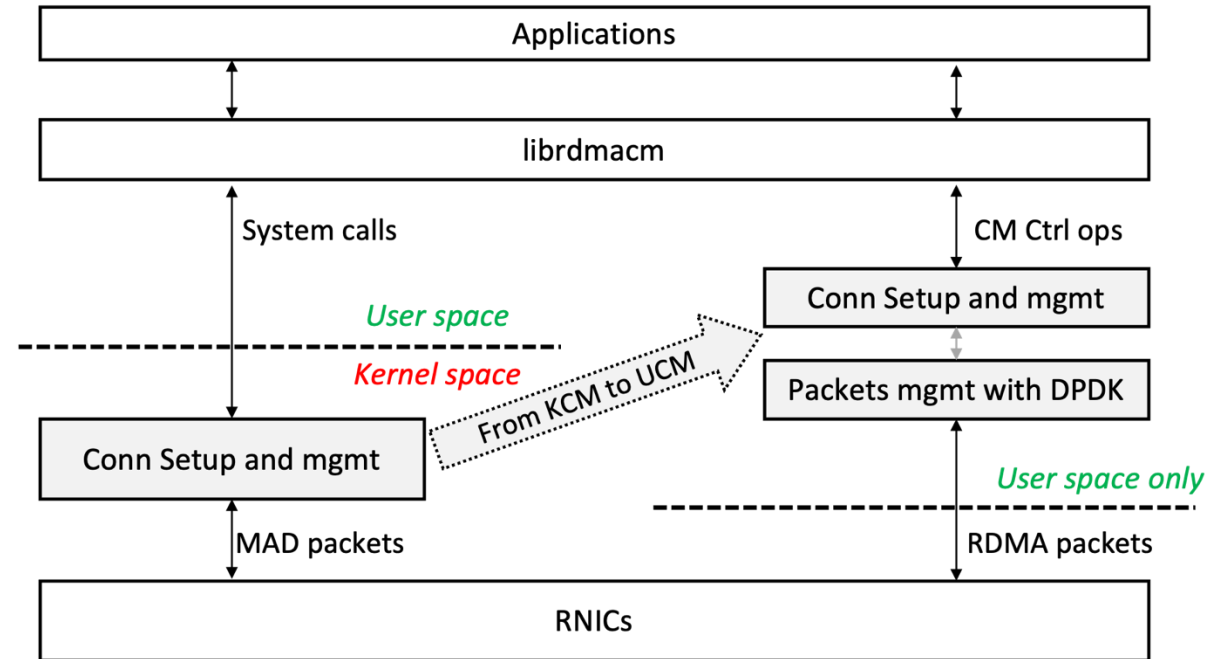


- UCM Framework

- **Work module**: setup and manage RDMA connections in user-space
- **Monitor module**: offer multi-method monitoring approaches for developers

# UCM design --- *UCM work module*

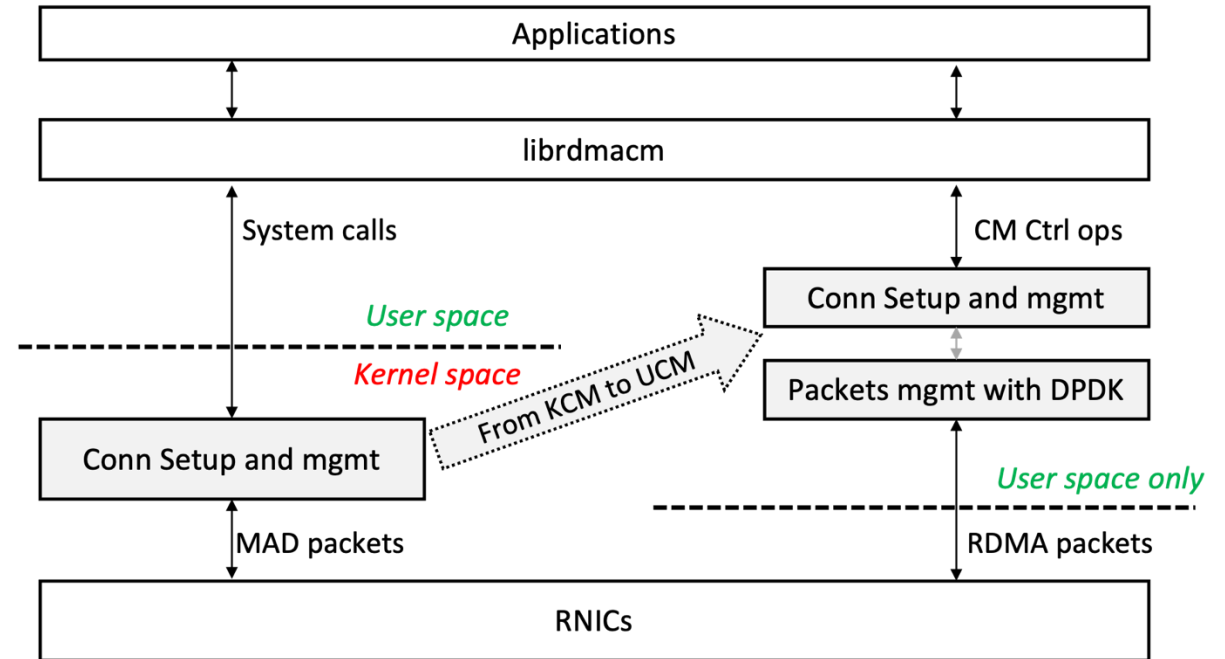
- How to bypass kernel?



- To support multi-thread, how to deal with thread communication?

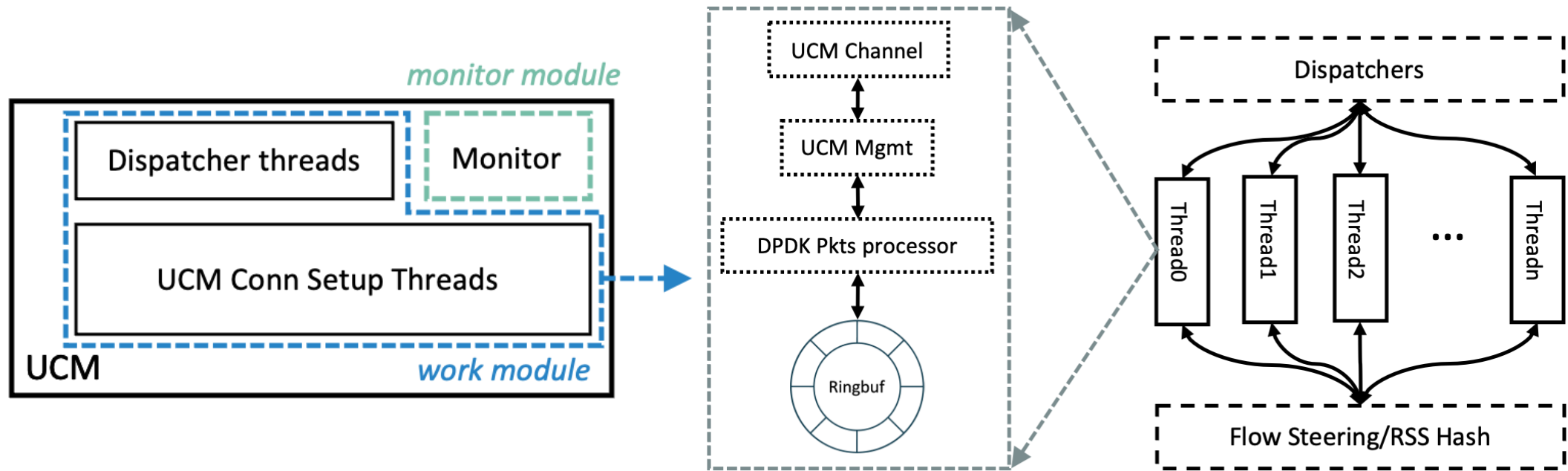
# UCM design --- *UCM work module*

- How to bypass kernel?
  - Conn Mgmt in LibOS
  - on-loading packets processing with DPDK
- To support multi-thread, how to deal with thread communication?





# UCM design --- *UCM work module*



- To support multi-thread, UCM leverages **multi-thread lock-free management**
  - NIC features : Flow Steering, RSS hash
  - Make sure that each connection's related information is accessed and managed by only one thread

# UCM design --- *UCM monitor module*

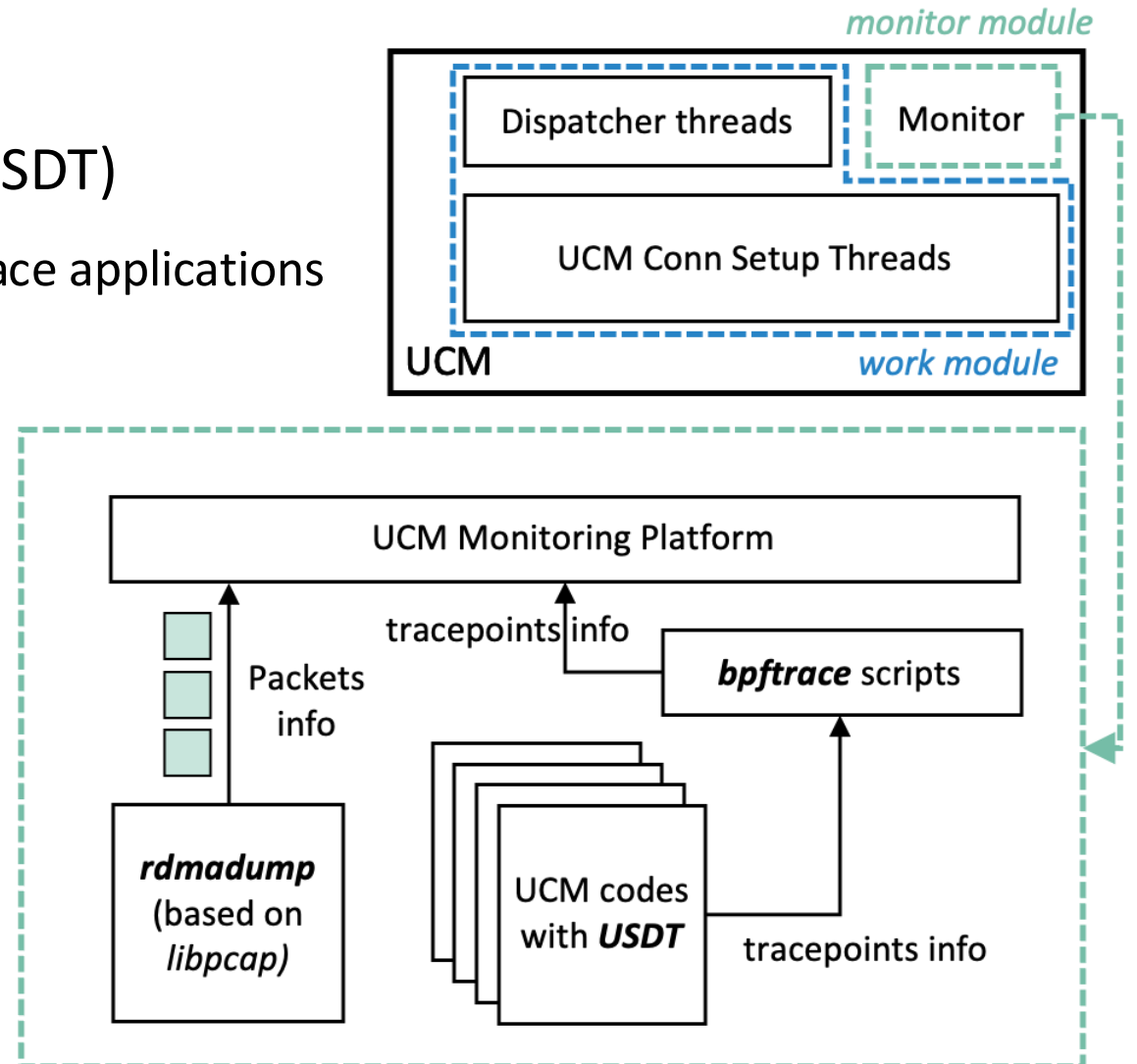
- Monitoring approaches

- User Statically-Defined Tracing (USDT)

- Add dynamic *tracepoints* to user-space applications
- Dynamic Monitoring with USDT

- Self-defined capturing tool

- Based on *libpcap*



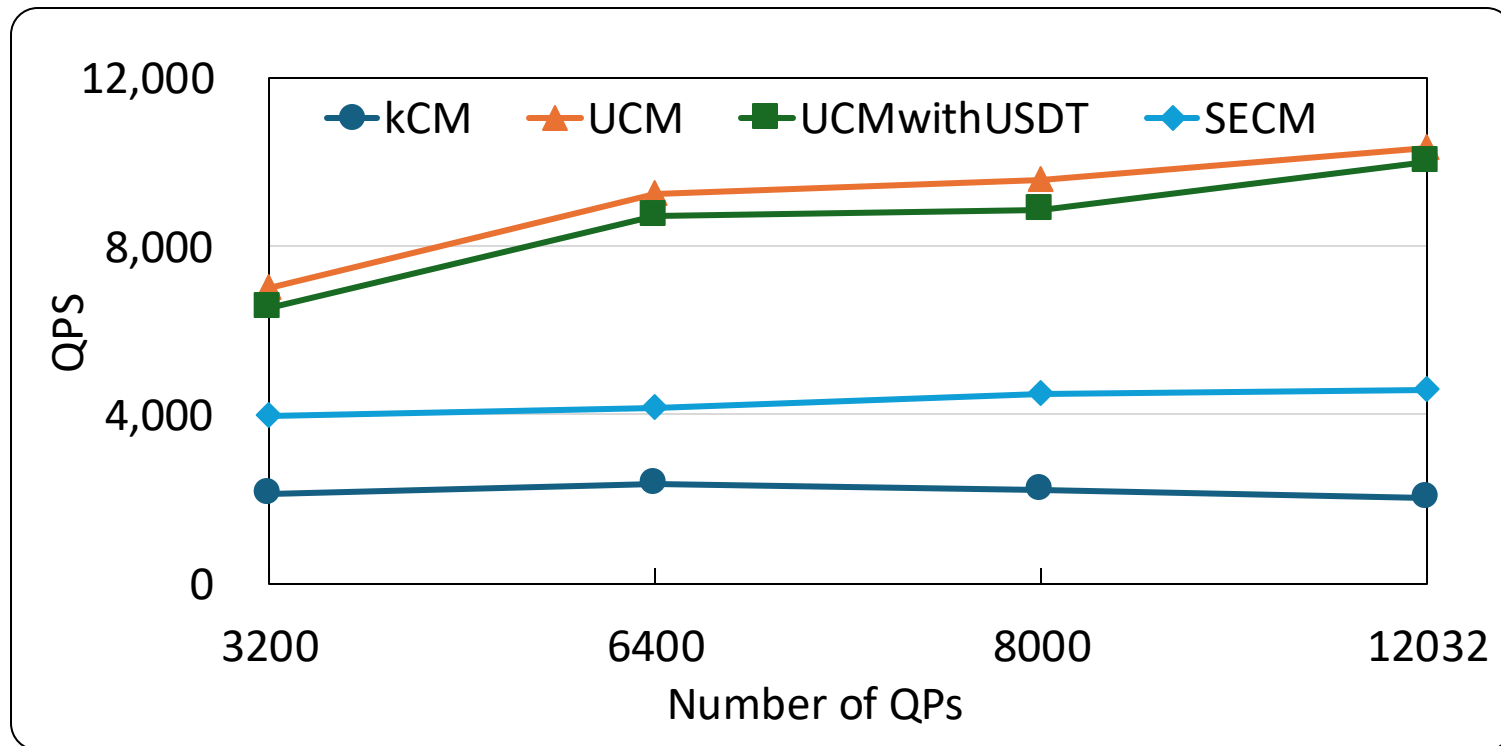
# Evaluation

---

- End-to-end Testbed
  - Hardware: Mellanox ConnectX-6 Dx EN
  - Software: cmttime, perftest, Mellanox OFED 5.8 driver,
  - Comparisons: **UCM**, kCM (original RDMA), SECM
- Goals
  - Compare UCM's **connection setup speed** under different scenarios with the *sota* approaches (single-threaded, multi-threaded, and extreme application)
  - Evaluate the impact of UCM's **maintenance overheads** on production applications
  - Showcase **successful maintenance experiences** with UCM

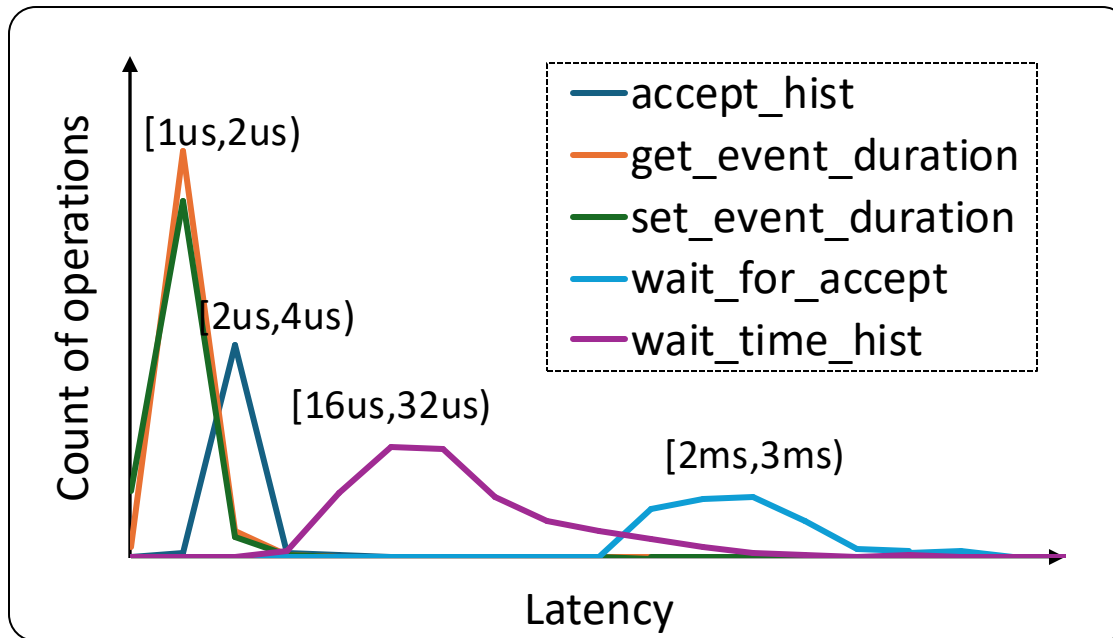
# Evaluation#1--- *UCM work performance (RPC apps)*

- Deploy UCM with an enterprise RPC framework in ByteDance
  - UCM's QPS (number of QPs per second) performs **3.3-5.1x** that of kCM and **1.8-2.2x** of SECM.
  - The **extra overhead** introduced by enabling USDT is only **3.2%-7%**.

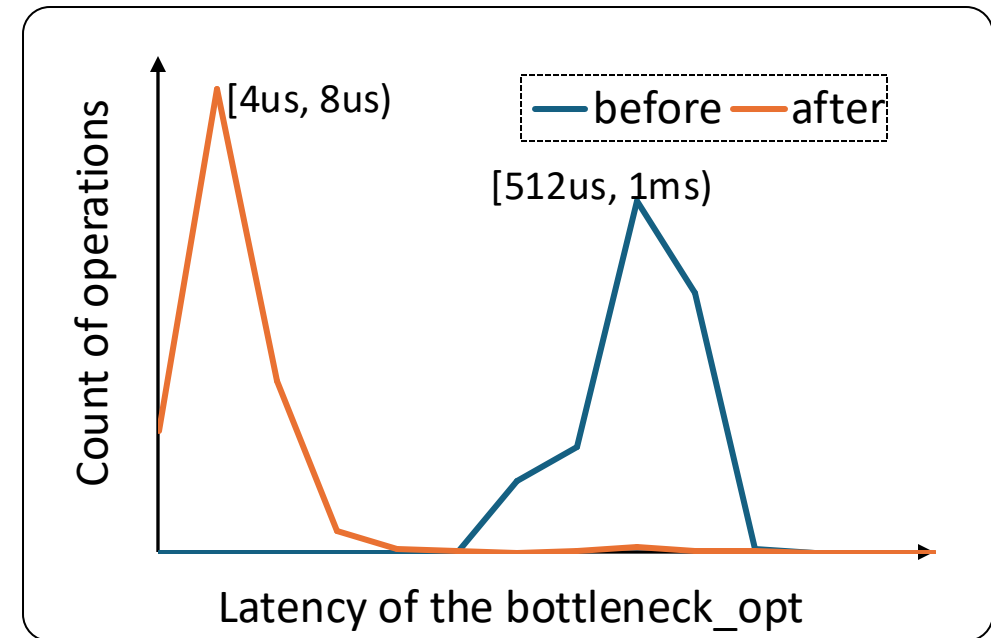


# Evaluation#2--- *UCM monitoring effort*

- Monitoring for optimization
  - latency of every operation in online connections
  - performance after optimization



Latency Monitoring for Conn Setup Steps



Latency Comparison for a Single Step

# Conclusion & Future work

---

## ➤ Related work

- KRCore<sup>1</sup> : a shared connection pool
- SECM<sup>2</sup> : parallel connection setup

Dependency	KRCore	SECM	UCM
User app	✓	✓	
RDMA library	✓	✓	✓
RNIC	✓		

## ➤ **UCM: The first pure user-space RDMA connection management framework.**

- ✓ Dramatically accelerated RDMA CM setup efficiency
- ✓ Better observability for production operations

## ➤ We hope UCM will inspire more new possibilities for optimizing the RDMA protocol stack at the software (user space) level.

[1] KRCORE: a microsecond-scale RDMA control plane for elastic computing. (ATC 2022)

[2] SECM: Securely and efficiently connections setup using RDMA-CM. Computer Networks 250 (2024)



# Thanks!

## Q&A

Contact : [shj@hnu.edu.cn](mailto:shj@hnu.edu.cn)