

Storage Developer Conference September 22-23, 2020

NVMe-oF on RDMA Performance Challenges and Solutions in Commodity Servers

Yamin Friedman & Rob Davis Nvidia networking business unit

#### SSDs Create a Storage Networking Performance Bottleneck



#### SD @

#### 200Gb Network Speeds and NVMe-oF on SD@ RoCE to the Rescue



# What is NVMe over Fabrics (NVMe-oF)?

- Extending NVMe protocol over a fabric
  - NVMe commands and data structures are transferred end to end
  - Enables SAN features for NVMe based storage systems
- Bypassing legacy stacks for performance
- First products and early demos all used RDMA
  - Standard supports multiple transports
- Impressive performance gains



https://www.theregister.co.uk/2018/08/16/pavilion\_fabrics\_performance/

2020 Storage Developer Conference. © Nvidia Networking. All Rights Reserved.



#### SD@

#### What is RDMA?



12.00 10.00 100GbE Throughput (GB/sec) 8.00 RoCE 6.00 4.00 2.00 0.00 No RDMA **RDMA** With RDMA **2x Better Bandwidth** Half the Latency

33% Lower CPU

See MS demo: <u>https://www.youtube.com/watch?v=u8ZYhUjSUol</u>

Microsoft Storage Spaces Throughput

SD@



#### Applications That Need Even More Performance



2020 Storage Developer Conference. © Nvidia Networking. All Rights Reserved.

#### SD@

#### System interrupt overload

### System interrupts and RDMA

20

- In order to handle networking traffic the kernel uses interrupts
- In NVMe-oF storage flows there are multiple interrupts per IO
- When the system is overloaded with interrupts it interferes with the ability to perform RDMA

#### **Test setup**

- CPU: Intel(R) Xeon(R) Platinum 8176M CPU @ 2.10GHz 28 cores
- Linux version: 5.7.0-rc3 with added shared CQs
- Network adapter: Mellanox ConnectX-5 Ex

#### Interrupt overload

SD@



### **Dynamic Interrupt Moderation (DIM)**

#### **Interrupt Moderation**

 Interrupt moderation is a mechanism for aggregating more work per system interrupt

 RDMA completion queues can be configured to arm a system interrupt only once enough work has been gathered or a timeout reached

 There is no single configuration that is optimal for a specific system and time

# **Dynamic interrupt moderation**

20

 In order to effectively use interrupt moderation we need a method that adapts to the current state

 DIM monitors the current state of the system and determines the optimal parameters moment to moment

# **DIM algorithm**

- Statistics collected:
  - Number of completions polled
  - Number of interrupts
- Three stage algorithm:
  - Optimize for number of completions
  - Optimize for ratio of completions to interrupts
  - Reduce moderation if ratio is low

#### **4K read comparison**



2020 Storage Developer Conference. © Nvidia Networking. All Rights Reserved.

#### SD@

#### **4K write comparison**



4K read latency

SD@

#### **Shared Completion Queues**

# Interrupt moderation across applications

- As presented DIM is very effective in reducing interrupts per CQ
- What happens when there are many applications each with their own set of CQs?

#### **Shared CQs**

- In the Linux kernel driver CQ interrupt handling is provided as a service for all applications
- In this model each application opening its own set of CQs is an inefficient use of system resources
- With shared CQs we provide an API that provides full functionality while reducing system overhead

#### 4K read 4 disk comparison



4K read latency using 4 disks

SD<sub>20</sub>



#### 4K read 4 disk comparison

#### 2500 3000000 3000 3000000 2500000 Puo 2500 2500000 2000 second Better sec 2000000 2000 2000000 ettei KIOPs nterrupts per 1500 Interrupts per 1500000 usec 1500 1500000 Ň 1000000 1000 1000 1000000 500000 500 500 500000 0 21 23 25 27 17 19 15 Number of cores 11 13 15 17 19 21 23 25 27 g Number of cores With DIM interrupts With DIM interrupts With DIM and shared CQs interrupts With DIM and shared CQs interrupts - With DIM KIOPS With DIM average latency With DIM and shared CQs KIOPS

4K write latency using 4 disks

SD<sub>20</sub>

2020 Storage Developer Conference. © Nvidia Networking. All Rights Reserved.

4K writes KIOPS using 4 disks

#### **Concluding remarks**

#### Conclusions

- Discovered hardware limitation on commodity servers for handling interrupts during RDMA
- Provided two features to upstream Linux to overcome this limitation

#### Key takeaways

- While optimizing our storage solutions, bottlenecks may be found throughout the hardware and software stack
- Solutions may need to be very flexible and perhaps cross abstraction layers

#### **Contact information**

- Yamin Friedman: <u>yaminf@nvidia.com</u>
- Rob Davis: <u>rdavis@nvidia.com</u>

# Please take a moment to rate this session.

Your feedback matters to us.