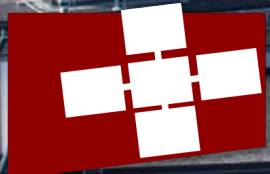


B. ROTHENBERGER, K. TARANOV, A. PERRIG, T. HOEFLER

# ReDMArk: Bypassing RDMA Security Mechanisms

Usenix Security 2021



## Collaborators for this work



Benjamin Rothenberger  
PhD Student  
Netsec



Konstantin Taranov  
PhD Student  
SPCL

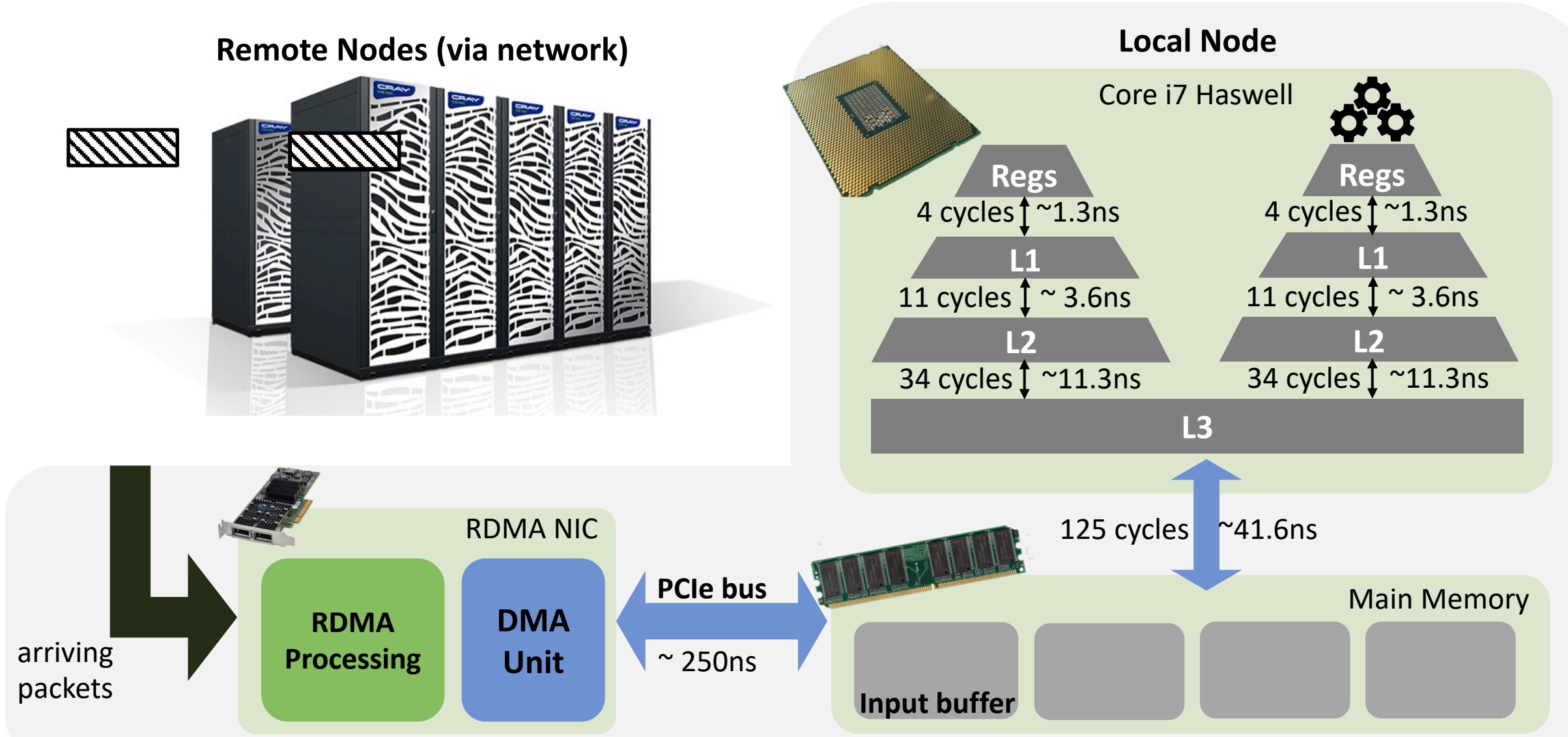


Adrian Perrig  
Professor  
Netsec



Torsten Hoefler  
Professor  
SPCL

# Data Processing in Modern RDMA Networks



# RDMA is a Trending Topic in HPC and Cloud Systems

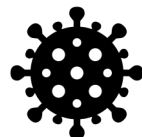
Octopus'17  
 XSTORE'20  
 Hermes'20  
 A1'20  
 ScaleRPC'19  
 TH-DPMS'16  
 FileMR'14  
 FaSST'16  
 HydraDB'15  
 DrTM+R'16  
 Dare'15  
 Wukong'16  
 FaSST'16  
 Crail'19  
 Storm'19  
 Hyperloop'18  
 DSLR'18  
 HERD'14  
 RDMP-KV'20  
 SparkRDMA'14  
 APUS'17  
 NAM-DB'17  
 RAMCloud'15  
 ccNUMA'18  
 Grappa'15  
 DaRPC'22  
 DrTM'15  
 CoRM'21  
 Catfish'19  
 Derecho'19  
 C-Hint'14

## RDMA

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



**Vulnerabilities?**



**Exploits?**



**Mitigations?**

# ReDMARK Overview

Implemented 6 attacks

4 adversary models

10 vulnerabilities

8 mitigations

	T1	T2	T3	T4	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	M1	M2	M3	M4	M5	M6	M7	M8
A1	◐	●	●	○	○	○	○	○	◐	◐	○	●	○	○	+	×	×	×	×	✓	×	✓
A2	◐	●	●	○	○	○	○	○	◐	◐	○	●	○	○	+	×	×	×	×	✓	×	✓
A3	●	◐	●	○	●	◐	◐	●	○	○	●	◐	●	◐	×	✓	×	✓	✓	×	×	×
A4	●	○	○	○	○	○	○	○	●	○	○	○	○	○	×	×	×	×	×	×	+	×
A5	●	◐	●	○	○	○	○	○	○	○	●	○	○	○	×	×	+	×	×	×	+	×
A6	○	○	○	●	○	○	○	○	○	○	●	○	○	◐	×	×	+	×	×	×	×	×

● enables attack      ○ facilitates attack      ○ does not affect attack  
 ✓ mitigates attack      + increases attack complexity      × does not mitigate attack

*weak rkey*    *static init.*    *shared key gen.*    *weak mem. rand.*    *lin. inc. QPN*    *fixed starting PSN*    *lim. starting PSN*    *no enc./auth.*    *single PD*    *ODP enabled*    *rand. QPN*    *rand. rkey*    *HW counters*    *mem. win. type 2*    *multiple PDs*    *enc./auth. in IB*    *resource const.*    *in-network filt.*

# Adversary Model

(T1) An attacker with a normal end host

- can connect to RDMA services
- issue messages over these connections

(T2) An attacker with a compromised end host

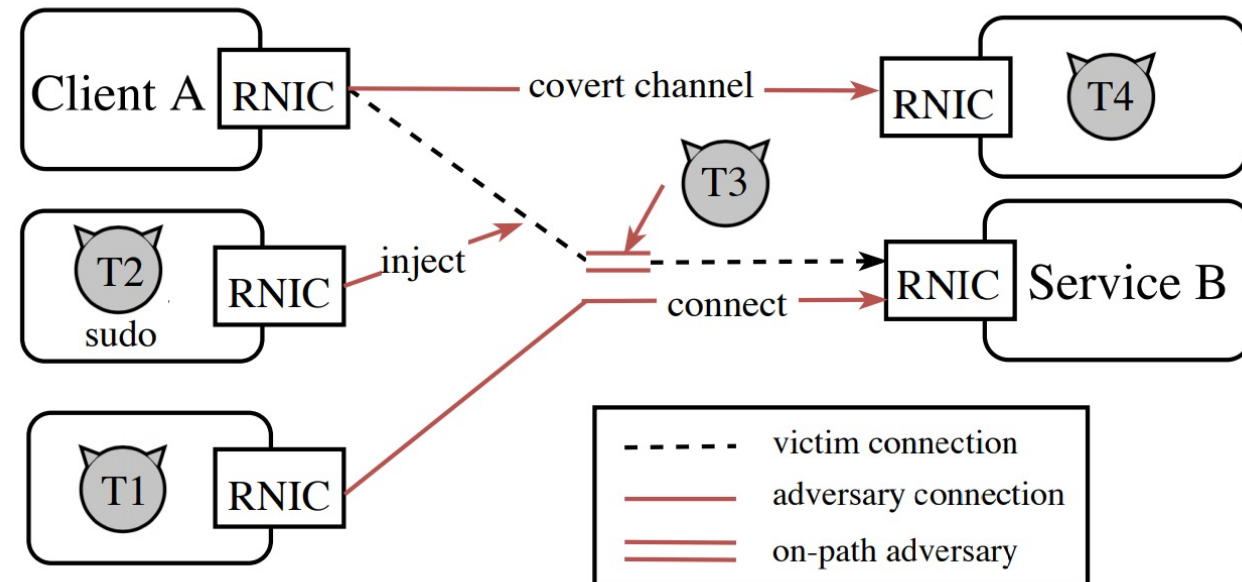
- fabricate and inject messages

(T3) An in-network attacker (e.g., malicious switch)

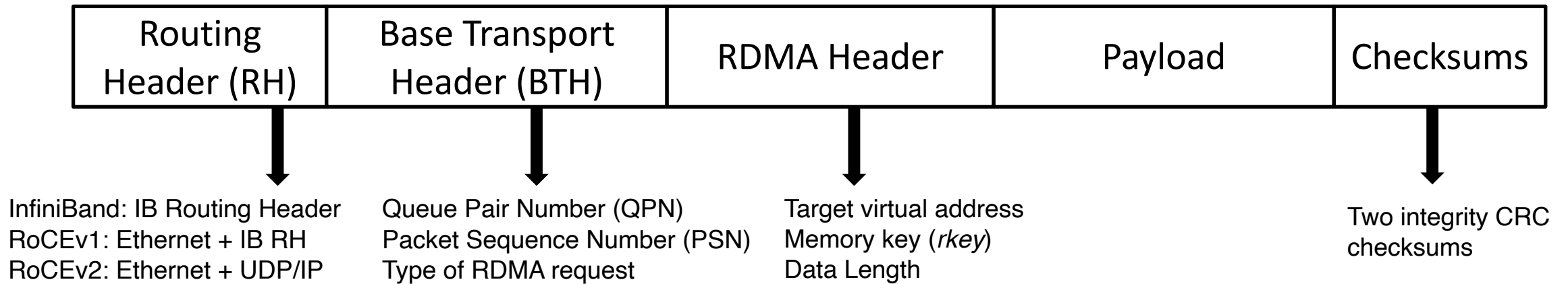
- on-path
- eavesdrop, modify

(T4) A malware-based attack

- use RDMA for data exfiltration (e.g., as covert channel)



# RDMA Write Packet Format and Packet Processing



# RDMA Write Packet Format and Packet Processing

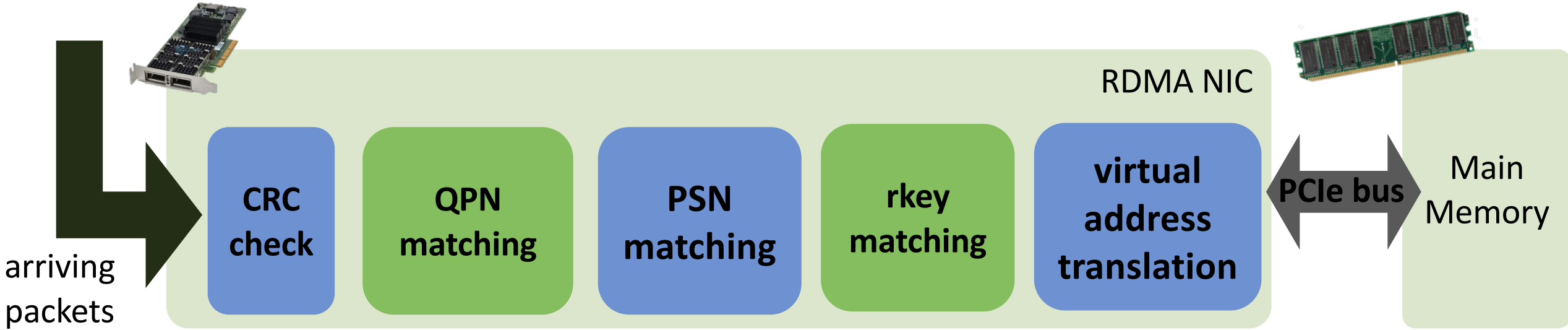


InfiniBand: IB Routing Header  
RoCEv1: Ethernet + IB RH  
RoCEv2: Ethernet + UDP/IP

Queue Pair Number (QPN)  
Packet Sequence Number (PSN)  
Type of RDMA request

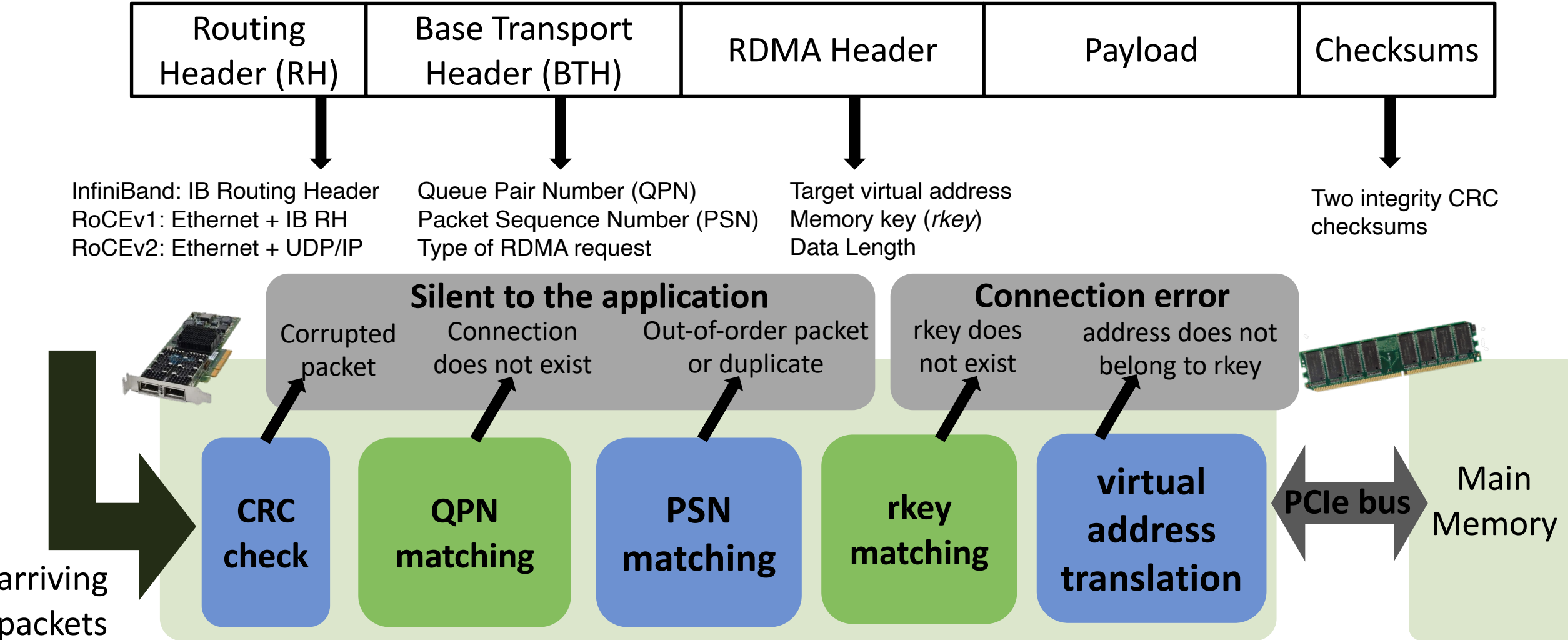
Target virtual address  
Memory key (*rkey*)  
Data Length

Two integrity CRC checksums

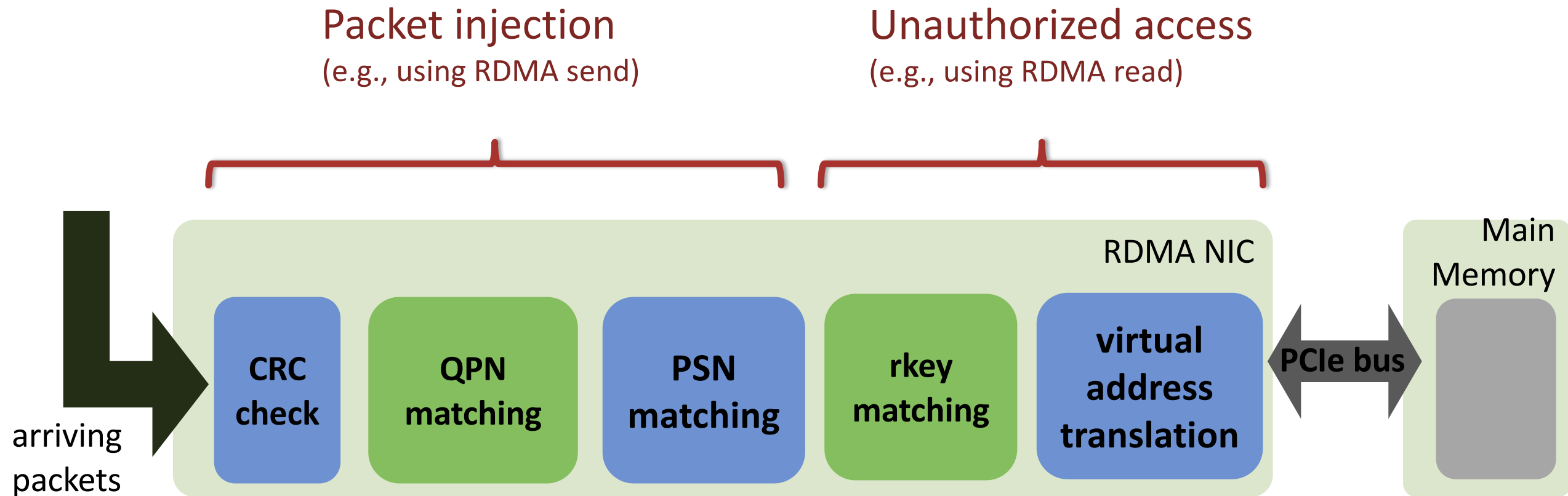




# RDMA Write Packet Format and Packet Processing



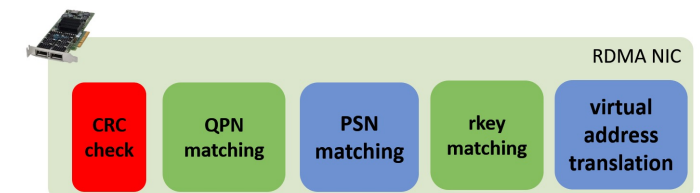
# Bypassing RDMA Processing Checks



# Towards Packet Injection -- CRC Check

- **Observations**

- Neither encryption nor authentication are used in today's RDMA protocols
- CRC checksums are used for packet integrity checks  
*but have known seeds and polynomials  
and can easily be computed by an adversary*



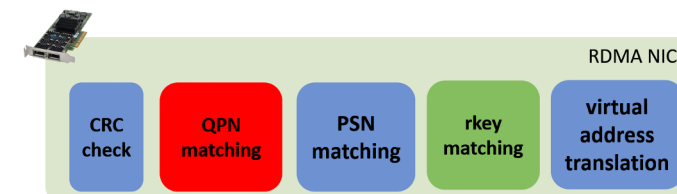
## Towards Packet Injection -- QPN Matching

### ■ Observations

- Queue pair numbers are 24 bits (< 17M possible QPNs!)
- In practice: they are allocated sequentially!
  - predicting preceding or subsequent QPNs is trivial

### ■ Device analysis

Model	Driver	Arch.	QPNs
Broadcom NetXtreme-E BCM57414	bnxt_re	RoCEv2	sequential
Broadcom Stingray PS225 BCM58802	bnxt_re	RoCEv2	sequential
Mellanox ConnectX-3 MT27500	mlx_4	IB/RoCEv1	sequential
Mellanox ConnectX-4 MT27700	mlx_5	IB/RoCEv2	sequential
Mellanox ConnectX-5 MT27800	mlx_5	IB/RoCEv2	sequential
softRoCE	rxce	RoCEv2	sequential



## Towards Packet Injection – PSN matching

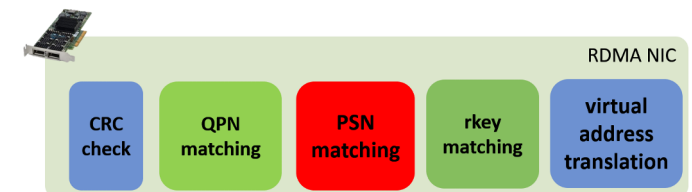
### ■ Observations

- Packet Sequence Number (PSN) is also 24 bits
- PSN can be selected by the entity that creates an RDMA connection

### ■ Connection establishment via IB verbs versus RDMA connection manager

- RDMA connection manager assigns a random PSN to the connection
- Establishing a connection using InfiniBand verbs leaves the option to the developer
- Most analysed open-source RDMA systems tend to use IB verbs with a static PSN (simplicity?)

System	Connection
Infiniswap [11]	Manager
Octopus [21]	Native
HERD [12]	Native
RamCloud [25]	Native
Dare [28]	Native
Crail [30, 31]	Manager



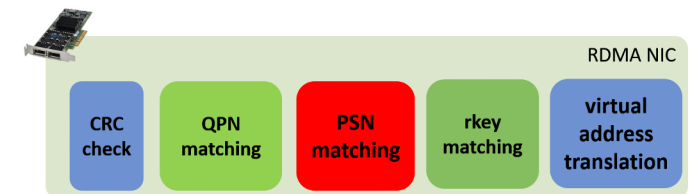
# Towards Packet Injection

## ■ Approach

- Bypassing the first three checks allows us to inject RDMA send packets (no RDMA header)
- Our PoC injection tool can inject up to 1.6 Mpps (Mellanox ConnectX-5)
  - takes roughly 11s to enumerate the full 24 bit PSN

## ■ Observations

- Injecting RDMA packets with invalid PSN does not break the connection
- Duplicate packets are dropped (and acknowledged)
  - “silent” packet replacement is possible!
- Injecting  $2^{24}$  packets makes PSN counter wrap and can hide the attack from the application



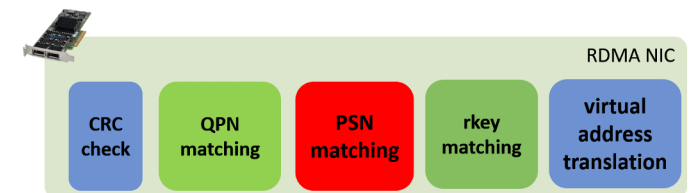
# Misuse Packet Injection for Denial-of-Service

## ■ Approach

- Packets that passed the first three checks but contain protocol errors can force the QP into an error state → breaks the QP connection!
- Our tool can inject up 1.6 Mpps
  - known PSN: we can scan 1.6 M connections per second (and disconnect!)*
  - unknown PSN: enumerate a full PSN in ~ 11s (QPN is sequential)*

## ■ Observations

- QPN randomness is crucial to increase the attack complexity for packet injection
- Example: victim with 1,000 RDMA connections with a random QPN, our tool is expected to break one of the connections in 48h



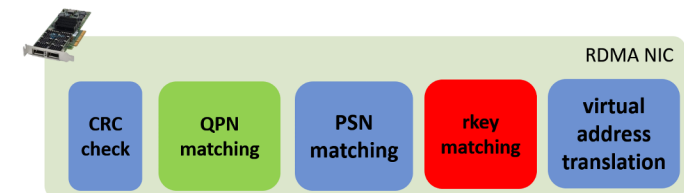
## Towards Unauthorized Access – Guessing rkeys

### ■ Observations

- rkeys are used as 32 bit access control tokens
- but: the rkey generation is highly predictable (less than 3 bits of entropy!)

### ■ Other problems

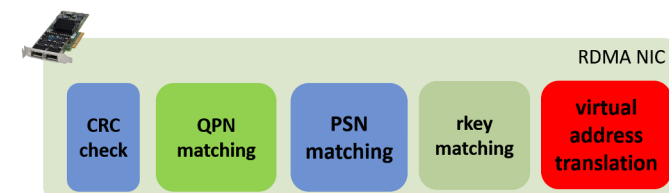
- Static initialization for key generation: the NIC generates the same keys after a reboot
- Same constant protection domain for all QPs:  
allows to access memory even without impersonation using any QP connection
- Shared key generator state: applications use the same network interface even if they use different protection domains





## Towards Unauthorized Access – Guessing Addresses

- **Virtual addresses are 64 bits**
  - Linux typically only uses 48 bits
  - Developers tend to use page-aligned memory for performance → 36 bits!
- **Consecutive allocation of memory regions**
  - Subsequent objects in memory are allocated in consecutive addresses with respect to a random address base
  - Example: InfiniSwap\*
    - Infiniswap is a remote swapping device for Linux*
    - Uses posix\_memalign in a loop to allocate register buffers of 1GB*
    - Allows an attacker to predict the position of a newly allocated buffer*



\* J. Gu, Y. Lee, Y. Zhang, M. Chowdhury, and Kang G. Shin. 2017. Efficient memory disaggregation with INFINISWAP. In *Proceedings of the 14th USENIX Conference on Networked Systems Design and Implementation (NSDI'17)*.

# Implementing Unauthorized Memory Access

## Approach

- An attacker (M) can connect to a RDMA service to get an address and a rkey of its communication buffer
- Assuming the buffers are allocated sequentially, the attacker can guess addresses and rkeys of other buffers on the service
- All 6 analysed open-source RDMA systems were vulnerable to this attack
- Attack is even simpler for an in-network attacker (eavesdrop rkey and buffer addresses)

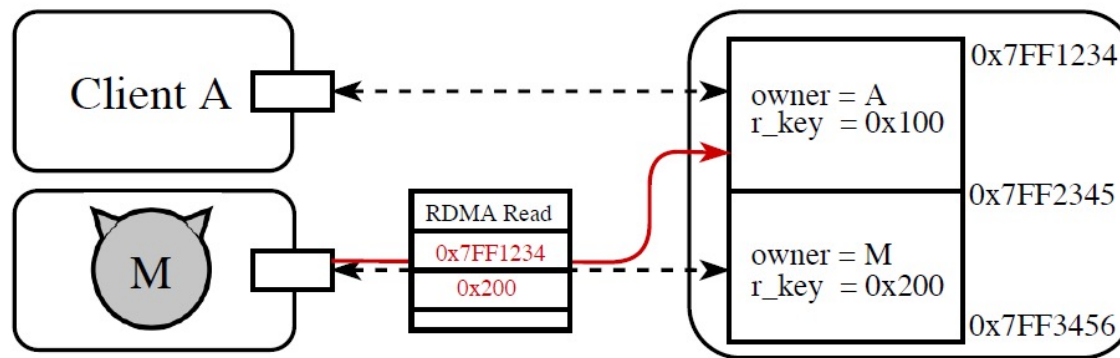
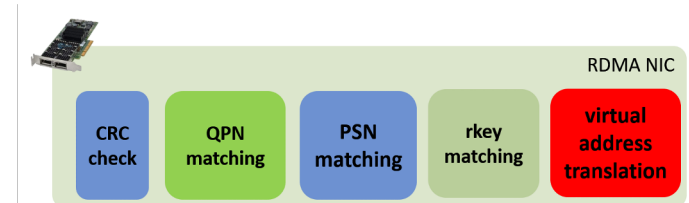


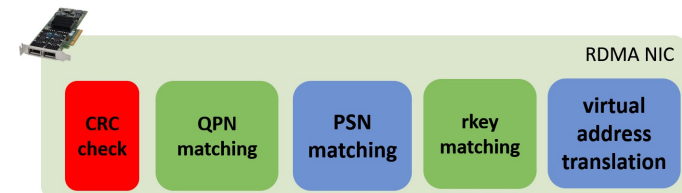
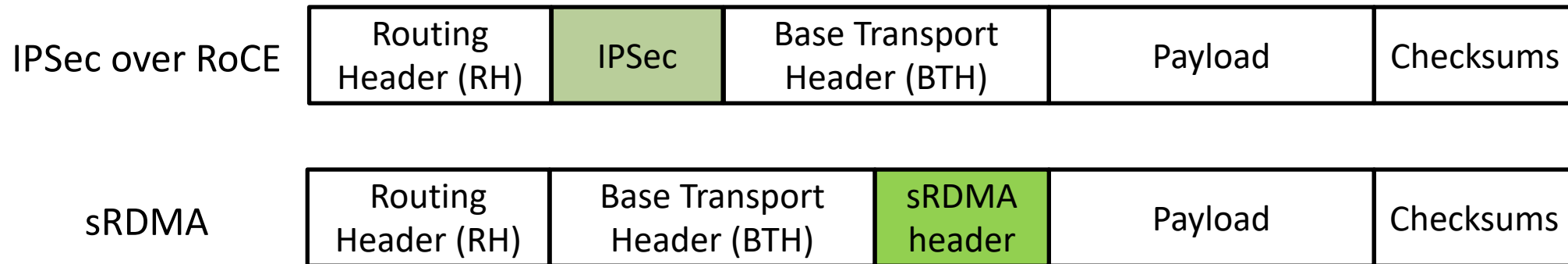
Figure 4: Unauthorized memory access on the same host.



# Mitigation Mechanisms

# Prevent Packet Injection and Packet Alteration

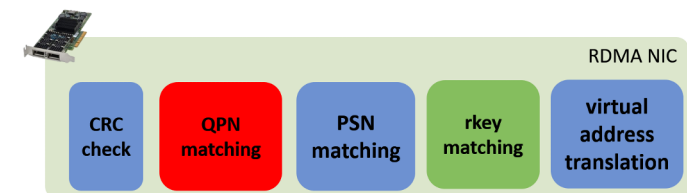
- **Use a secure transport with authentication**
  - IPsec for RoCE (e.g., Mellanox Connect-X 6 DX)
  - sRDMA\* for InfiniBand and RoCE



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

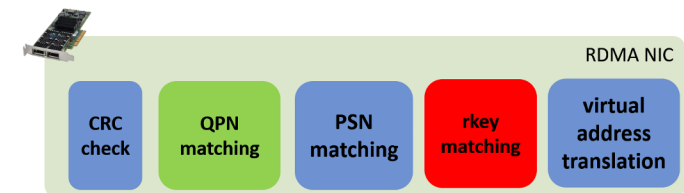
## QPN & PSN Randomization

- **We propose a software-based algorithm for QPN randomization in the paper**
  - Approach: RDMA allows creation of connection stubs that get a QPN assigned without actually connecting to a RDMA service
- **PSN randomization**
  - Use RDMA CM → randomly generates a PSN
  - But: RDMA CM exchanges connection parameters in plaintext
  - Solution: Native IB verbs interface with a random PSN



# Rkey Randomization

- **We propose a software-based Rkey randomization algorithm for short-term mitigation**
- **Use multiple Protection Domains (PDs)**
  - share PDs between trusted connections
- **Use Memory Windows Type 2**
  - can be pinned to a specific QPN
- **sRDMA proposes crypto-based memory protection**



## Additional Content in the Paper

- **Additional attacks**

- QP exhaustion
- Performance degradation
- Using RDMA for data exfiltration

- **Mitigations**

- Short-term and long-term mitigation mechanisms
- Example: software-based algorithms for QPN and rkey randomization

## Thank you for your attention!

- ReDMARK provides an in-depth analysis of current RDMA security
- We discovered 10 vulnerabilities / design flaws
- We implemented 6 attacks under 4 different threat models
- We tested 6 open-source systems
- We propose 8 mitigation techniques

ReDMARK implementation:

