



PeRF

: Preemption-enabled RDMA Framework

Sugi Lee¹, Mingyu Choi¹, Ikjun Yeom², Younghoon Kim²

Acryl Inc¹, Sungkyunkwan University²

01 Intro

RDMA

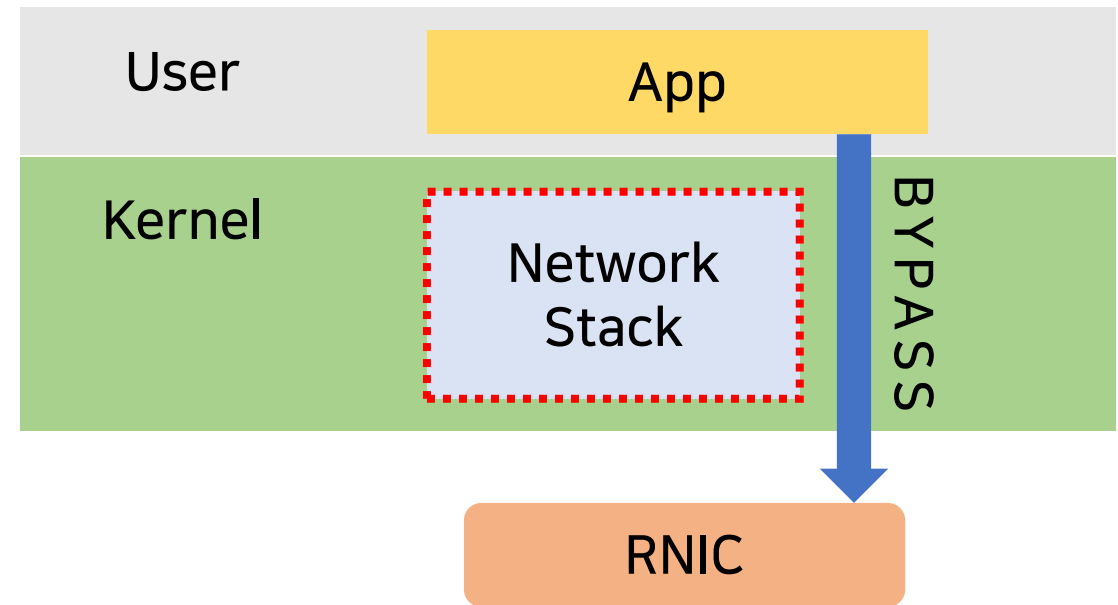
- High throughput & ultra-low latency via zero-copy operations
- Big-data analysis, machine learning, distributed storage, etc.

➔ Thanks to the **kernel-bypass** feature

At the same time...

- Problems such as security, scalability, and **performance isolation**

➔ Due to the lack of control in kernel space (**kernel-bypass**)



02 Background

RNIC's Scheduling Mechanism

Setup

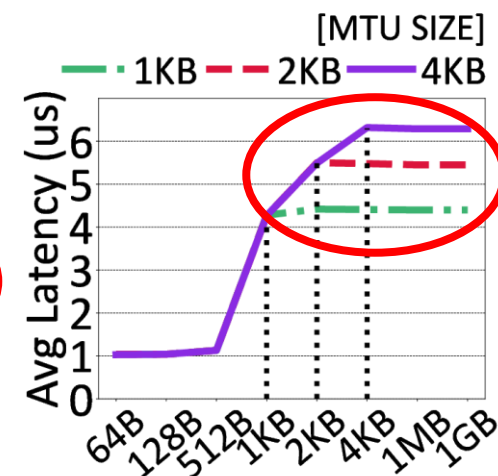
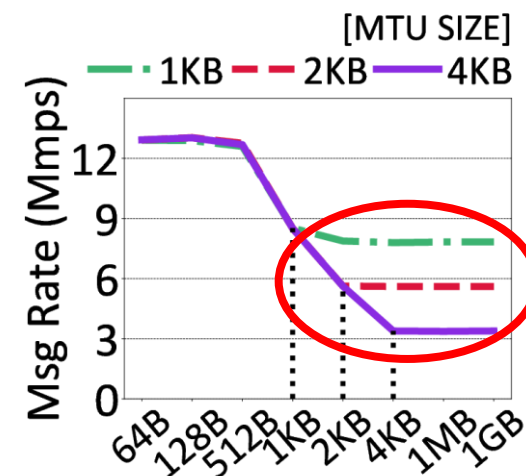
- 16-core Intel i7-11700K 3.6GHz
- 32GB RAM
- ConnectX-6 100 Gbps RoCE

Msg-level

QP-level

- 2 apps sharing an RNIC
 - One app sending **small** messages (16B)
 - The other app sending **large** messages as background traffic (64B~1GB)

→ Performance remains stable once it reaches MTU



Evenly processing MTU-segmented packets from various QPs

02 Background

RNIC's Scheduling Mechanism

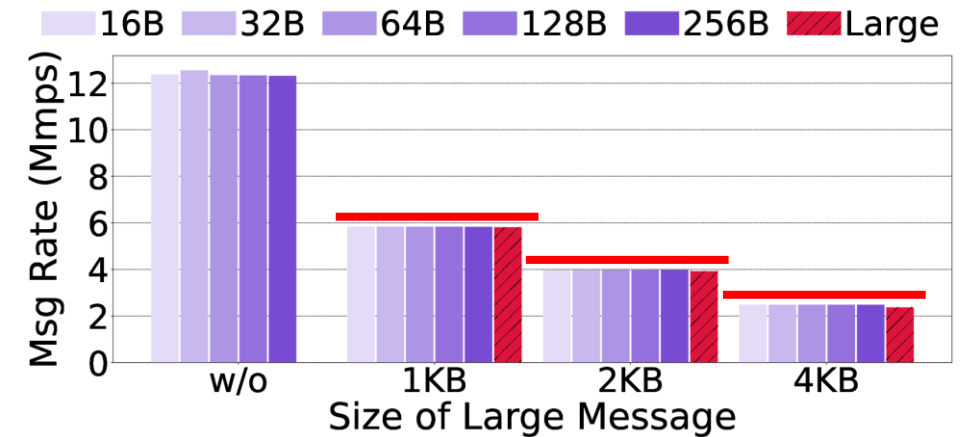
Setup

- 16-core Intel i7-11700K 3.6GHz
- 32GB RAM
- ConnectX-6 100 Gbps RoCE

Msg-level

QP-level

- 5 apps sending batches of **small messages** (16, 32, 64, 128 and 256B each) with single QP
- Adding an app sending **larger messages** (1, 2, or 4KB) with single QP



→ Larger messages limits the msg rate of other apps (bandwidth bottleneck)

Round-Robin Scheduling with Multi-QPs

02 Background

Application Type

Msg-level

	Bandwidth-intensive	Message-intensive	Delay-sensitive
Message Size	Large	Small	Small
Sparsity	Low	High	Low

QP-level



Single QP	B_App _{single}	M_App _{single}	D_App _{single}
Multiple QPs	B_App _{multi}	M_App _{multi}	D_App _{multi}

In this paper,

B_App



M_App

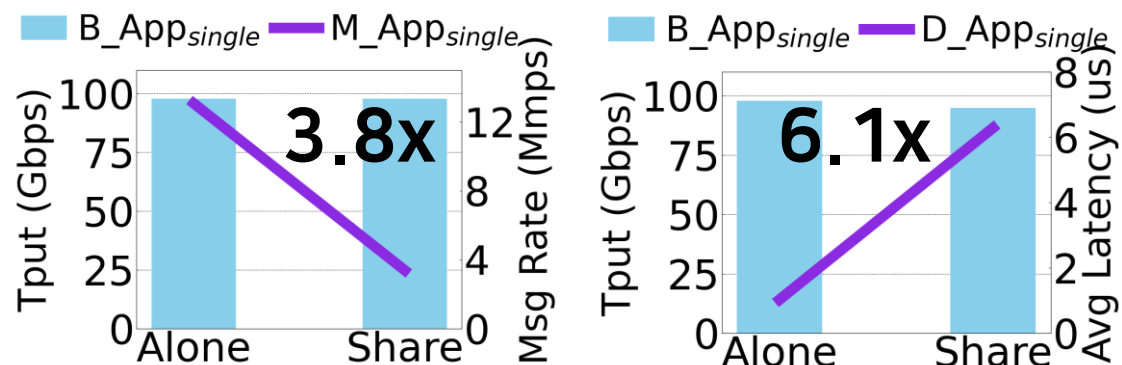


D_App



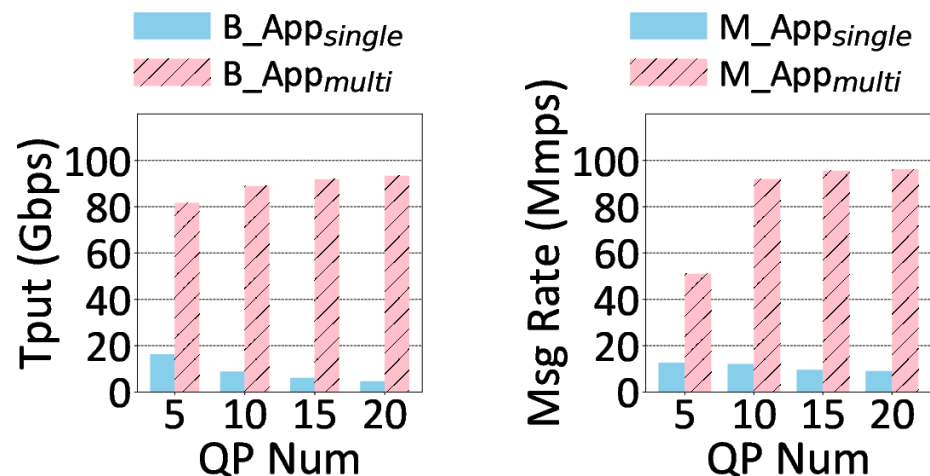
02 Background

Message-level Anomaly



- B_App harms both M_App and D_App.

QP-level Anomaly



- Multi-QP app harms single-QP app.

RNIC does not support Multi-Tenancy

Prior Works for Performance Isolation

	HW-based	SW-based
Pros	✓ Strict performance isolation	✓ Flexible performance isolation
Cons	- Unable to handle dynamic changes	- Reservation-based resource allocation (non-work-conserving) - Unavoidable performance degradation

We propose, **PeRF**

A Preemption-enabled RDMA Framework

- ✓ software-based (flexible)
- ✓ bare-metal performance RDMA (work-conserving)
- ✓ no need for estimation of network resources
- ✓ transparent to applications

03 Key Idea

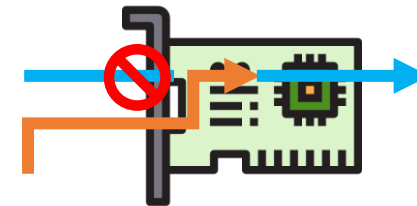
Preemption

preemption : the ability of an operating system to temporarily **interrupt** a currently running task to run another task.

Preemption in RDMA?

“Interrupting a message and prioritizing another in RNIC”

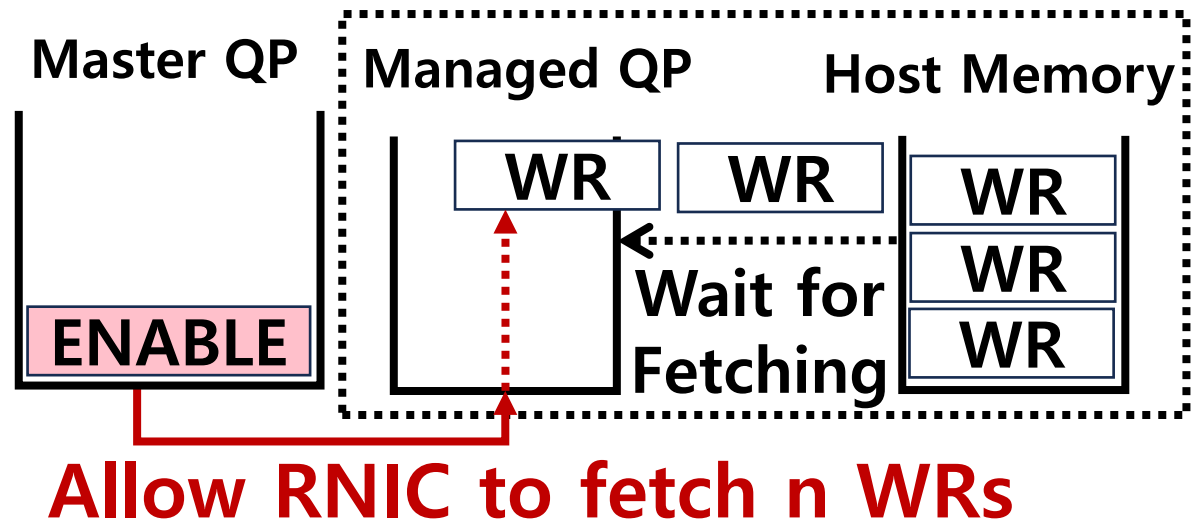
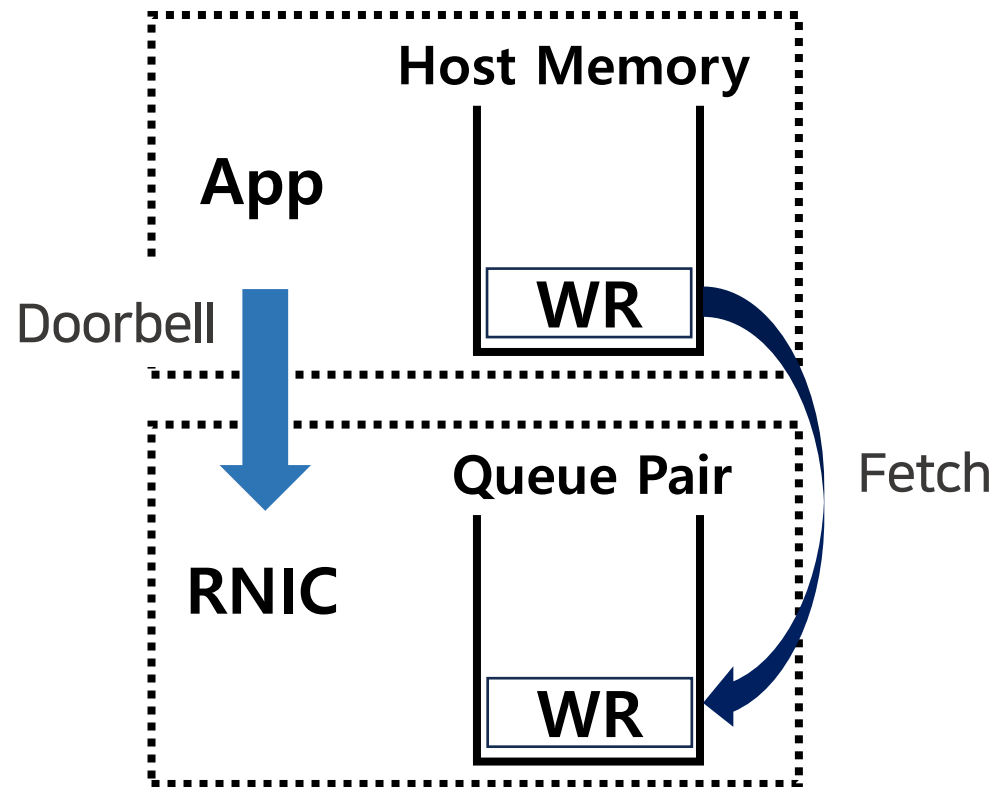
by leveraging a **managed QP** and a combination of **specialized WRs**



03 Key Idea

Preemption Mechanism

Managed QP and ENABLE_WR

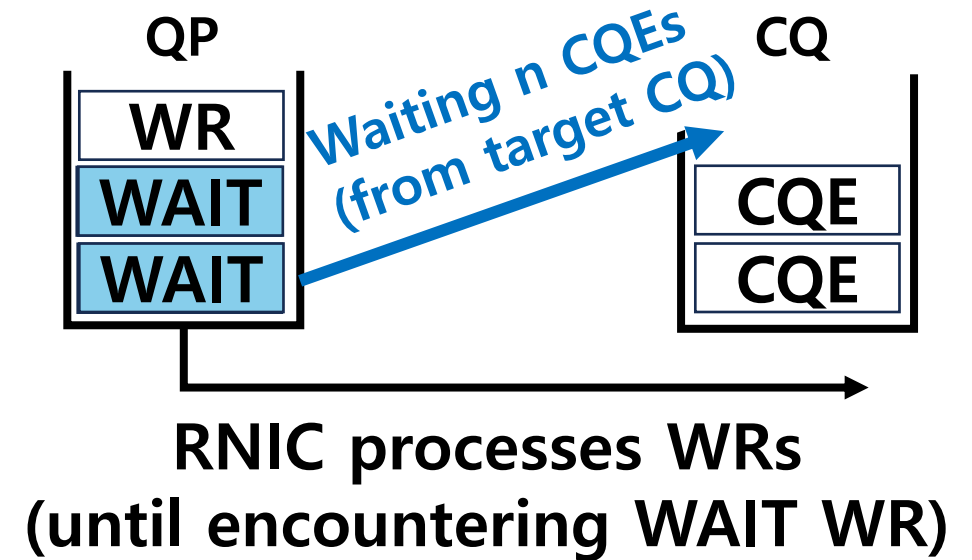
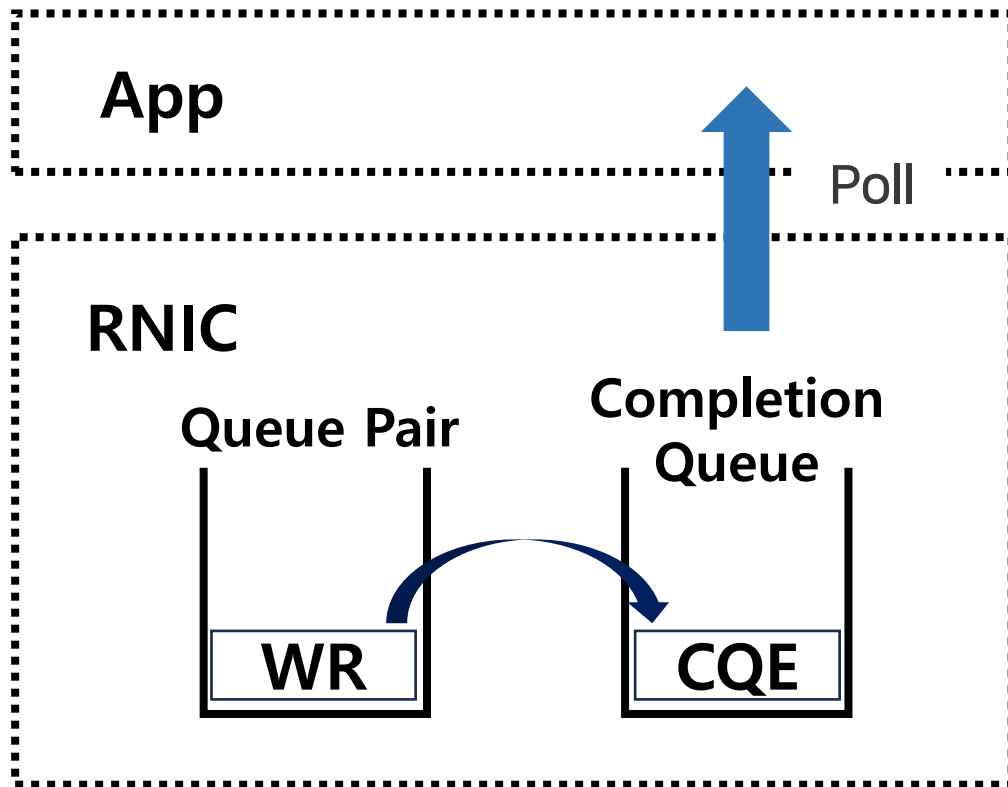


ENABLE_WR

03 Key Idea

Preemption Mechanism

CQE and WAIT_WR



WAIT_WR

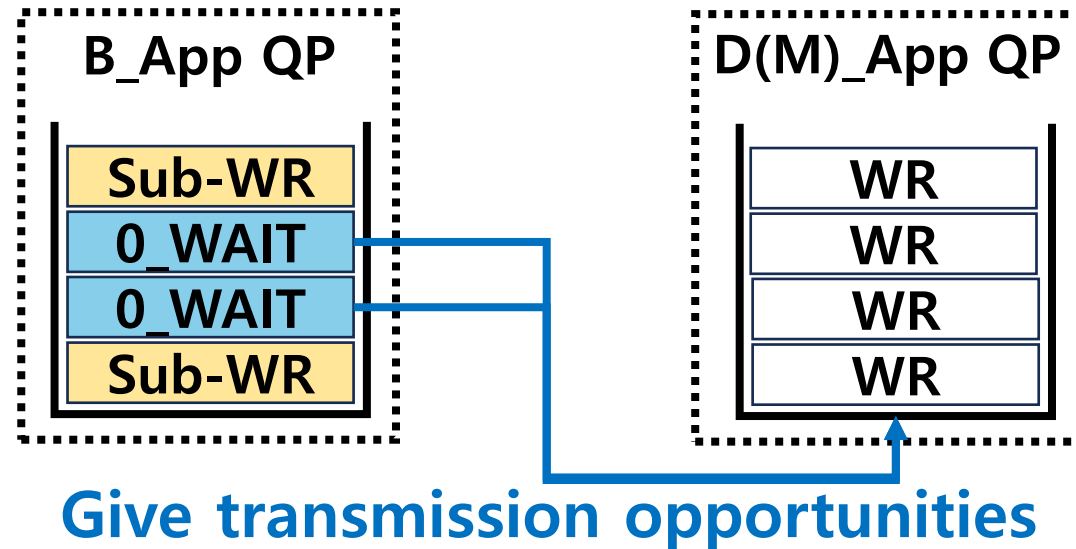
03 Key Idea

Preemption Mechanism (Message-level)

Msg-level

QP-level

- Sub-WR : subdivided unit of a larger work request
- O_WAIT : a WAIT_WR that immediately completes upon execution (wait_cqe_num=0)

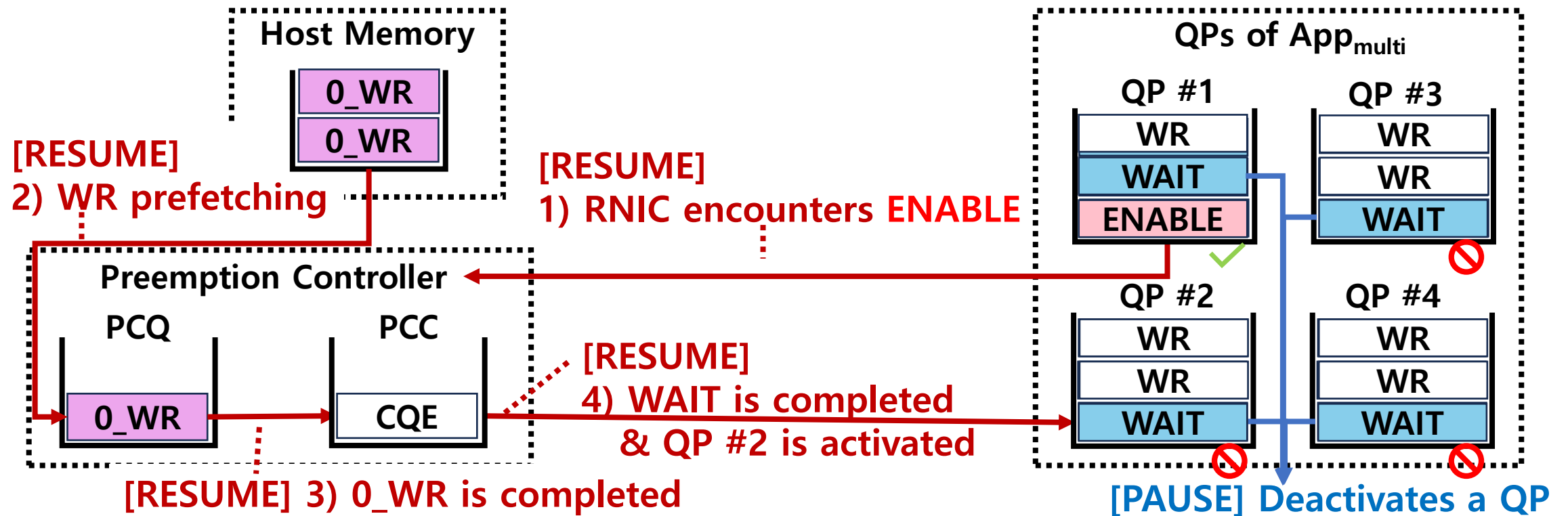


Transmission Interrupt

03 Key Idea

Preemption Mechanism (QP-level)

- 0_WR : a work request with 0 byte message



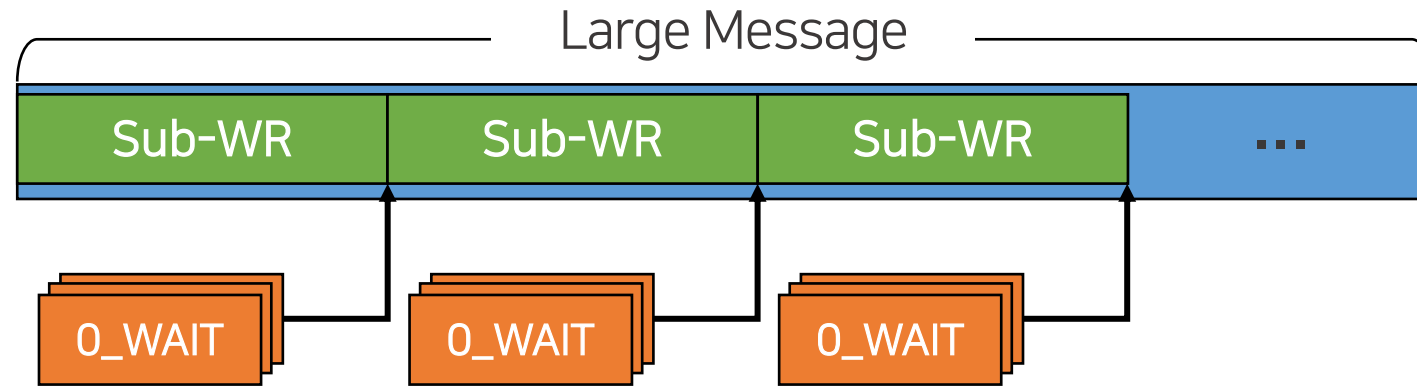
PAUSE/RESUME Operations

Message-level Scheduling

Msg-level

QP-level

Large Message Scheduling Engine (LMSE)



- D(M)_App is competing with B_App
- Split B_App's message into **smaller chunks (Sub-WRs)**
- Post ***O_WAIT*** WRs between message chunks

Grant more opportunities for small message flows

04 Design

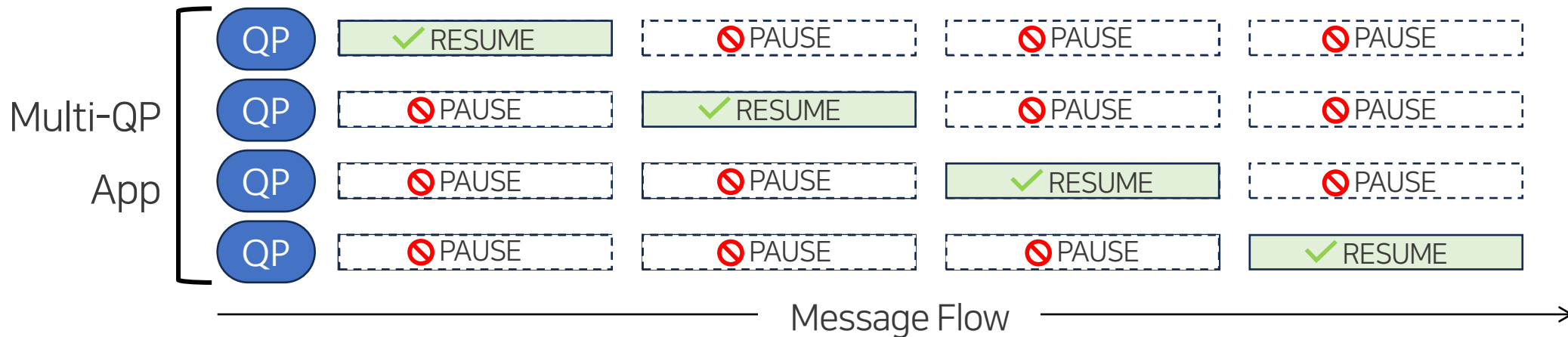
QP-level Scheduling

Msg-level

QP-level

Multi-QP Scheduling Engine (MQSE)

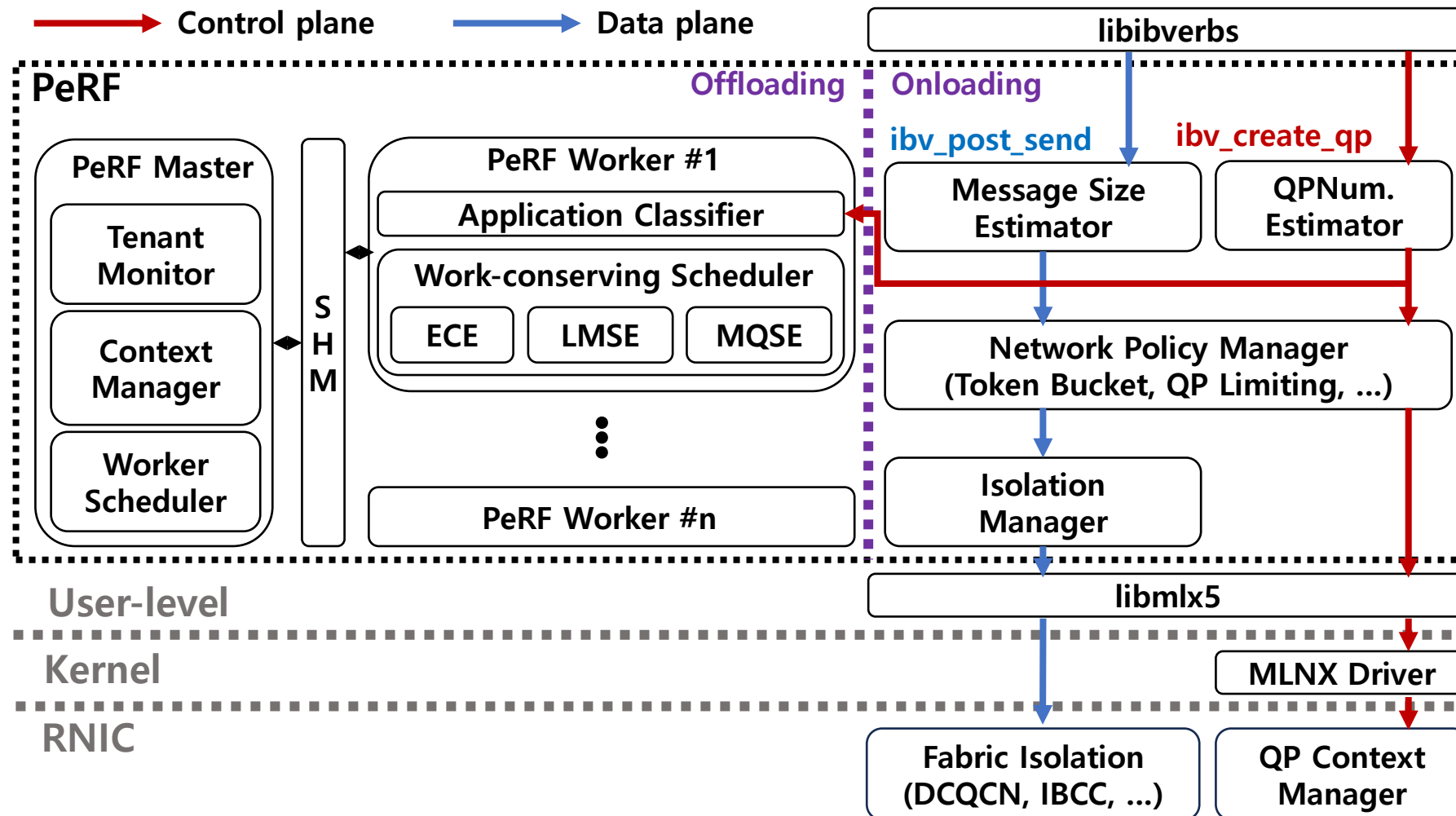
- When competing with a multi-QP App
- Control the number of active QPs with **PAUSE/RESUME** operations



Restrict the number of activated QPs of multi-QP app

04 Design

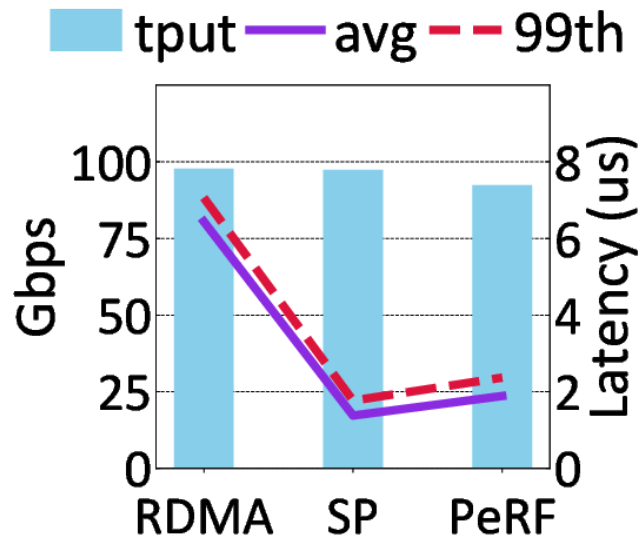
PeRF Architecture



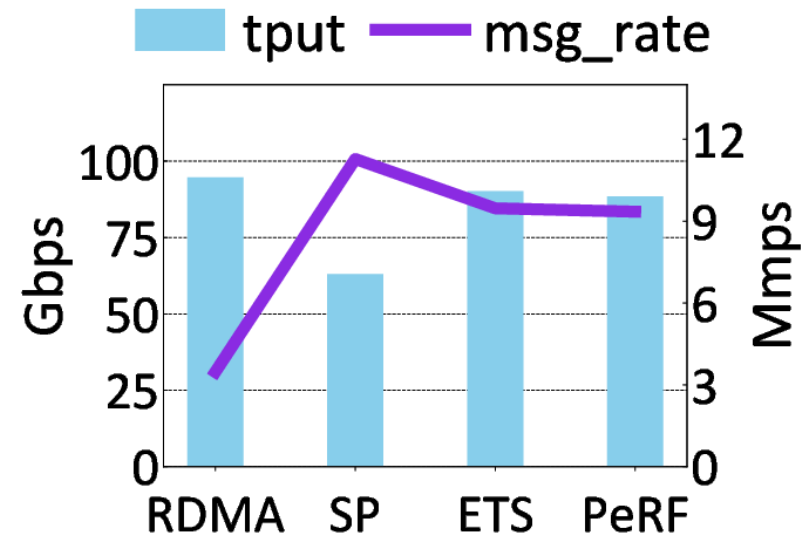
05 Evaluation

Baseline Benchmark

- Strict Policy (SP) : priority queue (set to prioritize D_App)
- Enhanced Transmission Selection (ETS) : weighted round-robin (B_App : M_App = 16 : 1)



B_App_{single} vs D_App_{single}



B_App_{single} vs M_App_{single}

PeRF performs nearly as HW-based solutions

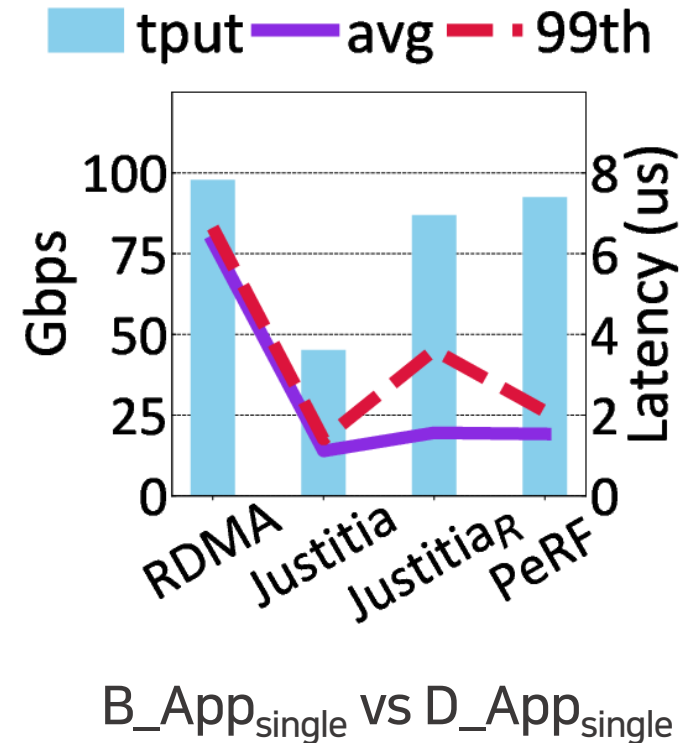
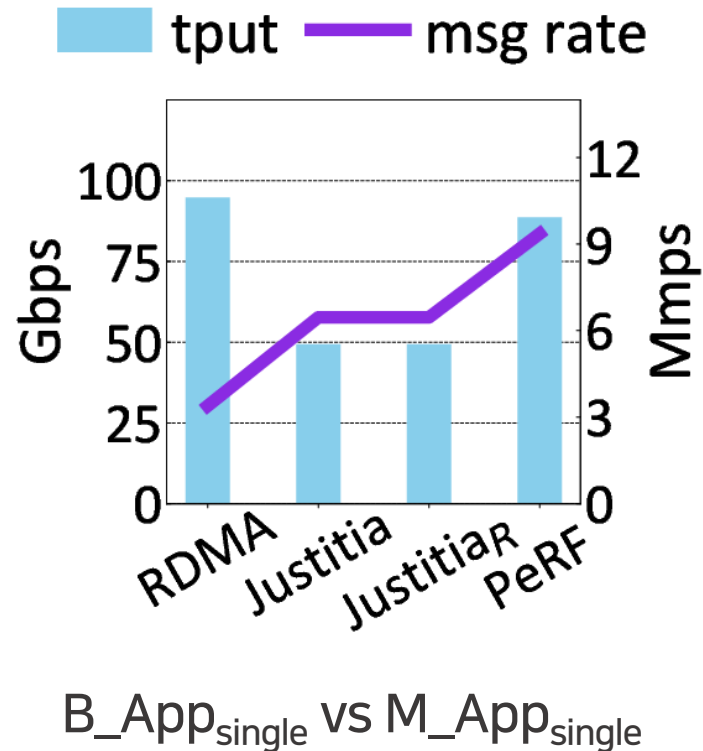
05 Evaluation

Message-level Isolation

Msg-level

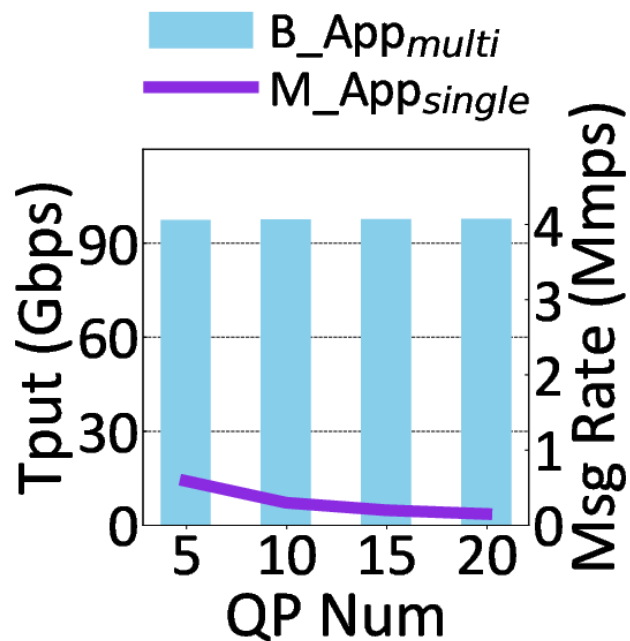
QP-level

- Justitia : a token-based resource allocation solution
- Justitia_R : a relaxed version of Justitia (latency threshold = 15 us)

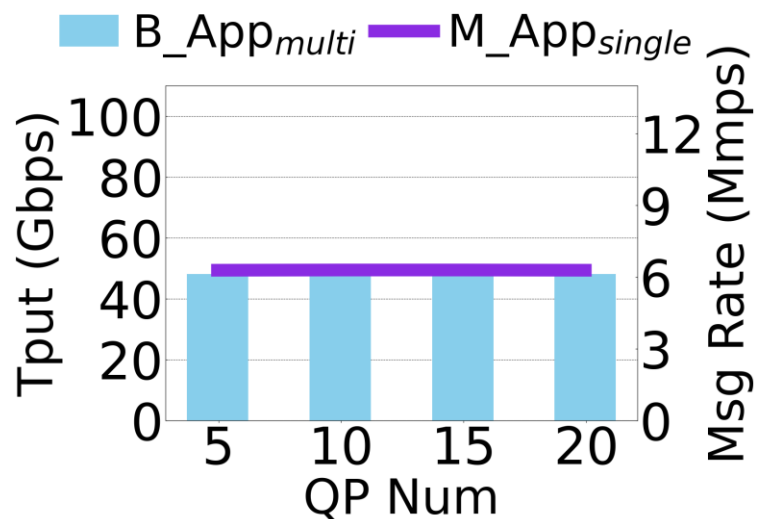


05 Evaluation

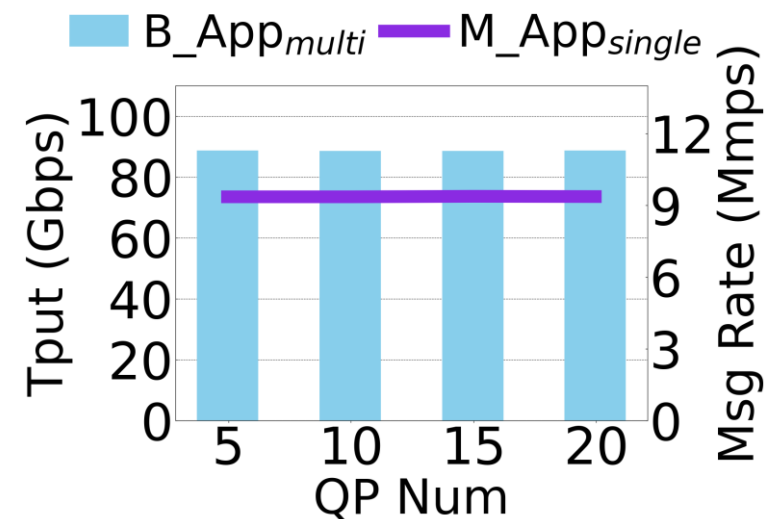
QP-level Isolation



RDMA



Justitia



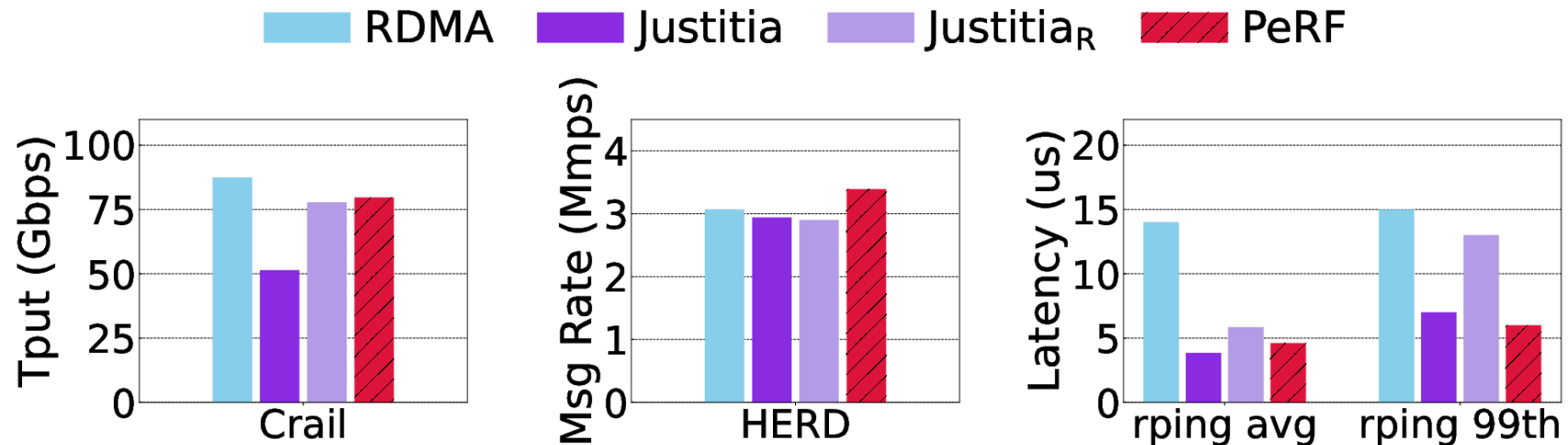
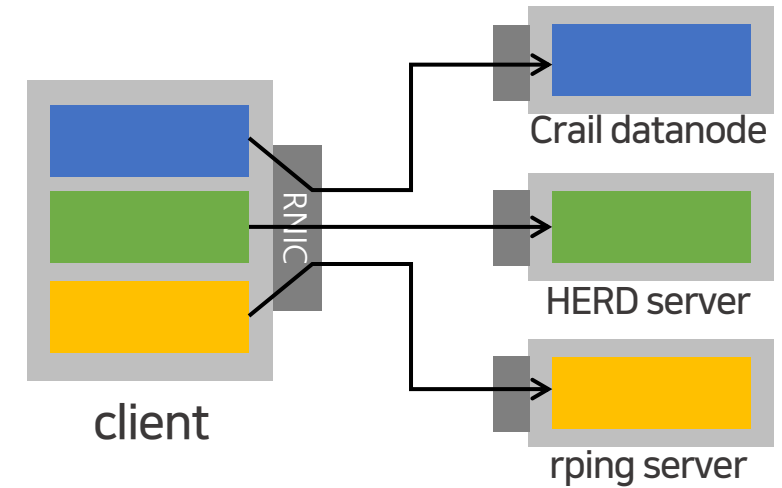
PeRF

Msg-level QP-level

05 Evaluation

Real World Applications

Application	Type	Message Size	Batch
Apache Crail	B_App	1GB (1MB chunk)	1
HERD	M_App	5% PUT 50B 95% GET 17B	32
rping	D_App	16B	1



Other Experiments

1. Other Commercial RNICs Test
2. Support for SEND and READ Operations
3. Scalability Test
4. Weighted Policy Test
5. Congested Networks Test

... and more!

06 Conclusion

PeRF

- ✓ Traditional RDMA solutions struggle with **performance anomalies** and **inefficient resource utilization** in **multi-tenant** environments.
- ✓ PeRF uses a novel RNIC **preemption mechanism** and **work-conserving packet/QP scheduling** to address these challenges.
- ✓ PeRF ensures **software-based performance isolation** without sacrificing the high performance of RDMA.

Thanks & Q/A

github: <https://github.com/acryl-aaai/perf>
contact: ① sglee0323@gmail.com
② mingyuchoi514@gmail.com