

# Receiver-Driven Congestion Control for InfiniBand

Yiran Zhang, Kun Qian, Fengyuan Ren Tsinghua University

## InfiniBand Interconnect

- Lossless and low-latency communication
- Credit-based link-level flow control to guarrantee losslessness
- The side-effect: congestion spreading
  - Head-of-line blocking: waiting for credits

• Increasing queueing delay and damaging the completion time of applications





## **IB** Congestion Control in IB Specification

Switch behavior

 Congestion detection and identification: exceeding the selected threshold & there are available credits -> a root of congestion -> FECN notification

• Channel Adapter (CA) behavior

• Rate regulation: On receiving CNP, regulates the sending rate with a fixed moving step according to a CCT table -> multiple round-trip times to converge





## Traffic in HPC Systems

- Applications: Communication-intensive, I/O-intensive → short-lived MPI messsages long-lived I/O traffic
- Shared network infrastracture and file systems →
  increasingly diverse traffic patterns

Coexisting diverse traffic imposes great challenges to congestion control





## **Experimental observations**

- Simulations
- Mixture of short messages (F3-F17) and long-lived messages (F1 F2)
- Congestion port: PO



Burst last for about 1ms





50th International Conference on Parallel Processing (ICPP) August 9-12, 2021 in Virtual Chicago, IL

## Experimental observations

• Overall throughput performance



**S2** 

P2

**F1** 

**F2** 

**H1** 

**H2** 

**S1** 

Burst

**R1** 

**R2** 

F17



## Insights (1)

• Periodical updated credits may confuse congestion detection.



(a) Available credits at port P2





## Insights (2)

- Periodical updated credits may confuse congestion detection.
- The sluggish rate adjustment of end-to-end IB CC mismatches the fast hop-by-hop credit-based flow control.



## **Insights Summary**

- Periodical updated credits may confuse congestion detection.
- The sluggish rate adjustment of end-to-end IB CC mismatches the fast hop-by-hop credit-based flow control.
- Congestion-unaware rate increase contradicts with rate decrease: parameter tuning is tough





### **Design: Overview**



RR CC





50th International Conference on Parallel Processing (ICPP) August 9-12, 2021 in Virtual Chicago, IL

## Key idea (1): Receiver-Driven Congestion Identification

- The receiver-side has the packets marking pattern to guide congestion identification:
  - Congested flows: *all* packets passing through the congestion root and are marked with FECN
  - Victim flows: only *fractional* packets have available credits to send and may be marked with FECN.





## Key idea (2): Receiver-Driven Rate Regulation

• The receiver-side has the *receiving rate* information to guide rate regulation:

The *receiving rate* is the maximum rate that the congested flow can pass through the bottleneck link.



- Endpoint congestion
- In-network congestion
- o Coexisting multiple congestion points



### Design: Receiver – Congestion Identification and Notification

Algorithm 1 Receive side CA behavior

- 1: // for each Queue Pair or Service Level
- 2: while during each period *T* do
- 3: **if** received packet with FECN bit set **then**
- 4:  $fecn\_count \leftarrow fecn\_count + 1$
- 5:  $recv\_count \leftarrow recv\_count + 1$
- 6: else
- 7:  $recv\_count \leftarrow recv\_count + 1$
- 8: end if
- 9: end while
- 10: **if** period T expires **then**
- 11:  $P_{FECN} \leftarrow fecn\_count/recv\_count$
- 12:  $receive\_rate \leftarrow recv\_count * MTU/T$
- 13:  $recv\_count \leftarrow 0$
- 14:  $fecn\_count \leftarrow 0$
- 15: **if**  $P_{FECN} \ge 0.95$  **then**
- 16: Send CNP with BECN bit set and *receive\_rate*
- 17: **else**
- 18: **if** *receive\_rate < line\_rate* **then**
- 19: Send CNP without BECN bit set
- 20: **end if**
- 21: **end if**

CONFER 22: end if

PROCES

INTL



#### Receiving rate calculation

#### Generating Congestion Notification Packets to guide rate increase/decrease explicitly (BECN bit)

ce on Parallel Processing (ICPP) in Virtual Chicago, IL



## Design: Sender – Rate Regulation

Algorithm 2 Send side CA behavior

- 1: // for each Queue Pair or Service Level
- 2: if received CNP with BECN bit set then
- 3:  $target_rate \leftarrow current_rate$
- 4:  $current\_rate \leftarrow receive\_rate * d$
- 5: **else**
- 6: if consecutive 3 CNPs without BECN bit set then
- 7:  $target_rate \leftarrow target_rate + R_{AI}$
- 8:  $current\_rate \leftarrow (current\_rate + target\_rate)/2$
- 9: **else**
- 10:  $current\_rate \leftarrow (current\_rate + target\_rate)/2$
- 11: **end if**
- 12: **end if**



Simple yet effective rate adjustment rules for congested flows and uncongested/victim flows:



lacksquare

Rate decrease according to piggybacked receiving rate







• Validation in the typical scenario (1)



#### Fast congestion elimination & short convergence time



50th International Conference on Parallel Processing (ICPP) August 9-12, 2021 in Virtual Chicago, IL



**S2** 

P

**H**3

**F1** 

**F2** 

**H**1

**H2** 

**S1** 

Burst

H4

**R2** 

F17

(H17

(H16)

• Validation in the typical scenario (2)







50th International Conference on Parallel Processing (ICPP)

August 9-12, 2021 in Virtual Chicago, IL

 $H1 \xrightarrow{F2} P2 \xrightarrow{P1} P1 \xrightarrow{P0} F17$   $F3 \xrightarrow{r} F17$   $H3 \xrightarrow{H4} \dots \xrightarrow{H16} H17$ 

**F1** 

**S2** 

**S1** 

**R1** 

**R2** 

• Large-scale simluations

1024-node Fat-tree (QDR rate), mixed MPI and I/O messages generated by MPI servers and I/O servers.





50th International Conference on Parallel Processing (ICPP)

In-Cooperation

August 9-12, 2021 in Virtual Chicago, IL

• Other properties



Fairness

Overhead of CNP



50th International Conference on Parallel Processing (ICPP)

August 9-12, 2021 in Virtual Chicago, IL



## Conclusion

- Investigating IB CC performance under typical mixed traffic scenario.
- Proposing RR CC for InfiniBand network, with much better average messages latency, fewer parameter tuning and compatible with commodity InfiniBand switches.

## Thank you!

zyr17@mails.tsinghua.edu.cn







50th International Conference on Parallel Processing (ICPP) August 9-12, 2021 in Virtual Chicago, IL

