# Autonomous NIC Offloads

Boris Pismenny    Haggai Eran    Aviad Yehezkel
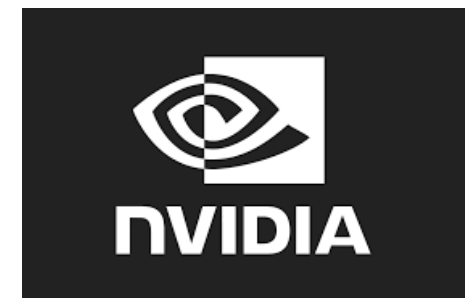
Liran Liss    Adam Morrison    Dan Tsafrir

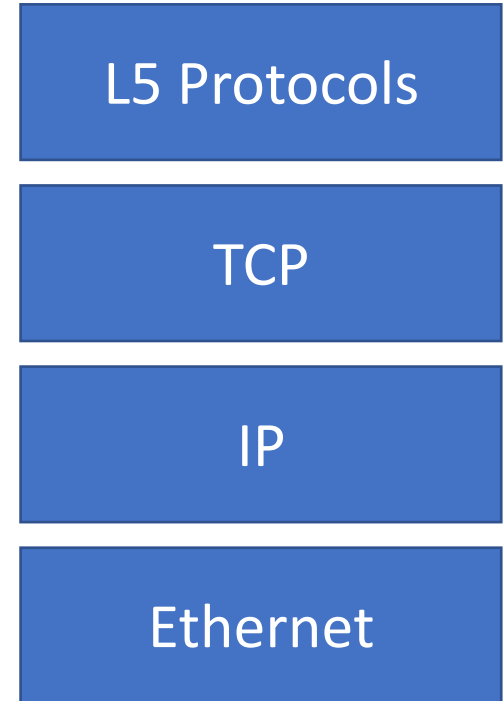How to accelerate application layer (L5) computations transparently to software TCP/IP?

# Offloading data-intensive layer-5 protocols

**L5P examples**

- tls
- nvme-tcp
- http
- grpc
- thrift
- iscsi
- nbd

**Computation examples**

- encryption
- decryption
- digest
- copy
- pattern matching
- (de)serialization
- (de)compression

L5 Protocols

TCP

IP

Ethernet

# Offloading data-intensive layer-5 protocols

**L5P examples**

- tls
- nvme-tcp
- http
- grpc
- thrift
- iscsi
- nbd

**Computation examples**

- encryption
- decryption
- digest
- copy
- pattern matching
- (de)serialization
- (de)compression

| L5 Protocols |
|:---:|
| TCP |
| IP |
| Ethernet |

# Offloading data-intensive layer-5 protocols

**L5P examples**

- tls
- nvme-tcp
- http
- grpc
- thrift
- iscsi
- nbd

**Computation examples**

- encryption
- decryption
- digest
- copy
- pattern matching
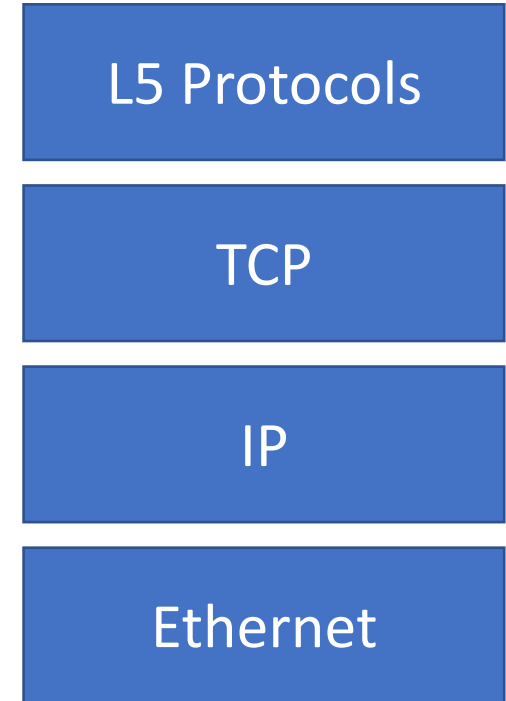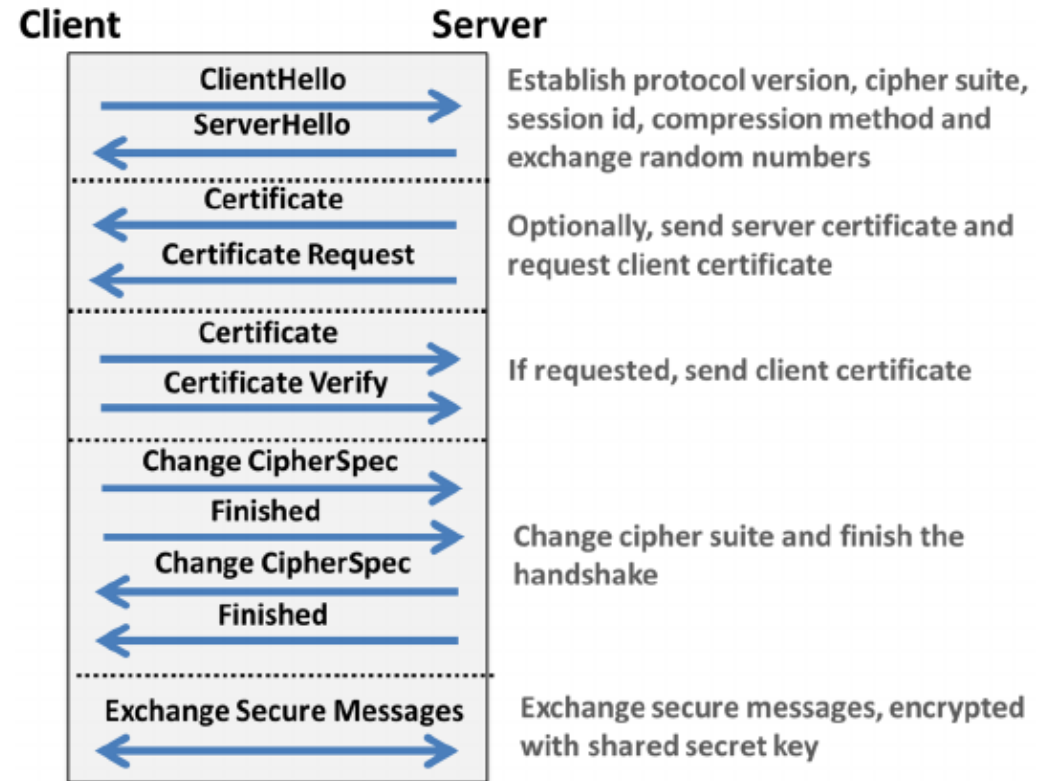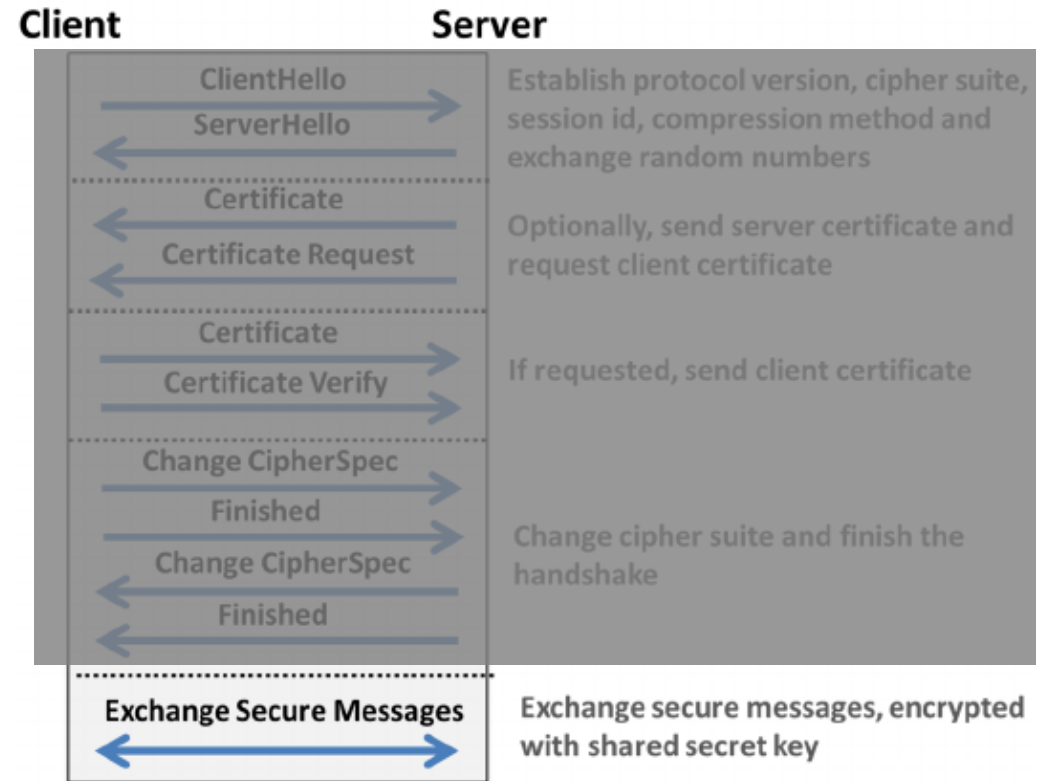- (de)serialization
- (de)compression

| L5 Protocols |
| --- |
| TCP |
| IP |
| Ethernet |

4

# What is TLS?

- Most popular way to encrypt TCP traffic
- 2 stages
  - Handshake
  - Data transfer



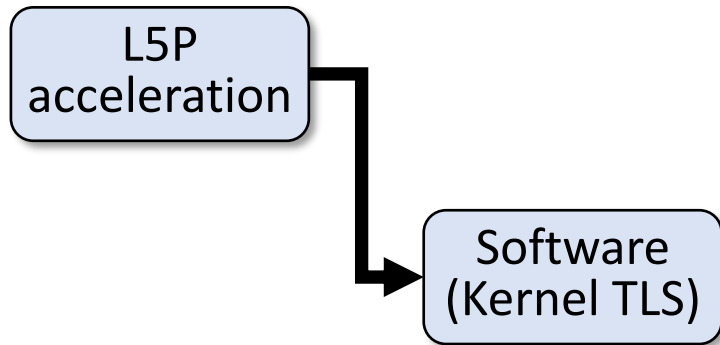| Client | | Server |
| --- | --- | --- |
| ClientHello → | | Establish protocol version, cipher suite, session id, compression method and exchange random numbers |
| ServerHello ← | | |
| Certificate ← | | Optionally, send server certificate and request client certificate |
| Certificate Request ← | | |
| Certificate → | | If requested, send client certificate |
| Certificate Verify → | | |
| Change CipherSpec → | | Change cipher suite and finish the handshake |
| Finished → | | |
| Change CipherSpec ← | | |
| Finished ← | | |
| Exchange Secure Messages ↔ | | Exchange secure messages, encrypted with shared secret key |

# What is TLS?

- Most popular way to encrypt TCP traffic
- 2 stages
    - Handshake
    - Data transfer
- We focus on data transfer
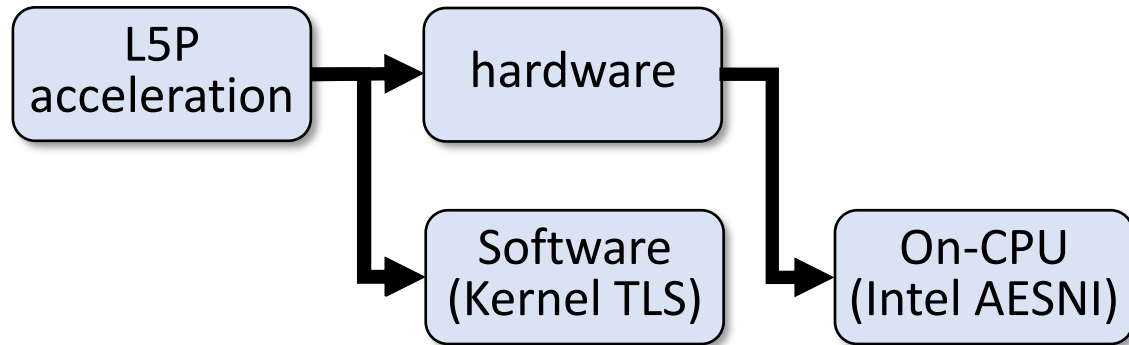
# Design Space

L5P acceleration

# Design Space

```
┌─────────────────┐
│       L5P       │
│   acceleration  │───┐
└─────────────────┘   │
                      │   ┌─────────────────┐
                      └──▶│     Software    │
                          │   (Kernel TLS)  │
                          └─────────────────┘
```

| Pros | Cons |
|------|------|
| No additional hardware | Can't avoid data intensive computations |

# Design Space

```
┌─────────────┐      ┌─────────────┐
│     L5P     │─────▶│  hardware   │──────┐
│acceleration │──┐   └─────────────┘      │
└─────────────┘  │                        ▼
                 │   ┌─────────────┐   ┌─────────────┐
                 └──▶│  Software   │──▶│   On-CPU    │
                     │(Kernel TLS) │   │(Intel AESNI)│
                     └─────────────┘   └─────────────┘
```

| Pros | Cons |
|------|------|
| Uses fast CPU registers and cache memory | can consume >50% core to compute |
| Low overhead | |

# Design Space

```
┌──────────────┐      ┌──────────────┐      ┌──────────────┐
│     L5P      │ ───► │   hardware   │ ───► │   Off-CPU    │ ───┐
│ acceleration │      │              │      │              │    │
└──────┬───────┘      └──────────────┘      └──────────────┘    │
       │                     │                     │            │
       ▼                     ▼                     ▼            ▼
┌──────────────┐      ┌──────────────┐      ┌──────────────┐
│   Software   │      │   On-CPU     │      │    Other     │
│ (Kernel TLS) │      │(Intel AESNI) │      │ (Intel QAT)  │
└──────────────┘      └──────────────┘      └──────────────┘
```

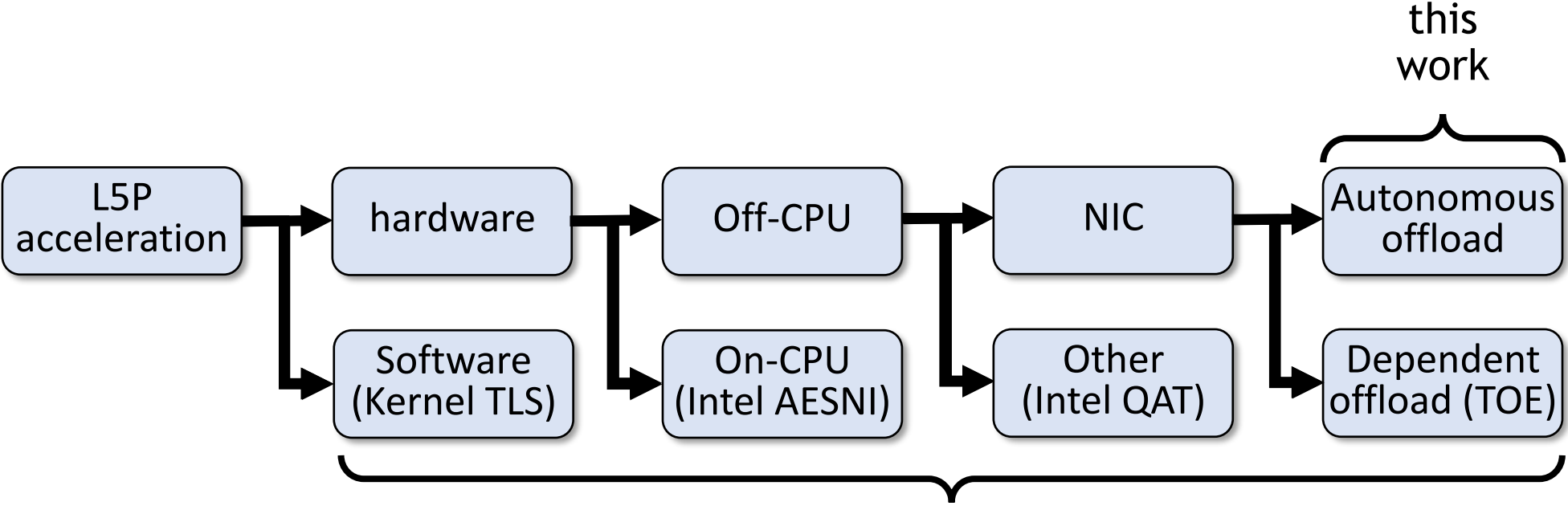| Pros | Cons |
|------|------|
| CPU overhead is independent of data size | Significant parallelism required to outperform on-CPU acceleration |

# Design Space

```
┌─────────────┐      ┌─────────────┐      ┌─────────────┐      ┌─────────────┐
│     L5P     │──┐──▶│   hardware  │──┐──▶│   Off-CPU   │──┐──▶│     NIC     │──┐
│acceleration │  │   └─────────────┘  │   └─────────────┘  │   └─────────────┘  │
└─────────────┘  │   ┌─────────────┐  │   ┌─────────────┐  │   ┌─────────────┐  │   ┌─────────────┐
                 └──▶│  Software   │  └──▶│   On-CPU    │  └──▶│    Other    │  └──▶│  Dependent  │
                     │(Kernel TLS) │     │(Intel AESNI)│     │(Intel QAT)  │     │offload (TOE)│
                     └─────────────┘     └─────────────┘     └─────────────┘     └─────────────┘
```

existing

| Pros | Cons |
|------|------|
| Eliminates CPU overhead | Depends on offloading: TCP, IP, routing, QoS, firewall, etc. |

# Design Space

this work



| L5P acceleration | hardware | Off-CPU | NIC | Autonomous offload |
|---|---|---|---|---|
| | Software (Kernel TLS) | On-CPU (Intel AESNI) | Other (Intel QAT) | Dependent offload (TOE) |

existing

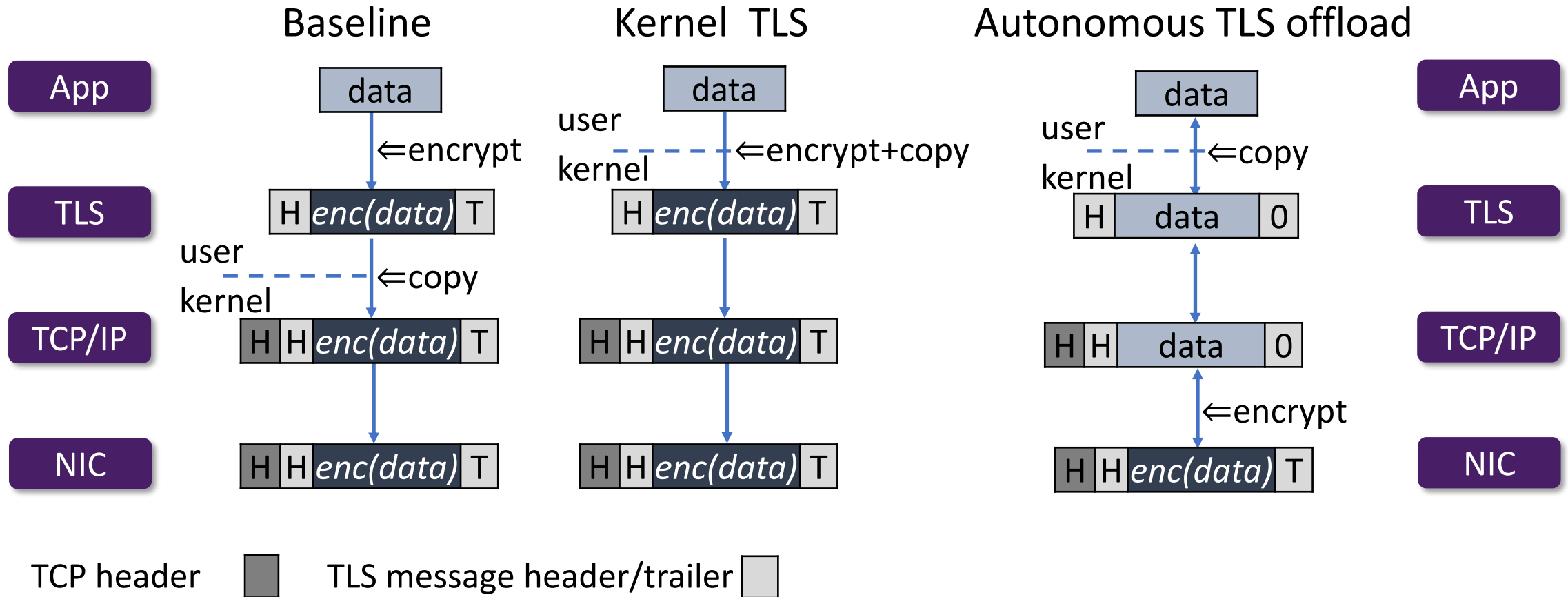| Pros | Cons |
|---|---|
| Eliminates CPU overhead | Overhead on recovery from reordering/loss |
| Works with software TCP, IP, routing, QoS, firewall, etc. | |

# Software specialization

Baseline

App

data

⇐encrypt

TLS

H *enc(data)* T

user
⇐copy

kernel

TCP/IP

H H *enc(data)* T

NIC

H H *enc(data)* T

TCP header ▮   TLS message header/trailer ▯

# Software specialization

**Baseline**

**Kernel TLS**

App

data

data

$\Leftarrow$encrypt

user
kernel
$\Leftarrow$encrypt+copy

TLS

| H | enc(data) | T |

| H | enc(data) | T |

user
kernel
$\Leftarrow$copy

TCP/IP

| H | H | enc(data) | T |

| H | H | enc(data) | T |

NIC

| H | H | enc(data) | T |

| H | H | enc(data) | T |

TCP header ▧    TLS message header/trailer ▢

**Kernel TLS enables**

- Cross layer optimization
- Direct communication between NIC and TLS layers

14

# Autonomous NIC offload: TLS



Baseline

Kernel TLS

Autonomous TLS offload

App

data

⇐encrypt

TLS

H *enc(data)* T

user
kernel

⇐copy

TCP/IP

H H *enc(data)* T

NIC

H H *enc(data)* T

data

user
kernel

⇐encrypt+copy

H *enc(data)* T

H H *enc(data)* T

H H *enc(data)* T

data

user
kernel

⇐copy

App

H data 0

TLS

H H data 0

⇐encrypt

TCP/IP

H H *enc(data)* T

NIC

TCP header �damp    TLS message header/trailer ▢

# Transmit offload in-sequence

- In-order offload is the simplest
  - Incrementally offload using NIC contexts

| TCP hdr 1 | TCP hdr 2 | TCP hdr 3 | TCP hdr 4 | TCP hdr 5 | TCP hdr 6 | TCP hdr 7 | TCP hdr 8 |

size,iv          size,iv          size,iv

## NIC contexts

**Static state**
- crypto keys

**Dynamic state**
- expected TCP seq
- current msg offset
- current msg size
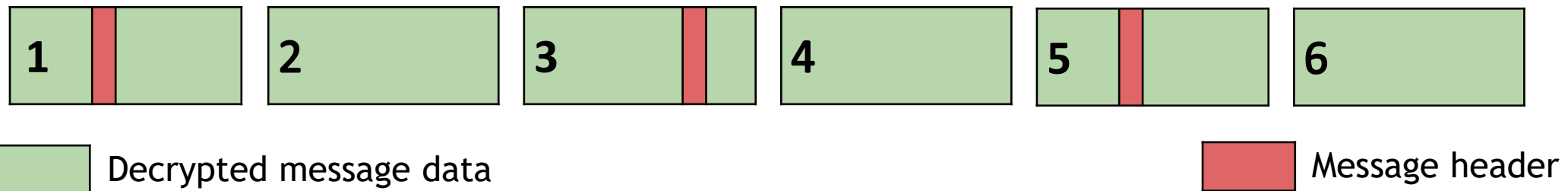- current msg IV
- msg ICV state

22

# Transmit offload out-of-sequence

- Wrong dynamic NIC context state
- Context recovery needs only the message prefix
  - Driver can get the prefix from software TLS



**NIC contexts**

Static state
- crypto keys

Dynamic state
- expected TCP seq
- current msg offset
- current msg size
- current msg IV
- msg ICV state

TCP hdr 1 | TCP hdr 2 | TCP hdr 3 | TCP hdr 4 | TCP hdr 5 | TCP hdr 6 | TCP hdr 7 | TCP hdr 8

size,iv    size,iv    size,iv

# Transmit offload out-of-sequence

- Wrong dynamic NIC context state
- Context recovery needs only the message prefix
  - Driver can get the prefix from software TLS

| TCP hdr 1 | TCP hdr 2 | TCP hdr 3 | TCP hdr 4 | TCP hdr 5 | TCP hdr 6 | TCP hdr 7 | TCP hdr 8 |
|---|---|---|---|---|---|---|---|

size,iv     size,iv     size,iv

**NIC contexts**

Static state
- crypto keys

Dynamic state
- expected TCP seq
- current msg offset
- current msg size
- current msg IV
- msg ICV state

# Transmit offload out-of-sequence

- Wrong dynamic NIC context state
- Context recovery needs only the message prefix
  - Driver can get the prefix from software TLS



NIC contexts

Static state
- crypto keys

Dynamic state
- expected TCP seq
- current msg offset
- current msg size
- current msg IV
- msg ICV state

# Transmit offload out-of-sequence

- Wrong dynamic NIC context state
- Context recovery needs only the message prefix
  - Driver can get the prefix from software TLS



- Reuse TCP transmit buffer for storing data
  - TCP ACKs release data in TLS record granularity

### NIC contexts

**Static state**
- crypto keys

**Dynamic state**
- expected TCP seq
- current msg offset
- current msg size
- current msg IV
- msg ICV state

# Receive offload in-sequence

- NIC offload Implementation is simple
  - Incrementally offload using NIC contexts

- Hardware reports one bit per packet
  - is packet decrypted and authenticated?



Decrypted message data      Message header

# Receive offload retransmission

- Retransmissions bypass offload
  - Software fallback

# Receive offload data reordering

- Record data reordering
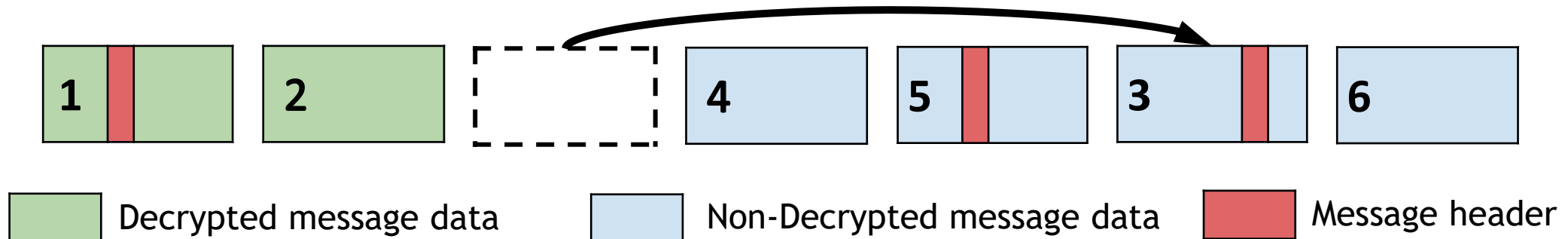  - Hardware skips to the next record
  - Continues offloading



| | 1 | | | | 3 | | | 4 | | 5 | | | 2 | | 6 | |

Decrypted message data   Non-Decrypted message data   Message header

# Receive offload data reordering

- Record data reordering
  - Hardware skips to the next record
  - Continues offloading

# Receive offload header reordering

- Record header reordering
  - Stops hardware NIC offloading
  - Software must recover NIC context to continue



| Decrypted message data | Non-Decrypted message data | Message header |

# Receive offload recovery problem

- NIC context recovery on receive is non-trivial:
  - Stopping packets to recover NIC context is impossible
    - Packets keep coming
  - Software alone cannot recover during traffic
    - Need to combine software and hardware

# Receive offload recovery solution

NIC context recovery relies on:

(1) Speculatively finding TLS message header magic pattern
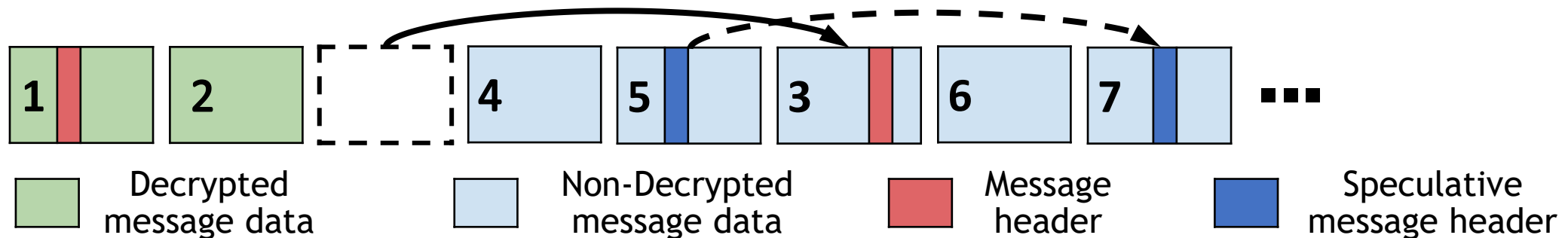   – TLS message type and version (0x170303)



Decrypted message data  
Non-Decrypted message data  
Message header  
Speculative message header

# Receive offload recovery solution

NIC context recovery relies on:
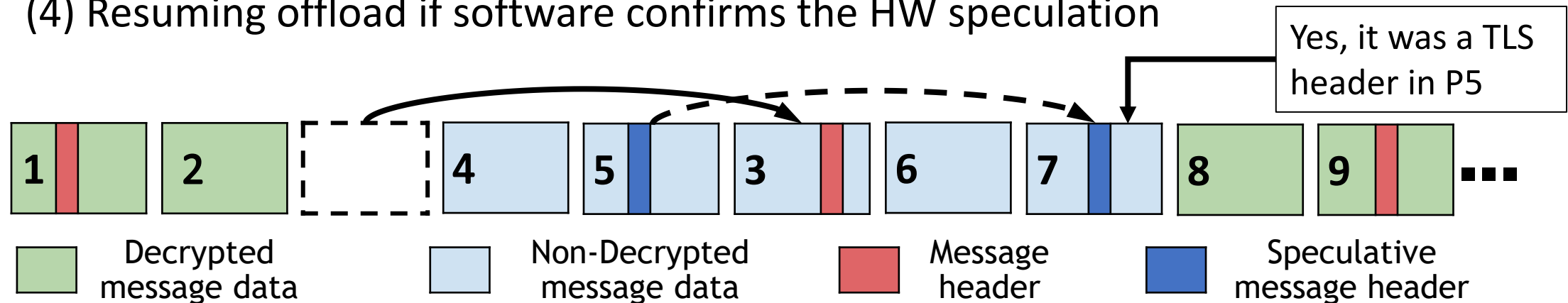
(1) Speculatively finding TLS message header magic pattern
- TLS message type and version (0x170303)

(2) Requesting software to confirm that this is indeed a TLS header, while

is it a TLS header?

| 1 | 2 | | 4 | 5 | 3 | 6 | 7 | ... |

Decrypted message data    Non-Decrypted message data    Message header    Speculative message header

# Receive offload recovery solution

NIC context recovery relies on:

(1) Speculatively finding TLS message header magic pattern
   – TLS message type and version (0x170303)

(2) Requesting software to confirm that this is indeed a TLS header, while

(3) Tracking subsequent messages using the message header's length field

# Receive offload recovery solution

NIC context recovery relies on:

(1) Speculatively finding TLS message header magic pattern
   – TLS message type and version (0x170303)

(2) Requesting software to confirm that this is indeed a TLS header, while

(3) Tracking subsequent messages using the message header's length field

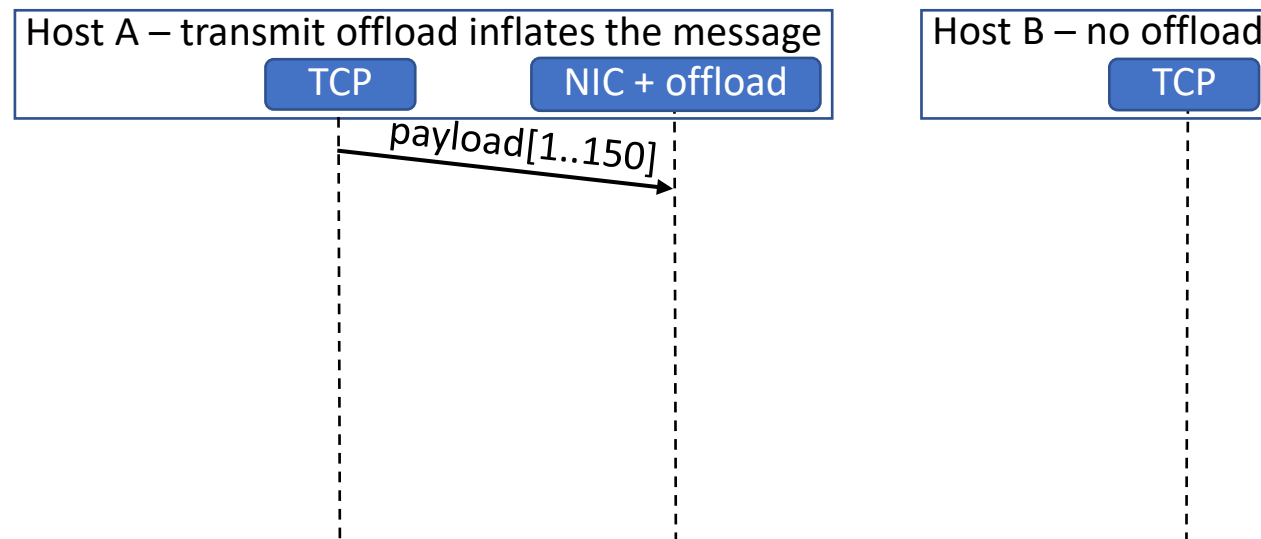(4) Resuming offload if software confirms the HW speculation



Yes, it was a TLS header in P5

Legend:
- Decrypted message data (green)
- Non-Decrypted message data (light blue)
- Message header (red)
- Speculative message header (dark blue)

# Autonomous offload properties

- What computations are autonomously offloadable?
  - Most computations, but not all

- What L5Ps are autonomously offloadable?
  - Many L5Ps, but not all

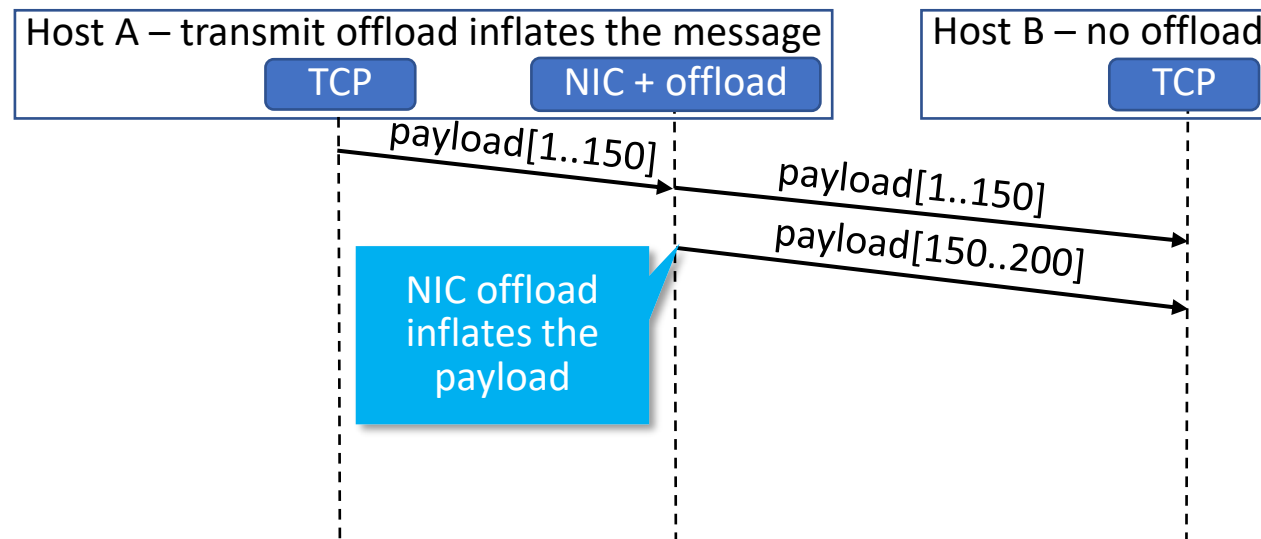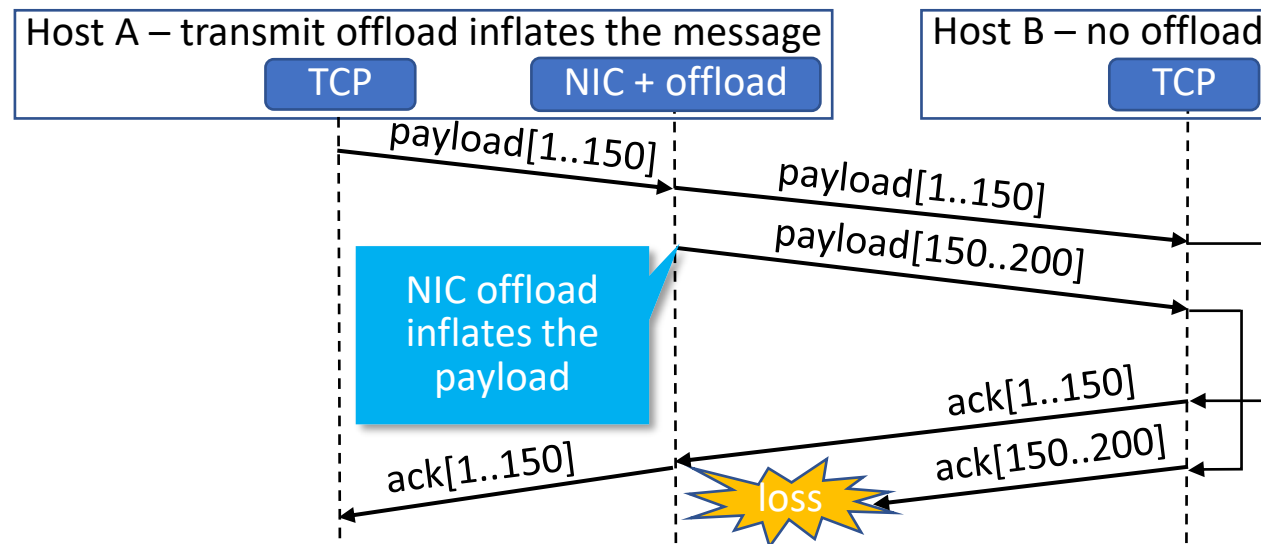# When computation is autonomously offloadable?

- On transmit, it must be size-preserving
  - This precludes transmit compression offloads

# When computation is autonomously offloadable?

- On transmit, it must be size-preserving
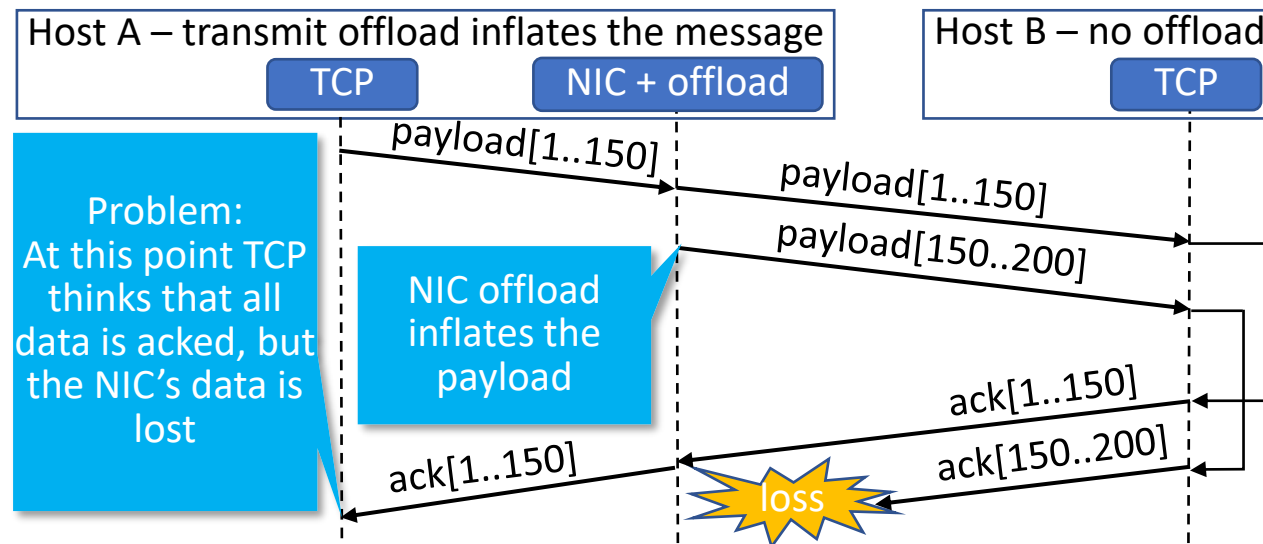  - This precludes transmit compression offloads

# When computation is autonomously offloadable?

- On transmit, it must be size-preserving
  - This precludes transmit compression offloads

# When computation is autonomously offloadable?

- On transmit, it must be size-preserving
  - This precludes transmit compression offloads



41

# When computation is autonomously offloadable?

- On transmit, it must be size-preserving
- It is computable on TCP packets of any size
  - This precludes some block ciphers (AES-CBC) which operate on 16B blocks

# When computation is autonomously offloadable?

- On transmit, it must be size-preserving

- It is computable on TCP packets of any size

- It uses constant-size message-independent state
  - It cannot depend on all stream payload
  - It can depend on message metadata (message sequence)

# When computation is autonomously offloadable?

- On transmit, it must be size-preserving

- It is computable on TCP packets of any size

- It uses constant-size message-independent state

- Many computations fit this requirement
  - encryption            – copy
  - decryption            – pattern matching
  - digest

# When L5Ps are autonomously offloadable?

- The protocol message header must contain:
    1. Message length field
    2. Plaintext magic pattern (version/opcode)

- Together these enable hardware-driven NIC context reconstruction

- Many protocols fit this requirement

    – http/2          – thrift
    – memcached       – grpc
    – iscsi           – nbd
    – smb

# Implementation
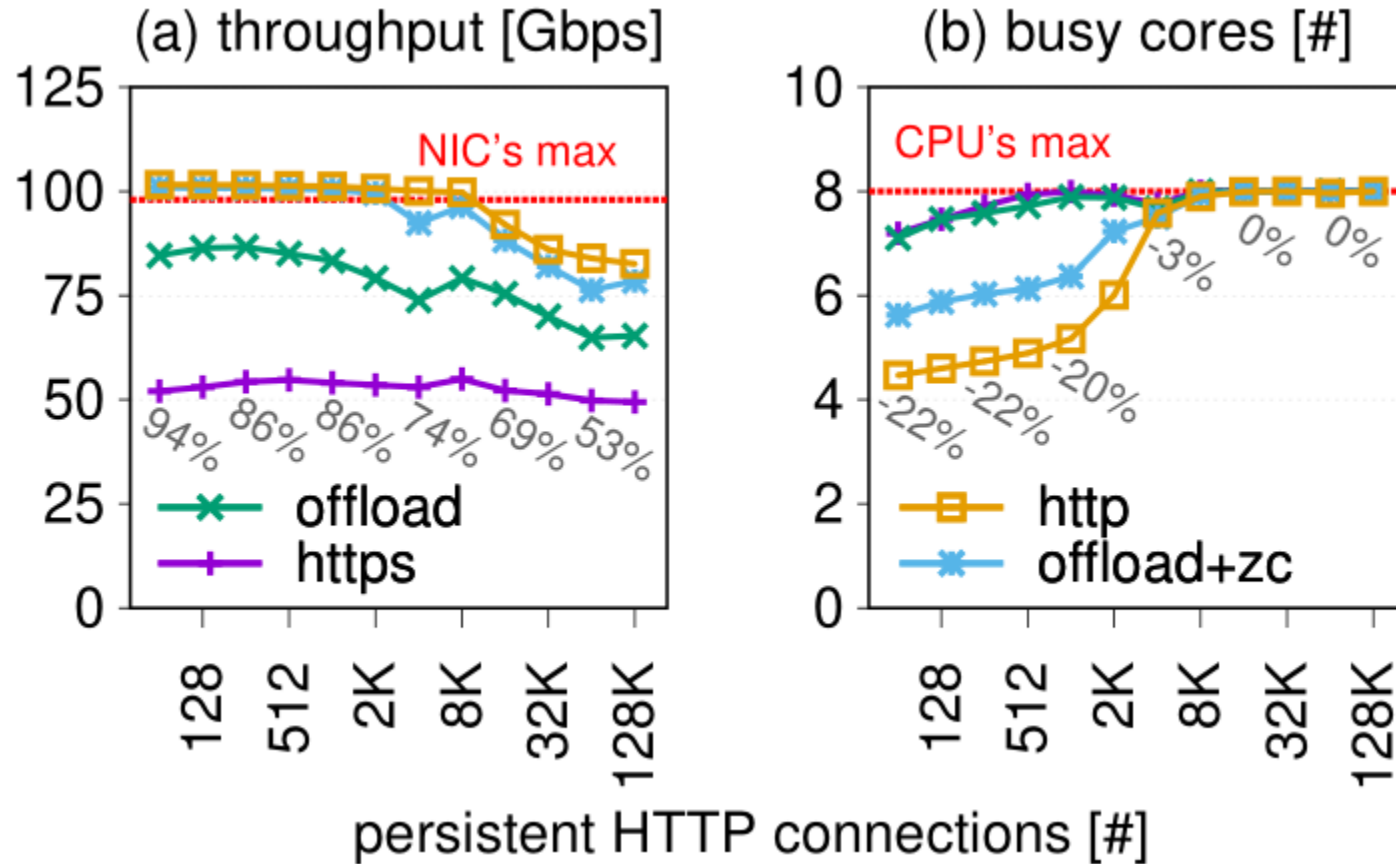
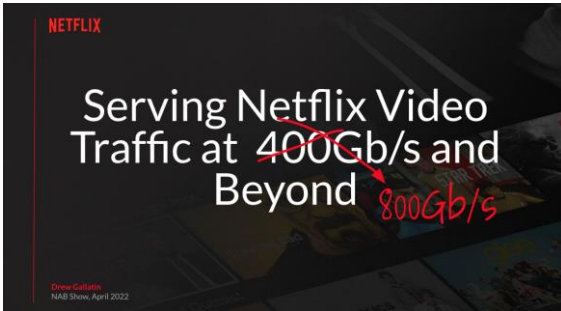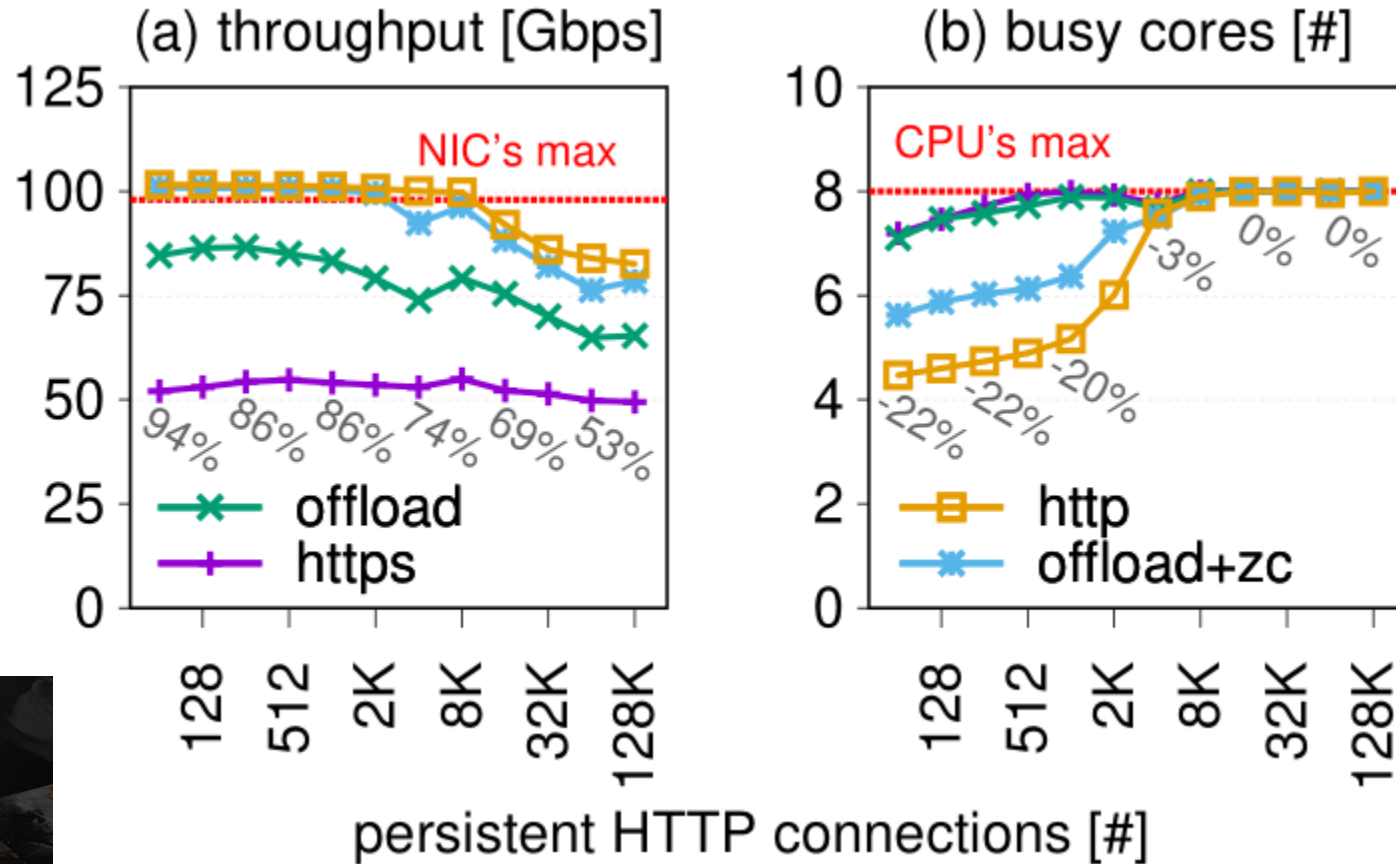- TLS crypto offload is available in Mellanox ConnectX6-Dx NICs:
    - OpenSSL:                                              1381 LoC (available upstream)
    - Linux kernel:                                    2223 LoC (available upstream)
    - Mellanox NIC driver:                     2095 LoC (available upstream)

# TLS sendfile scalability



(a) throughput [Gbps]

NIC's max

94%  86%  86%  74%  69%  53%

—✕— offload
—+— https

persistent HTTP connections [#]

(b) busy cores [#]

CPU's max

0%  0%
-3%
-20%
-22%  -22%

—□— http
—✕— offload+zc

persistent HTTP connections [#]

# TLS sendfile scalability

# Conclusion

- Autonomous NIC offloads is a framework for accelerating L5P computations efficiently while cooperating with software TCP/IP

- Autonomous NIC offloads is applicable to most protocols and computations

- Evaluation shows our approach improves throughput by up to 3.3x, and reduces CPU utilization by up to 60% and latency by up to 30%

# Thank you

Contact info:

      Boris Pismenny

      borispi@cs.technion.ac.il