# SCLP: Segment-oriented Connection-less Protocol for High-Performance Software Tunneling in Datacenter Networks

Ryota Kawashima†

Shin Muramatsu†

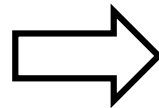Hiroki Nakayama‡

Tsunemasa Hayashi‡

Hiroshi Matsuo†

† Nagoya Institute of Technology, ‡ BOSCO Technologies, Inc.

# The Goal

Improving performance of overlay-based virtual networks

# Our Proposal

SCLP ⟹
- VXLAN (SCLP)
- Geneve (SCLP)
- …

Nagoya Institute of Technology    BOSCO Technologies

# Outline

❖ **Backgrounds**

❖ **Proposal**

❖ **Evaluation**

Nagoya Institute of Technology    BOSCO Technologies

# Network Virtualization

❖ **Multi-tenant Datacenter Networks**

- Each tenant can have its own <u>virtual networks</u>
- Each virtual network shares the physical network resources

# The Overlay-based Approach

## ❖ NVO3: Network Virtualization Overlays

- RFC 7364, 7365

NVE : Network Virtualization Edge



Physical network

# Outline

❖ **Backgrounds**

➢ Network Virtualization
➢ Tunneling Protocols
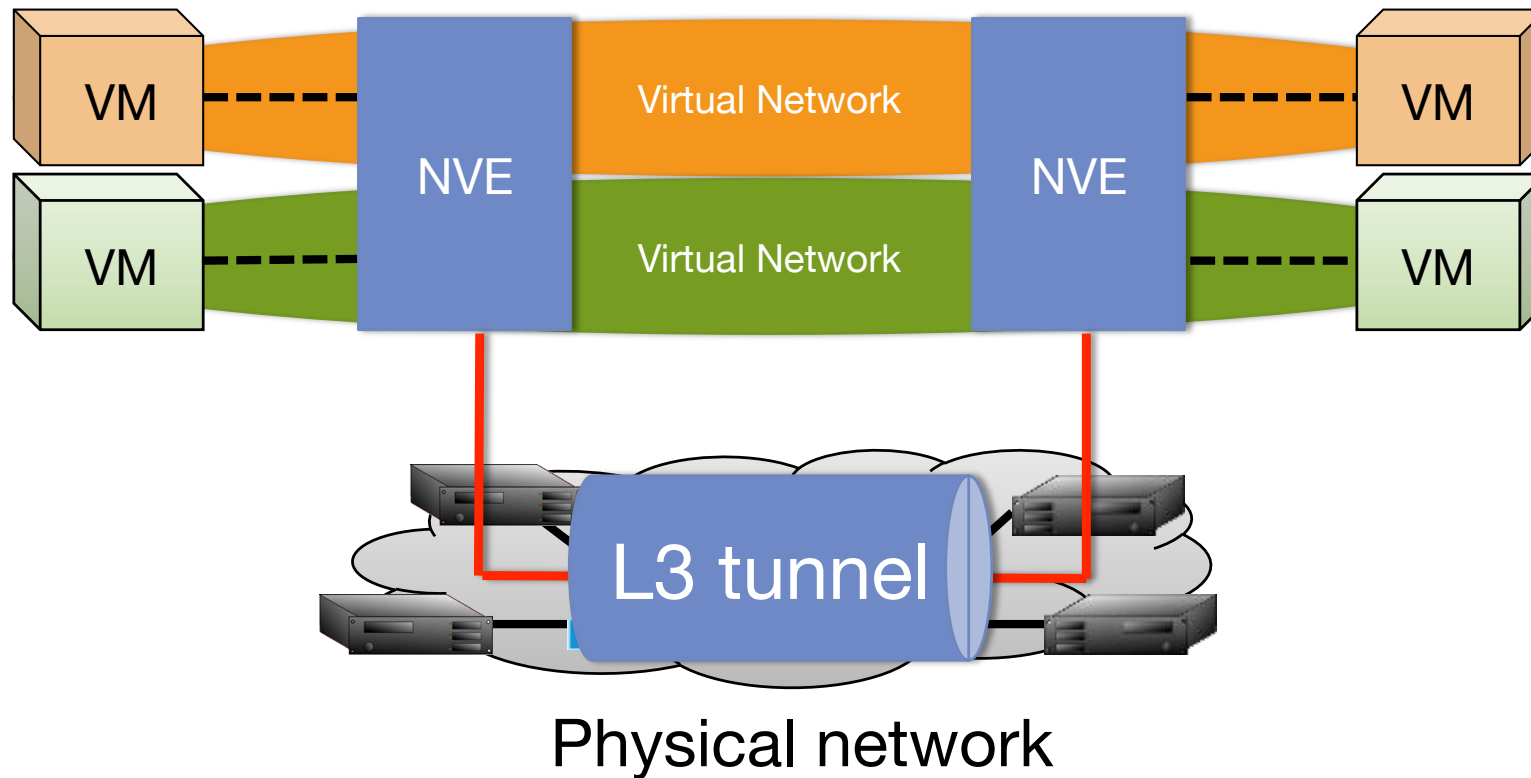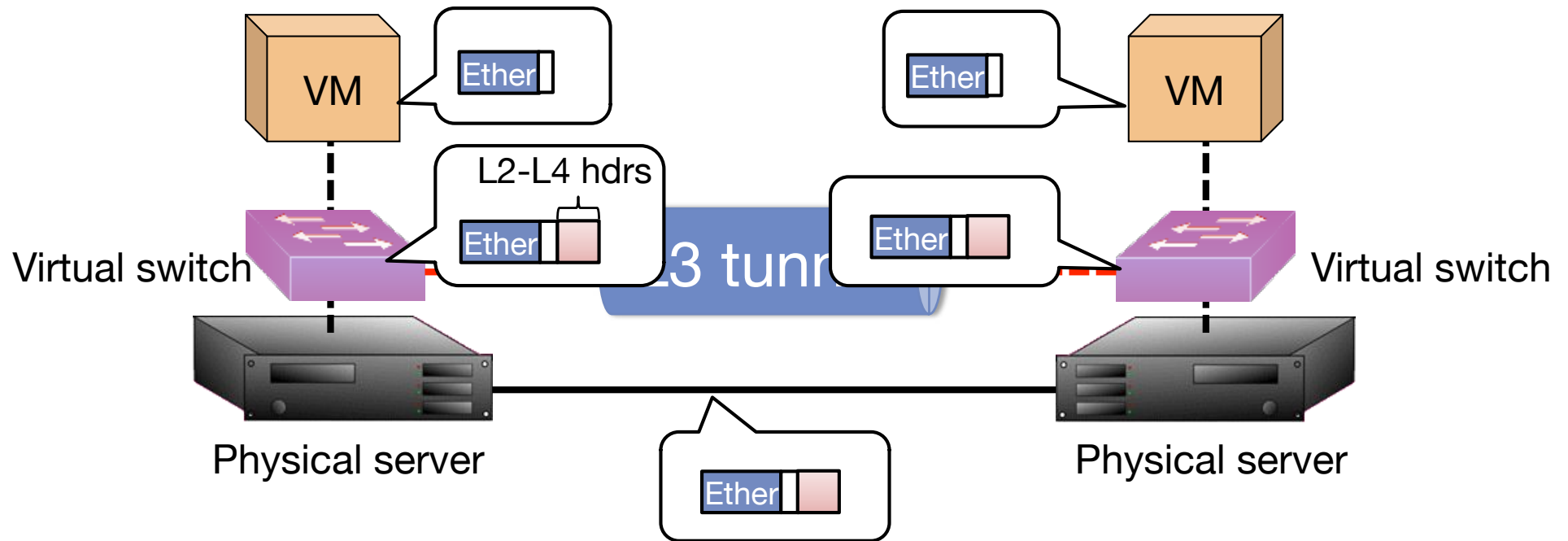➢ L4 protocol characteristics

❖ **Proposal**

➢ SCLP (Segment-oriented Connection-less Protocol)

❖ **Evaluation**

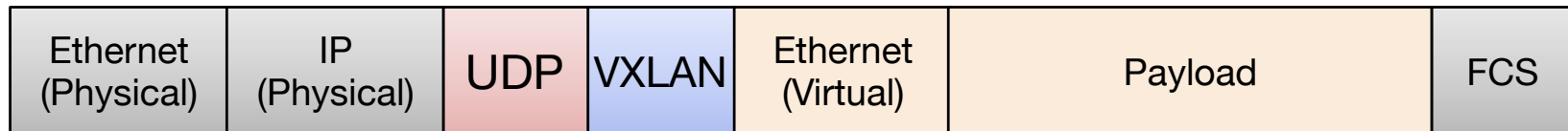➢ VM-to-VM communication using VXLAN over SCLP

# Tunneling Protocols

❖ **L2-in-L3 Tunneling**

# Major Tunneling Protocols

❖ **VXLAN (RFC 7348)**

- UDP based
- Linux kernel, OVS, VMware NSX, Cisco Nexus 1000V …

| Ethernet (Physical) | IP (Physical) | UDP | VXLAN | Ethernet (Virtual) | Payload | FCS |
|---|---|---|---|---|---|---|

❖ **NVGRE (RFC draft)**

- GRE based (no L4 protocol)
- Microsoft Hyper-V

| Ethernet (Physical) | IP (Physical) | NVGRE | Ethernet (Virtual) | Payload | FCS |
|---|---|---|---|---|---|

# Upcoming Tunneling Protocol

## ❖ **Geneve (RFC draft)**

- UDP based
- TLV option header
- H/W segmentation offload (future)

| Ethernet (Physical) | IP (Physical) | UDP | Geneve | Opt. | Ethernet (Virtual) | Payload | FCS |
|---|---|---|---|---|---|---|---|

Nagoya Institute of Technology  BOSCO Technologies

# Yet Another Tunneling Protocol

❖ **STT (Stateless Transport Tunneling, RFC draft)**

- Pseudo-TCP header
  - ➤ Exploiting TSO (TCP Segmentation Offload) feature
  - ➤ Semantics of header fields are modified
- VMware NSX

# Problems of Existing Protocols

❖ **Performance**

- VXLAN, NVGRE, Geneve

> Maximum throughput falls to one-half !

❖ **Compatibility**

- STT

> Middleboxes can discard STT packets !

Nagoya Institute of Technology    BOSCO Technologies

# Outline

❖ **Backgrounds**

➢ Network Virtualization

➢ Tunneling Protocols

➢ L4 protocol characteristics

❖ **Proposal**

➢ SCLP (Segment-oriented Connection-less Protocol)

❖ **Evaluation**

➢ VM-to-VM communication using VXLAN over SCLP

Nagoya Institute of Technology     BOSCO Technologies

# L4 Protocol Types

❖ **Message-oriented  (e.g. UDP)**

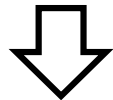- Packets are independent of each other

❖ **Segment-oriented (e.g. TCP)**

- Each packet has "byte-level sequence"
- <span style="color:red">Consecutive packets can be reassembled</span>

# Why is L4 Protocol Important ?



**Message-oriented**

A VM sends a large Ethernet frame

⬇

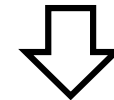The frame is encapsulated and divided to multiple packets

⬇

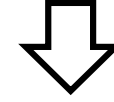Each packet is decapsulated and forwarded to a destination VM

⬇

The VM handles lots of frames
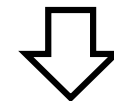
**Segment-oriented**

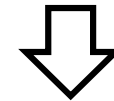A VM sends a large Ethernet frame

⬇

The frame is encapsulated and divided to multiple packets

⬇

Consecutive packets are reassembled

⬇

Each reassembled packet is decapsulated and forwarded to a destination VM

⬇

The VM handles fewer frames

Nagoya Institute of Technology    BOSCO Technologies

**13**

# Packet Structure Example (Tx)

# Packet Structure Example (Rx)

# Outline

❖ **Backgrounds**

- ➢ Network Virtualization
- ➢ Tunneling Protocols
- ➢ L4 protocol characteristics

❖ **Proposal**

- ➢ SCLP (Segment-oriented Connection-less Protocol)

❖ **Evaluation**

- ➢ VM-to-VM communication using VXLAN over SCLP

Nagoya Institute of Technology    BOSCO Technologies

# Our Proposal

❖ **SCLP: Segment-oriented Connection-less Protocol**

- L4 protocol
- Segment-oriented
- Connection-less
- Usage: Outer L4 protocol of existing tunneling protocols
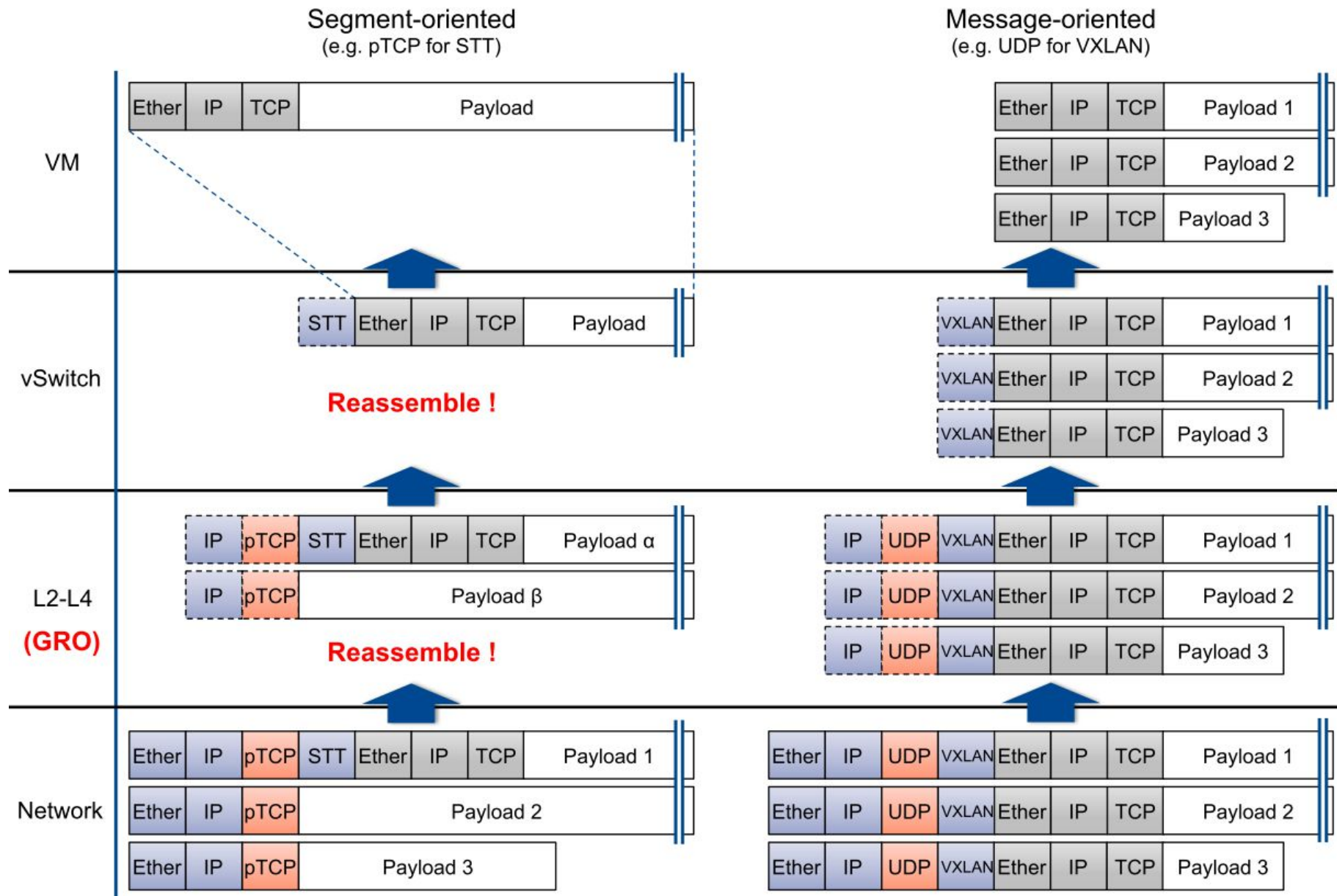  - ➢ e.g.) VXLAN over SCLP, Geneve over SCLP

| Ethernet (Physical) | IP (Physical) | UDP | VXLAN | Ethernet (Virtual) | Payload | FCS |
|---|---|---|---|---|---|---|

⬇

| Ethernet (Physical) | IP (Physical) | SCLP | VXLAN | Ethernet (Virtual) | Payload | FCS |
|---|---|---|---|---|---|---|

Nagoya Institute of Technology ▲▲ BOSCO Technologies

# Protocol Format

```
 0                   1                   2                   3
 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
```

| Source Port | Dest Port |
|---|---|
| Identification | F |
| Remaining | SCLP Checksum |

❖ **Identification:  Original payload ID**

❖ **F:             First Segment Flag**

❖ **Remaining:      Remaining payload size**

# How SCLP Works (Tx)

Original Payload
(3000 bytes)

Segmentation (MSS=1468)

'identification': 0x12345678

| SCLP | Payload 1 (1468 bytes) | | SCLP | Payload 2 (1468 bytes) | | SCLP | Payload 3 (64 bytes) |

'F': **1**
'remaining': 1532

'F': 0
'remaining': 64

'F': 0
'remaining': **0**

# How SCLP Works (Rx)

id    : -
size  : 0
offset : 0

Payload Buffer
(NULL)

id    : 0x12345678
size  : 3000
offset : 1468

① Received payload 1
(id = 0x12345678, length = 1468, F = 1, remaining = 1532)

1468
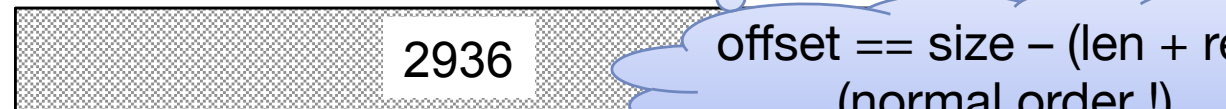
id    : 0x12345678
size  : 3000
offset : 2936

② Received payload 2
(id = 0x12345678, length = 1468, F = 0, remaining = 64)

2936

offset == size – (len + rem)
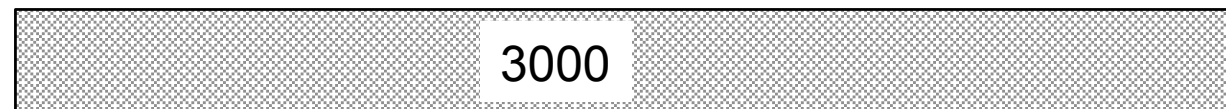(normal order !)

id    : 0x12345678
size  : 3000
offset : 3000

③ Received payload 3
(id = 0x12345678, length = 64, F = 0, remaining = 0)

3000

Nagoya Institute of Technology    BOSCO Technologies

# 2-Level Pre-reassembling

❖ **1st: GRO (Generic Receive Offload)**

```
┌─────────────────────┐
│   Protocol stack     │
└─────────────────────┘
           ↑
┌─────────────────────┐      ┌──────────────┐
│        GRO           │─────▶│ Reassembling │
└─────────────────────┘      └──────────────┘
           ↑
┌─────────────────────┐
│     NIC driver       │
└─────────────────────┘
           ↑
```

❖ **2nd: NVE's decapsulation processing**

```
┌─────────────────────┐
│         VM           │
└─────────────────────┘
           ↑
┌─────────────────────┐      ┌──────────────┐
│        NVE           │─────▶│ Reassembling │
└─────────────────────┘      │      &       │
           ↑                 │ Decapsulation│
┌─────────────────────┐      └──────────────┘
│   Protocol stack     │
└─────────────────────┘
           ↑
```

# Implementation

❖ **VXLAN over SCLP**

- CVSW component †
- Virtual NIC implementation of NVE
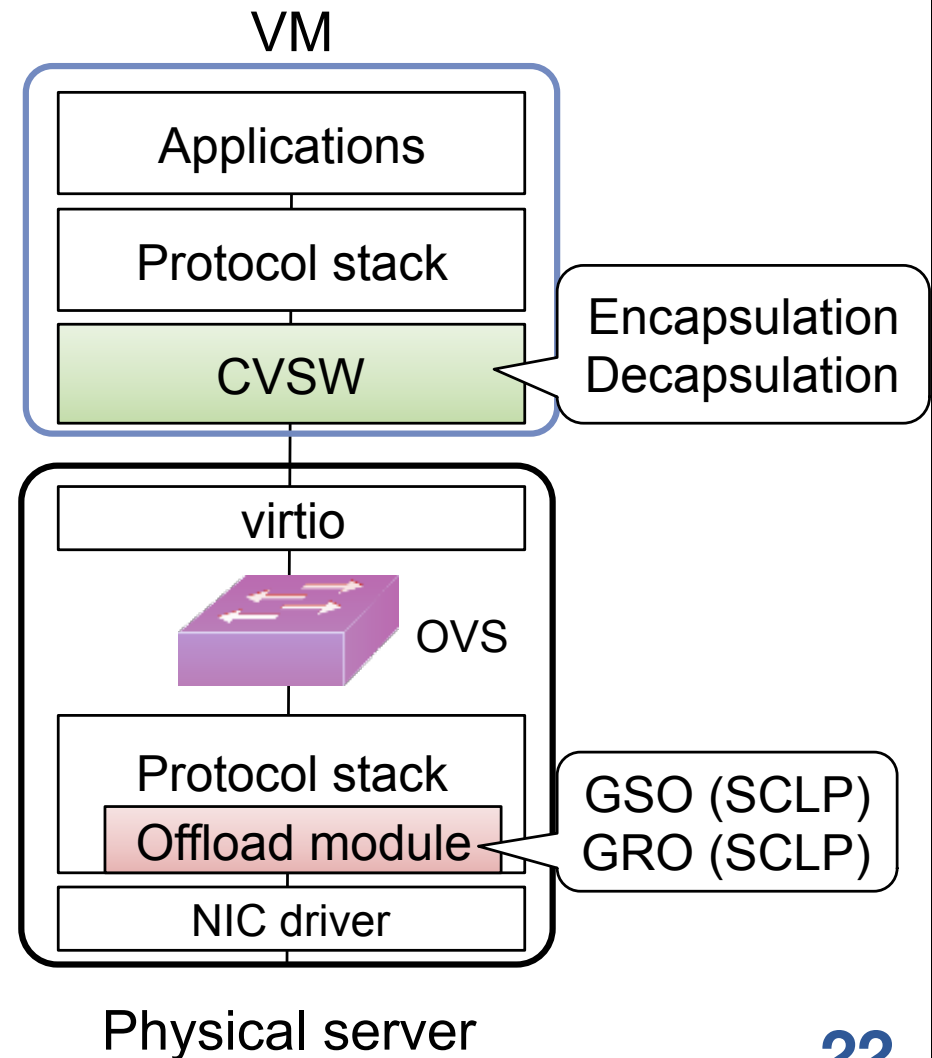
❖ **GSO/GRO offloading**

- Linux kernel module

† https://github.com/sdnnit/cvsw_net

VM

| Applications |
| --- |
| Protocol stack |
| CVSW |

Encapsulation Decapsulation

virtio

OVS

| Protocol stack |
| --- |
| Offload module |
| NIC driver |

GSO (SCLP)
GRO (SCLP)

Physical server

Nagoya Institute of Technology     BOSCO Technologies

# Outline

❖ **Backgrounds**

  ➢ Network Virtualization

  ➢ Tunneling Protocols

  ➢ L4 protocol characteristics

❖ **Proposal**

  ➢ SCLP (Segment-oriented Connection-less Protocol)

❖ **Evaluation**

  ➢ VM-to-VM communication using VXLAN over SCLP

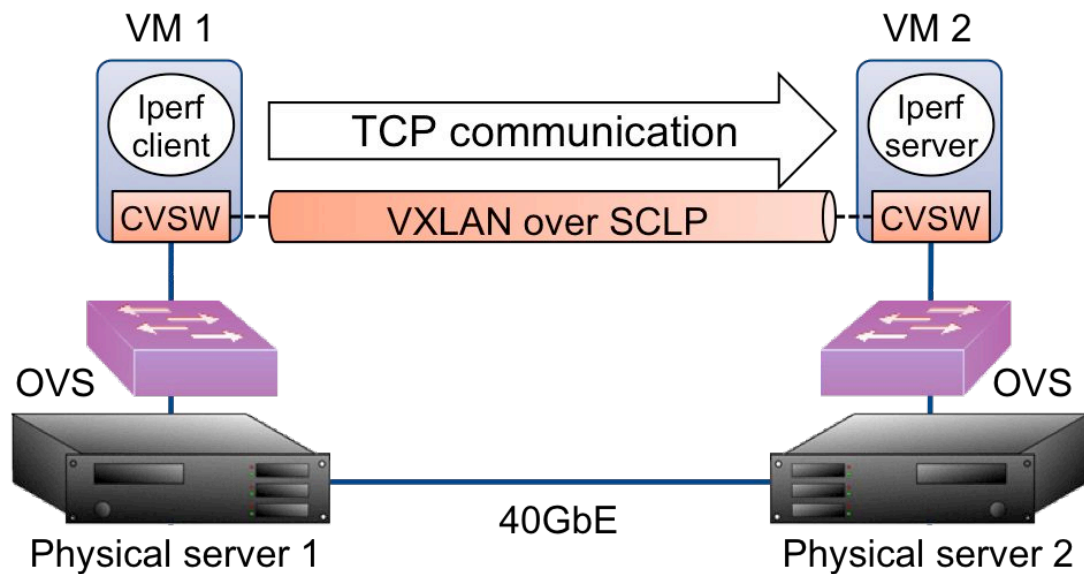Nagoya Institute of Technology  BOSCO Technologies

# Evaluation

❖ **Throughput of VM-to-VM communication with Iperf**

1. TCP communication
2. Effect of 2-level pre-reassembling

❖ **Competitors**

● VXLAN (over UDP)

● NVGRE

● STT

● Geneve (w/o HW offloading)
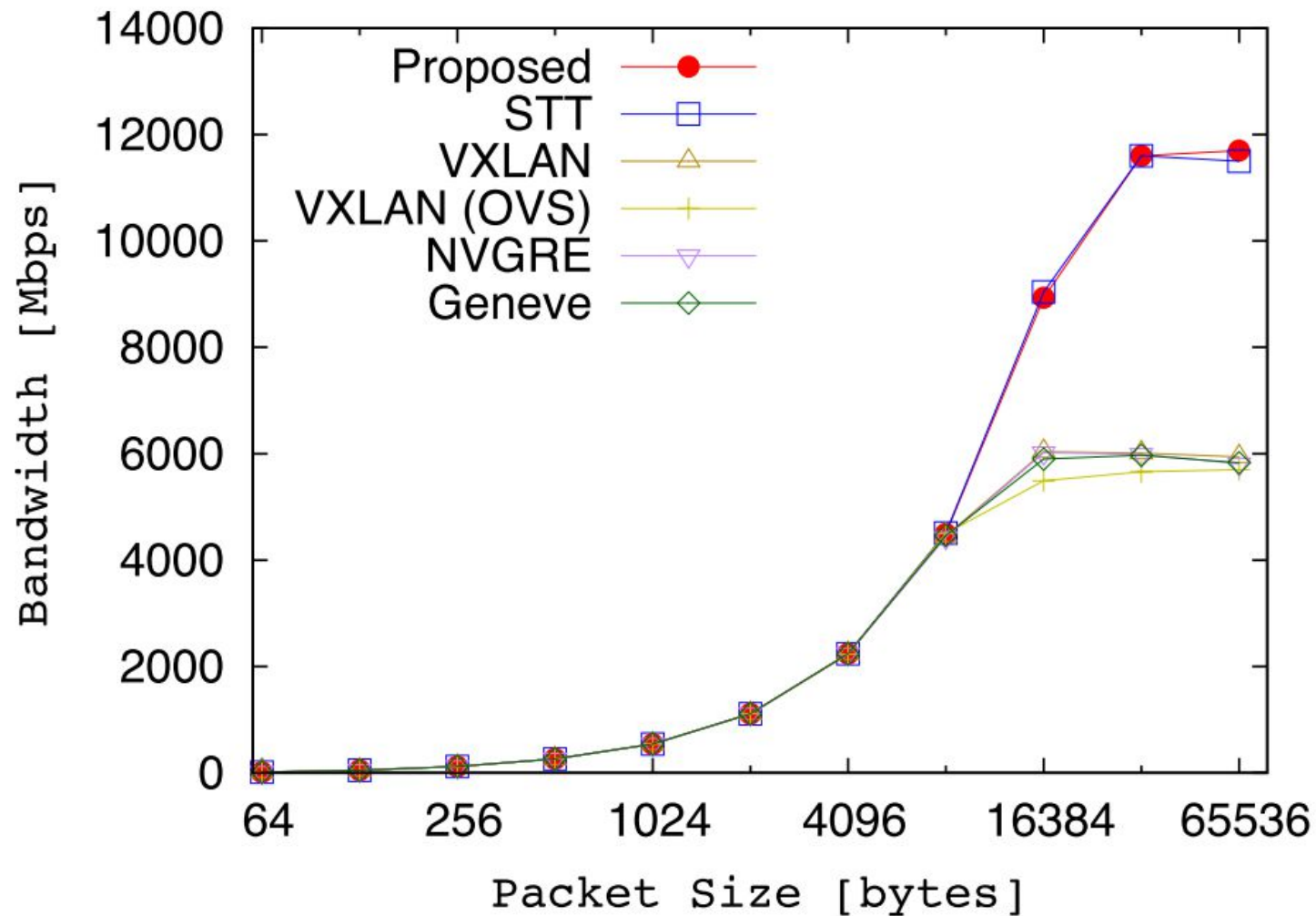
# Evaluation Environment



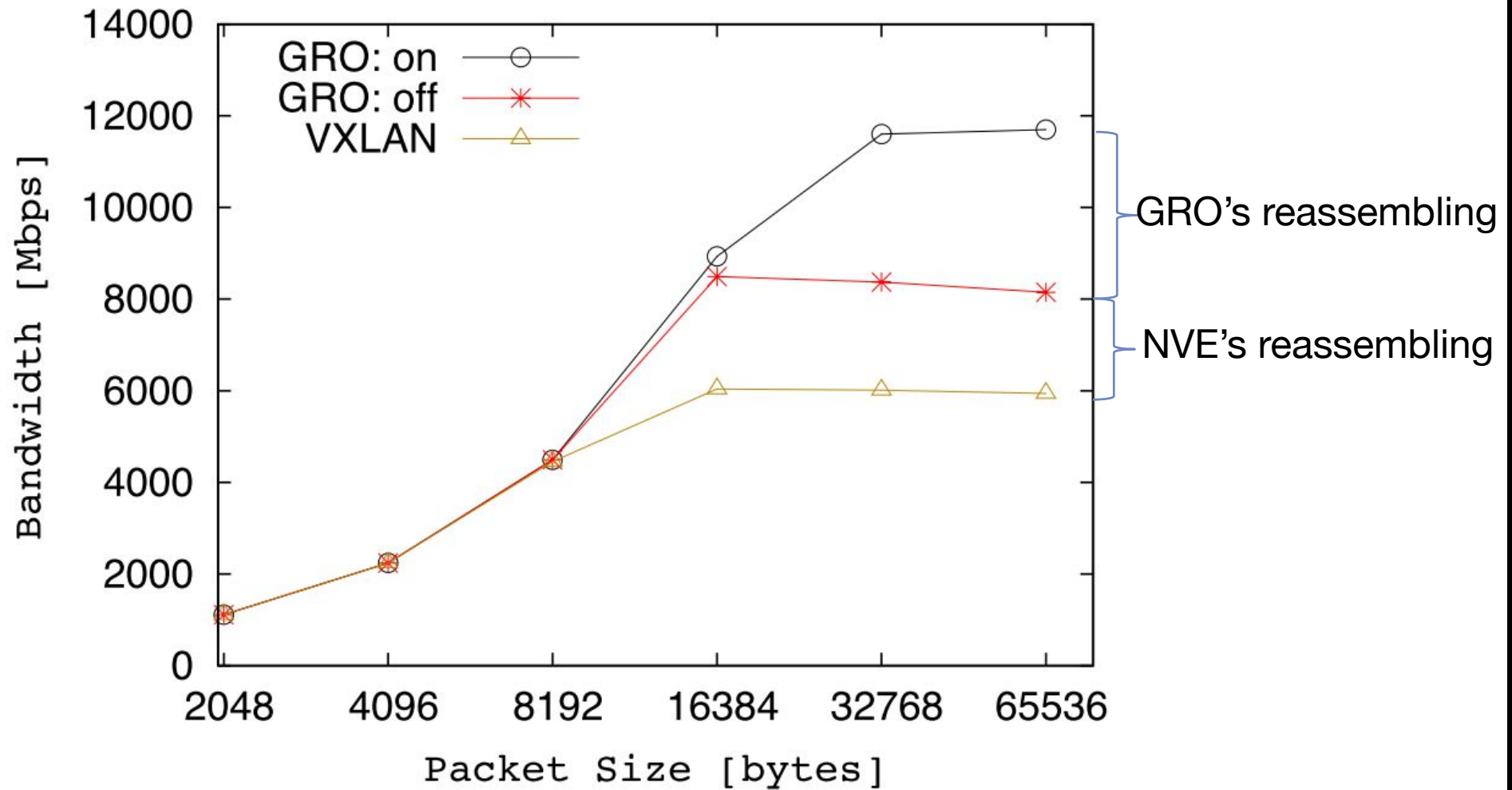| Virtual machines | VM 1 (Sender) | VM 2 (Receiver) |
|---|---|---|
| OS | CentOS 6.5 (2.6.32) | CentOS 6.5 (2.6.32) |
| CPU | 1 core | 1 core |
| Memory | 2 GBytes | 2 GBytes |
| Virtual NIC | CVSW (virtio-net) | CVSW (virtio-net) |
| MTU | adjusted | adjusted |
| Offloading features | TSO, UFO, GSO, GRO, CSUM | TSO, UFO, GSO, GRO, CSUM |

| Physical machines | Physical server 1 | Physical server 2 |
|---|---|---|
| OS | CentOS 6.5 (2.6.32) | CentOS 6.5 (2.6.32) |
| VMM | KVM | KVM |
| Virtual switch | Open vSwitch 2.3.0 | Open vSwitch 2.3.0 |
| CPU | Core i7 (3.60 GHz) | Core i7 (3.40 GHz) |
| Memory | 64 GBytes | 32 GBytes |
| MTU | 1500 bytes | 1500 bytes |
| Offloading features | TSO, GSO, GRO, CSUM | TSO, GSO, GRO, CSUM |
| Network | 40GBASE-SR4 | 40GBASE-SR4 |

Nagoya Institute of Technology    BOSCO Technologies

# Evaluation Results (TCP)

❖ **TCP communication**

# Evaluation Results (Pre-reassembling)

# Conclusion

❖ **Network virtualization**

- Overlay-based approach has become popular
- VXLAN is a de-facto tunneling protocol
- UDP-based tunneling has performance problems

❖ **Proposal: SCLP**

- Segment-oriented and connection-less L4 protocol
- 2-level pre-reassembling before decapsulation
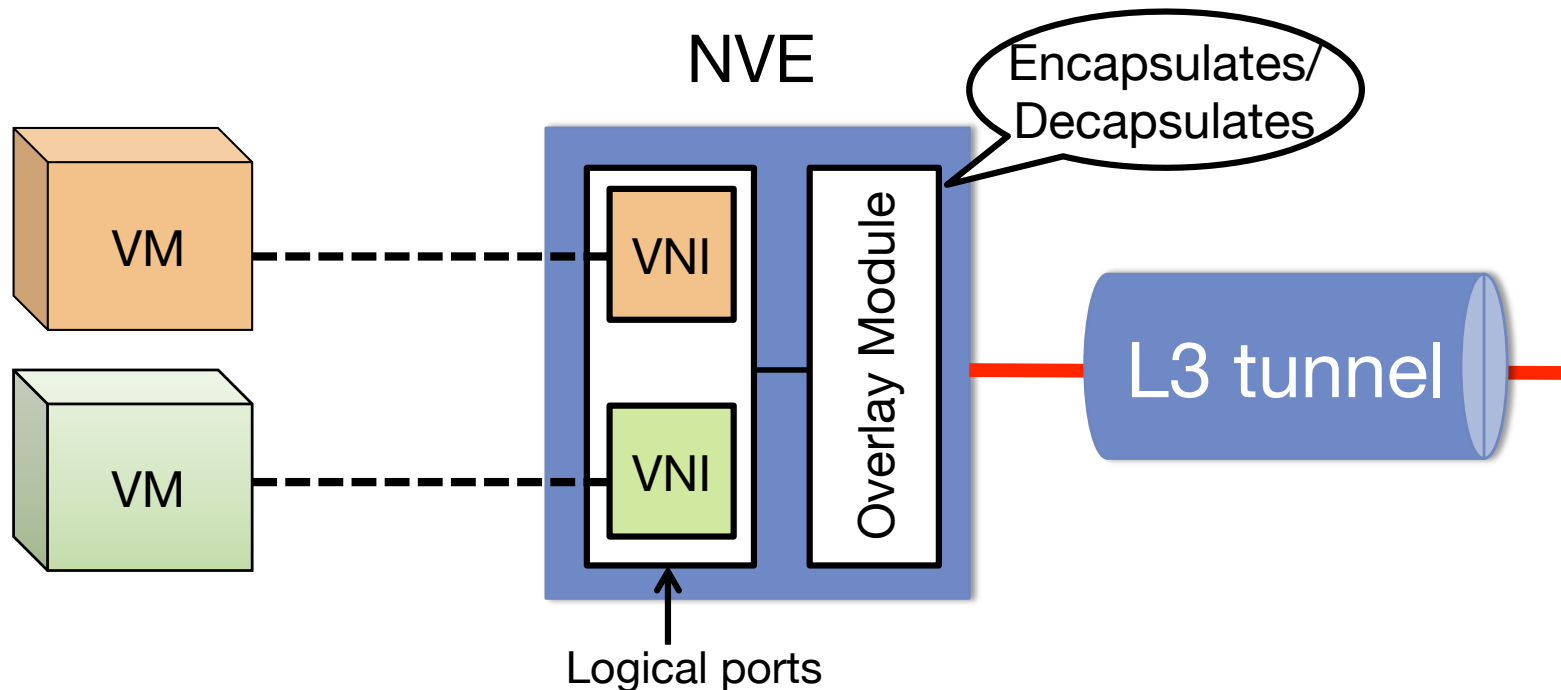- STT-comparable performance

❖ **Future work**

- Implementation of OVS-based SCLP
- Open source

Nagoya Institute of Technology    BOSCO Technologies
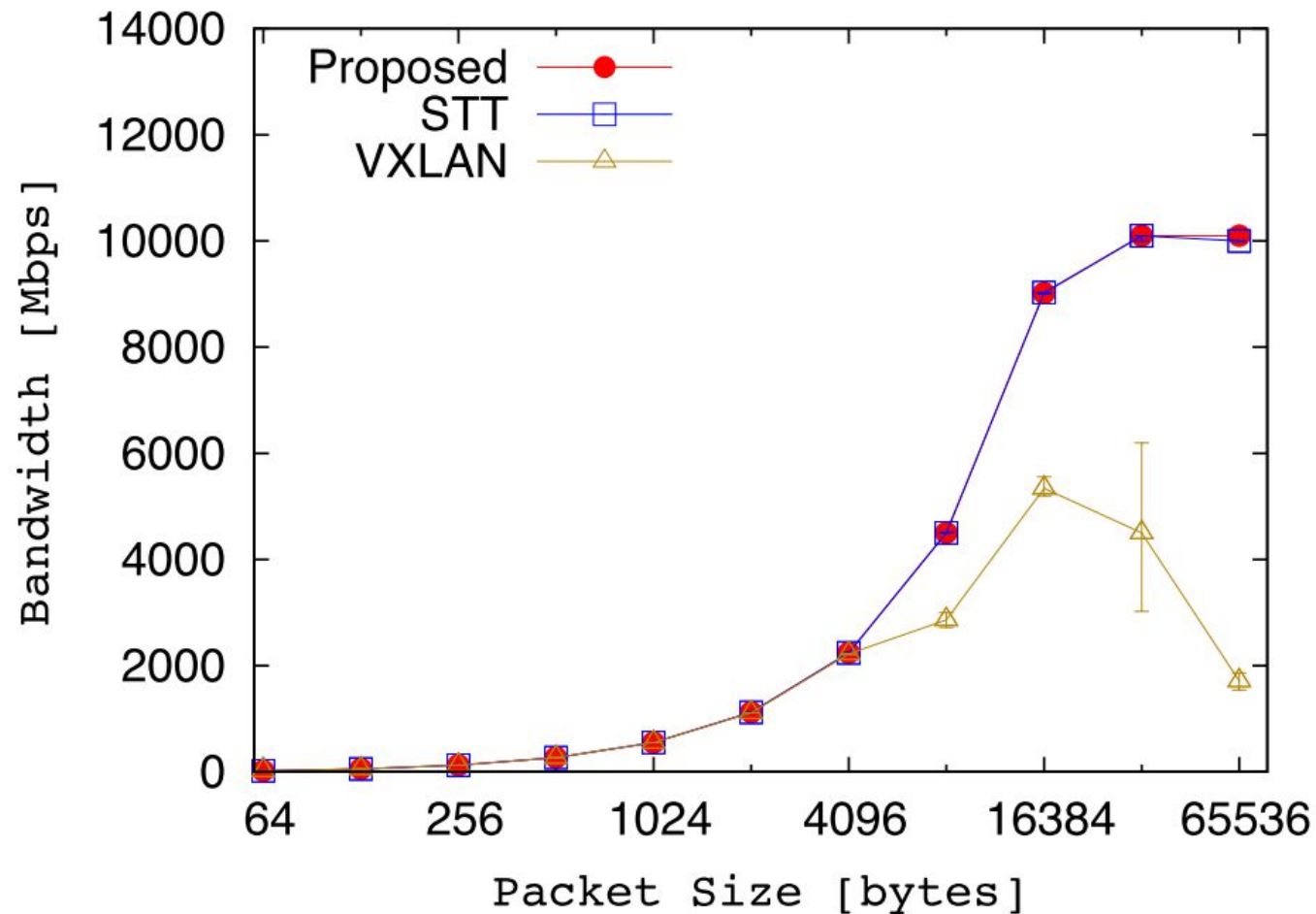
# NVE: Network Virtualization Edge

❖ **Tunnel End-Point**

- Physical switches
- Virtual switches
  - ➤ Open vSwitch (OVS), NSX switch, Hyper-V virtual switch

# Evaluation Results (UDP)

❖ **UDP communication**

Nagoya Institute of Technology    BOSCO Technologies

# Offloading Effects (STT)

| Offload | Tx / Rx | NIC / Kernel |
|---------|---------|--------------|
| TSO | Tx | NIC |
| GSO | Tx | Kernel |
| GRO | Rx | Kernel |



GRO effect !

Nagoya Institute of Technology    BOSCO Technologies