

### 2021 OFA Virtual Workshop

## **sRDMA: SECURE TRANSPORT FOR REMOTE DIRECT MEMORY ACCESS**

Konstantin Taranov

**ETH Zurich** 



### **RDMA IS A TREND IN PUBLIC CLOUDS**

#### Alibaba Builds High-Speed RDMA Network for AI and Scientific Computing

Alibaba Clouder June 3, 2019 💿 19,0

9 💿 19,036 🖾 0

Alibaba has built the RDMA high-speed network within its global and ultra-large data centers to support Al and scientific computing.

#### 200 Gigabit HDR InfiniBand Boosts Microsoft Azure High-Performance Computing Cloud Instances

The HDR InfiniBand Connected Virtual Machines Deliver Leadership-Class Performance, Scalability, and Cost Efficiency for a Variety of Real-World HPC Applications

November 18, 2019 12:00 PM Eastern Standard Time

CRN CRN

Oracle's Latest Exadata Machines Available As Public Cloud ...

The X8M introduces an RDMA networking fabric that clocks at 100 Gb/sec (the previous generation was 40 Gb/sec) and up to 96 terabytes of ... 15.10.2020

#### How RDMA Became the Fuel for Fast Networks

The networking technology that feeds the world's largest supercomputers and data centers spawned one of this year's top tech mergers and now it drives AI.

April 29, 2020 by RICK MERRITT

### **RDMA IS A TREND IN TOP SYSTEMS CONFERENCES**

		DrTM+H'	18	DSLR'18	NAM-DB'1	.7	
	FaRM'14	Wukong'	16	HERD'14	RAMCloud	'15	CoRM'21
Octopus'17	FileMR'14	FaSST'1	16	RDMP-KV'20	)		Catfish'19
XSTORE'20					ccNUMA'	18	
Hermes'20	нуйгар	B 12	R		Grappa'	15	Derecho'19
A1'20	DrTM+R'1	6 Cra	ail'19	SparkR	XDMA'14	C-	Hint'14
Scale	eRPC'19	Dare'15		Storm'19	DaRPC'14	DrT№	1'15
	TH-DPMS	'16 ŀ	lyper	loop'18	APUS'17		

designed for performance – lower latency, higher bandwidth, lower CPU utilization etc.

### **IS INFINIBAND SECURE?**



Rothenberger et al.: ReDMArk: Bypassing RDMA Security Mechanisms, Usenix Security'21

### FEATURES OF SECURE TRANSPORT

#### Integrity

• verifies that the data has not been forged or tampered with

#### Authentication

• ensures that the parties exchanging information are who they claim to be

#### Encryption

hides the data being transferred from third parties.

#### Replay attack protection

• Each unique packet is processed only one time.

### **INFINIBAND SECURITY CONSIDERATIONS**

RFC4297 — Remote Direct Memory Access (RDMA) over IP Problem Statement: "The RDMA protocols must permit integration with Internet security standards, such as IPsec and TLS."						
Manhee Lee et al., Security Enhancement in InfiniBand Architecture (IPDPS'05): "Replacing 32 bit checksum with a message authentication code"	April 2005					
<ul> <li>Yair Ifergan (Mellanox) 3 years ago</li> <li>Hi Nikhil,</li> <li>Based on a conversation I had with the Mellanox team, encryption for native IB, nor for RoCEv1 is not a standard (yet); None of the organizations chose to pick up the challenge yet. Typically, Mellanox promotes and involve in implementing standard protocols for bringing the best out of user's requirements and the community collaboration. Hope this helps.</li> </ul>	March 2017					
IPsec over RoCE "Offered by Mellanox ConnectX-6 DX network adapter"	August 2020					
Arjun Singhvi et al., 1RMA: Re-Envisioning Remote Memory Access (SIGCOMM'20): "A new RDMA protocol with a secure transport"	August 2020					

2005

### CAN APPLICATION-LEVEL SECURITY BE USED?

#### One-sided RDMA requests are completely performed by the NIC

- No CPU involvement on the destination machine
- The data must be encrypted/decrypted by the host CPU
- No protection against replay attacks

#### Two-sided communication is also offloaded to the NIC

- Packets cannot be discarded by the NIC
- Received data consumes resources of the connection
- The host CPU is responsible for verifying the received data negating RDMA advantages



InfiniBand is in urgent need of a secure transport

### **sRDMA – SECURE RDMA COMMUNICATION**

- sRDMA is lightweight security extension to RDMA which uses symmetric key cryptography to provide
  - Header Authentication
  - Packet Authentication
  - Payload encryption
  - Memory protection

#### sRDMA effectively prevents:

- Eavesdropping
- Spoofing attacks
- Replay attacks
- Man in the middle attacks

sRDMA is backward-compatible with InfiniBand Architecture (IBA), and can be easily adapted by

- Native InfiniBand
- RoCEv1
- RoCEv2

### **sRDMA – SECURE QP CONNECTION**

#### sRDMA introduces a new Secure Reliably Connected Queue Pair

• The application installs symmetric keys to a QP connection and required level of protection

```
//Recognized attribute masks by ibv_modify_qp
enum ibv_qp_attr_mask {
    ... // previous masks
    IBV_QP_SECURE = 1 << 26,
}</pre>
```

#### Supported security codes:

```
/*Message authentication codes */
   /* Hash-based MACs*/
   /* Header authentication*/
   IBV_HDR_HMAC_SHA1_160 = 0x1001, //EVP_sha1
   IBV_HDR_HMAC_SHA2_256 = 0x1003, //EVP_sha256
   IBV_HDR_HMAC_SHA2_512 = 0x1005, //EVP_sha512
   /* Packet authentication*/
   IBV_PCKT_HMAC_SHA1_160 = 0x2001, //EVP_sha1
   IBV_PCKT_HMAC_SHA2_256 = 0x2003, //EVP_sha256
   IBV_PCKT_HMAC_SHA2_512 = 0x2005, //EVP_sha512
```

```
//Recognized attributes by ibv_modify_qp
struct ibv_qp_attr {
    ... // previous fields
    struct ibv_secure_attr secure_attr;
}
```

// Secure attributes are the key and the security code
struct ibv\_secure\_attr {
 enum ibv\_qp\_crypto qp\_crypto; // crypto type
 uint8\_t sym\_key[MAX\_KEY\_LENGTH]; // symmetric key
};

```
/* Cipher-based MACs*/
/* Header authentication*/
IBV_HDR_CMAC_AES_96 = 0x3001, //EVP_aes-128-ocb
IBV_HDR_CMAC_AES_128 = 0x3002, //EVP_aes-128-ocb
IBV_HDR_CMAC_AES_256 = 0x3004, //EVP_aes-256-ocb
IBV_HDR_CMAC_POLY1305 = 0x3005, //EVP_chacha20_poly1305
/* Packet authentication*/
IBV_PCKT_CMAC_AES_96 = 0x4001, //EVP_aes-128-ocb
IBV_PCKT_CMAC_AES_128 = 0x4002, //EVP_aes-128-ocb
IBV_PCKT_CMAC_AES_256 = 0x4004, //EVP_aes-256-ocb
IBV_PCKT_CMAC_AES_256 = 0x4005, //EVP_chacha20_poly1305
```

/* Authenticated End	cry	/ption*/	
IBV_AEAD_AES_96	=	0x5001,	//EVP_aes-128-ocb
IBV_AEAD_AES_128	=	0x5002,	//EVP_aes-128-ocb
IBV_AEAD_AES_256	=	0x5004,	//EVP_aes-256-ocb
IBV_AEAD_POLY1305	=	0x5005,	//EVP_chacha20_poly1305

### **PACKET FORMAT**

RDMA	Routing Header (RH)	Base Transport Header (BTH)		Payload		Checksums	5
IPSec over RoCE	Routing Header (RH)	IPSec Base Transport Header (BTH)		ransport er (BTH)	Payload Che		Checksums
sRDMA	Routing Header (RH)	Base Tra Header	nsport (BTH)	sRDMA header	Рау	load	Checksums
sRDMA	Routing Header (RH)	Base Tra Header	Base Transport Header (BTH)		Рау	load	Checksun

sRDMA is an extension to InfiniBand architecture:

- Routing and checksums not affected
- Secure header is processed after processing of BTH

### **sRDMA VS IPSEC OVER ROCE**

- Larger headers in IPSec compared to sRDMA
- IPSec is only for RoCE, native InfiniBand is not supported
- IPSec over RoCE is unnatural: IPSec is not aware of QP connections
- IPSec cannot support additional features such as PD sharing or memory protection
- IPSec policies are installed per remote IP, sRDMA policies are installed per QP connection.
- IPSec policies are managed by OS, sRDMA policies are managed by an application.

### **BASE TRANSPORT HEADER (BTH)**

sRDMA	Routing Header (RH)	Base Transport Header (BTH)	sRDMA header	Payload	Checksums
-------	------------------------	--------------------------------	-----------------	---------	-----------

#### Changes to BTH

• We use 3 out of 7 reserved bits from BTH to indicate the presence of the secure header

bits bytes			31-24	23-16		·16	15-8	7-0
0-3			OpCode	SEM Pad TVer		TVer	Partition Key	
4-7	F	В	Reserved 6	Destination QP				
8-11	A		Reserved 7			PSN –	Packet Sequence Nu	ımber

#### • Secure header size

- sRDMA supports 7 different MAC sizes
- Value 0 is for backward-compatibility

Size (bits)	0	96	128	160	224	256	384	512
Value	0x0	0x1	0x2	0x3	0x4	0x5	0x6	0x7

### NONCE AND PACKET SEQUENCE NUMBER (PSN)

bits bytes			31-24	23-16		16	15-8	7-0	
0-3			OpCode	SEM Pad TVer		TVer	Partition Key		
4-7	F	В	Reserved 6		D			Destination QP	
8-11	Α		Reserved 7				PSN –	Packet Sequence Nu	ımber

#### IPSec uses nonce against replay attacks

- Nonce must never be reused
- Nonce can be predictable and be transmitted in clear

#### PSN is a part of BTH

- PSN is only 24 bit which get reused after 80 ms on modern network devices
- Mellanox ConnectX-5 can send up to 200 million packets per second!

#### sRDMA extends InfiniBand PSN counters to 64 bits

- Both sender and receiver maintain 64-bit counters, but they transmit 24 least significant bits (LSB)
- As PSNs are ordered, the endpoints can recover 64 bit sequence number from 24 LSB using sliding window

### **SRDMA - AUTHENTICATION AND SECRECY**

#### Header Authentication

$$mac_{hdr} = MAC_{K_{A,B}}(nonce_{A \to B} || RH || BTH)$$

#### Packet Authentication

 $mac_{pck} = MAC_{K_{A,B}}(nonce_{A \to B} || RH || BTH || PAYLOAD)$ 

#### Payload authenticated encryption

- Nonce, RH, and BTH are passed as Additional Authenticated Data
- Payload is encrypted and sent instead of plaintext

#### Overheads of AES-128 for N secure QP connections

	Key overhead	Nonce counter	Header
IPSec	16B * <i>N</i>	16B * N	32B
sRDMA	16B * <i>N</i>	10B * N	16B

#### uint8\_t pd\_key[MAX\_KEY\_LENGTH]; // pd protection key uint8\_t mem\_key[MAX\_KEY\_LENGTH]; // memory protection key ibv\_alloc\_secure\_pd(struct ibv\_context \*context,

uint32\_t type; // IBV\_SECURE\_PD | IBV\_SECURE\_MEM

struct ibv\_secure\_pd\_attr {

uint8\_t max\_memtree\_depth; // maximum depth of OWF tree

#### sRDMA proposes to install a key (K<sub>PD</sub>) to PD, and use this key to derive QP level keys

 The key is derived using pseudorandom function (PRF) based on adapter port addresses (APA) and QPN identifiers of the endpoints. Two endpoints derive the same symmetric key.

	Key overhead	Nonce counter	Header
IPSec	16B * N	16B * N	32B
sRDMA	16B * N	10B * N	16B
sRDMA + sPD	16B	10B* <i>N</i>	16B

};

$K_{A,B} = PR$	$RF_{K_{PD}}(APA_A)$	$\parallel QPN_A \mid$	$ APA_A $	$\parallel QPN_B)$
----------------	----------------------	------------------------	-----------	--------------------

#### sRDMA introduces secure protection domains (sPDs)

ibv\_secure\_pd\_attr\* attr);

//All resources created on this pd will be secure.

//allocate secure pd.

struct ibv\_pd \*

### **EXTENDED MEMORY PROTECTION**

#### Memory protection in IBA is based on rkey tags (32 bits)

- Each one-sided RDMA request must include rkey in its request
- Any endpoint with the rkey can access the memory
- For fine-grained access control, Memory windows type 2 can be pinned to a single QP

#### sRDMA proposes scalable crypto-based memory protection

Access to sub-region (SR) with addresses [START, END )

 $K_{SR} = PRF_{K_{MR}}(START_{SR} \parallel END_{SR})$ 

sRDMA does not introduce extra header and reuses the STH

 $mac_{hdr} = MAC_{K_{A,B}}(K_{SR}||mac_{hdr})$ 

Memory sub-delegatio	n	
	KMR	
	(SRI)	
	SR 1 SR 2	SR 3
m	m+4	m+8

### **PROCESSING OF INCOMING PACKETS**



### **IMPLEMENTATION OF sRDMA**

#### sRDMA is implemented on Broadcom Stingray PS225

- Eight-core ARM A72
- DDR4 8 GB DRAM
- Supports crypto-acceleration





### **IMPLEMENTATION OF sRDMA**



### **IMPLEMENTATION OF sRDMA**



### SOURCE AUTHENTICATION LATENCY



### SOURCE AUTHENTICATION LATENCY



### SOURCE AUTHENTICATION LATENCY



### **PACKET AUTHENTICATION LATENCY**



### **AEAD LATENCY**











### **READ BANDWIDTH**



### **READ BANDWIDTH**



### **READ BANDWIDTH**



### **sRDMA PAPER ALSO INCLUDES**

- Memory sub-delegation (clients can pass access like capabilities)
- Details on the implementation
- Additional experiments









Taranov et al. sRDMA -- Efficient NIC-based Authentication and Encryption for Remote Direct Memory Access, Usenix ATC'20



### 2021 OFA Virtual Workshop

# **THANK YOU**

Konstantin Taranov

**ETH Zurich** 

