

# Scallop: Scalable Video Conferencing Using SDN Principles

Oliver Michel, <u>Satadal Sengupta</u>, Hyojoon Kim, Ravi Netravali, Jennifer Rexford







- Essential application across industries
- Explosive growth since 2020









- Essential application across industries
- Explosive growth since 2020









#### Increasing adoption



- Essential application across industries
- Explosive growth since 2020









#### Increasing adoption





#### Increasing expectations







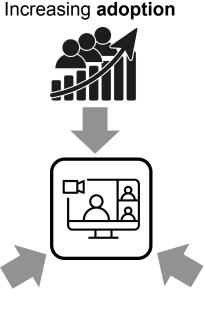
- Essential application across industries
- Explosive growth since 2020











Increasing expectations















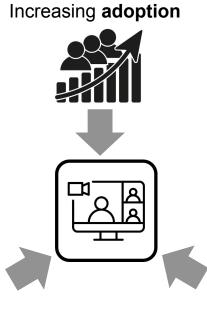
- Essential application across industries
- Explosive growth since 2020











Increasing expectations











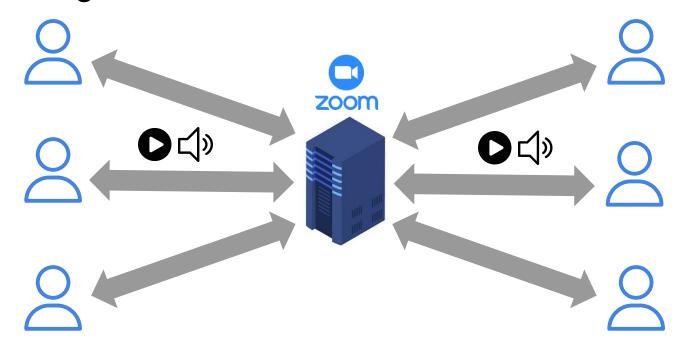


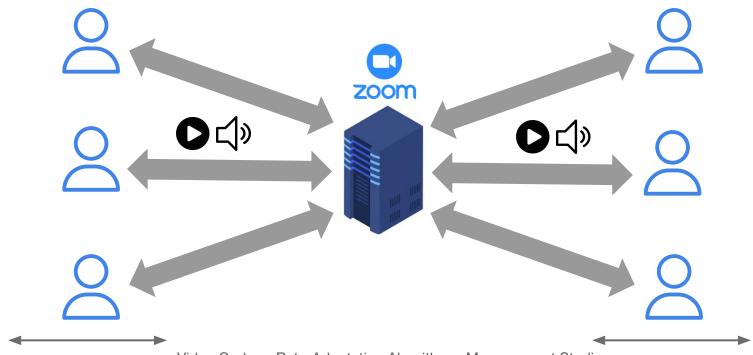


Challenge: Provide consistently high quality at scale

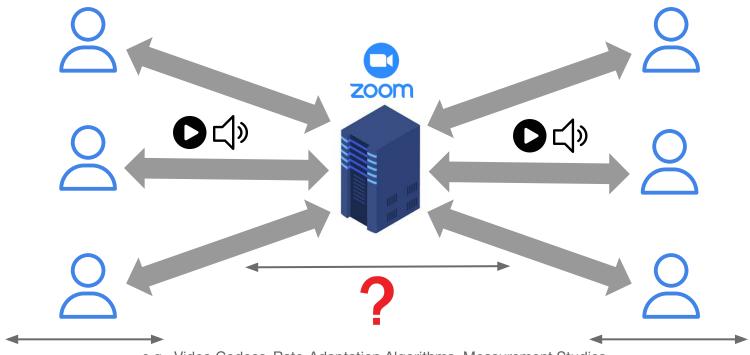




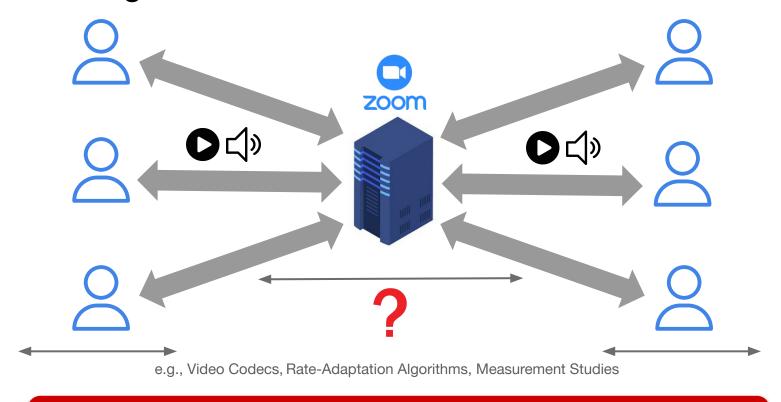




e.g., Video Codecs, Rate-Adaptation Algorithms, Measurement Studies



e.g., Video Codecs, Rate-Adaptation Algorithms, Measurement Studies



In this study: The application operators' perspective









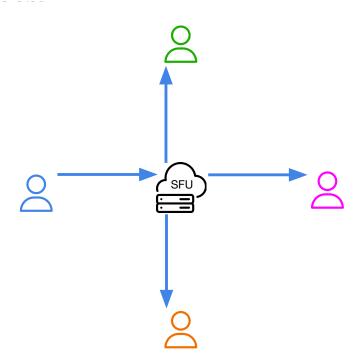










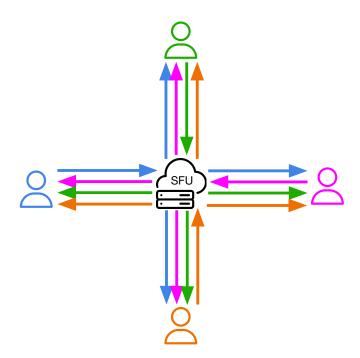


#### SFU roles:

(1) Relay audio and video streams





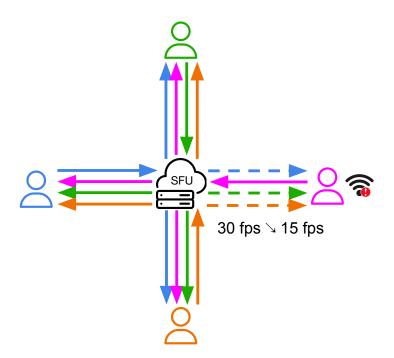


#### SFU roles:

(1) Relay audio and video streams





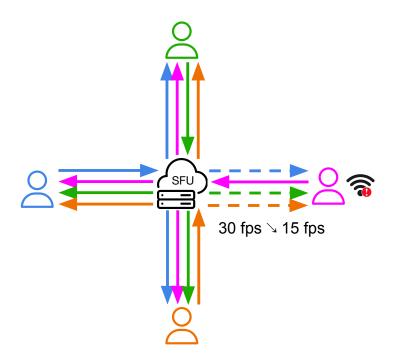


#### SFU roles:

- (1) Relay audio and video streams
- (2) Monitor and adapt media signals







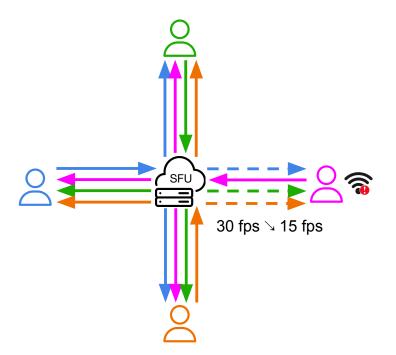
#### SFU roles:

- (1) Relay audio and video streams
- (2) Monitor and adapt media signals

#### SFUs hard to scale:







#### SFU roles:

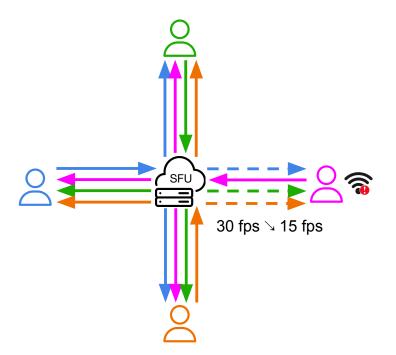
- (1) Relay audio and video streams
- (2) Monitor and adapt media signals

#### SFUs hard to scale:

(1) Workload hard to predict







#### SFU roles:

- (1) Relay audio and video streams
- (2) Monitor and adapt media signals

#### SFUs hard to scale:

- (1) Workload hard to predict
- (2) Quadratic scaling
  - $3 \rightarrow 4$  parts.  $\Rightarrow 9 \rightarrow 16$  streams

Dynamic, high-volume workload

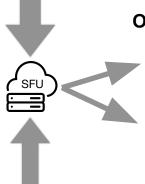


Dynamic, high-volume workload



Underprovisioning can affect quality massively

Dynamic, high-volume workload

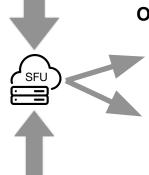


#### **Operators left with two options:**

- Massively over-provision
   → costly and wasteful
- Reactively autoscale
  - → risk harming QoE for users

Underprovisioning can affect quality massively

Dynamic, high-volume workload

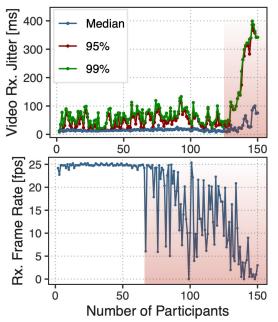


Underprovisioning can affect quality massively

#### Operators left with two options:

- Massively over-provision
   → costly and wasteful
- Reactively autoscale
  - → risk harming QoE for users

QoE in MediaSoup when increasing SFU load:

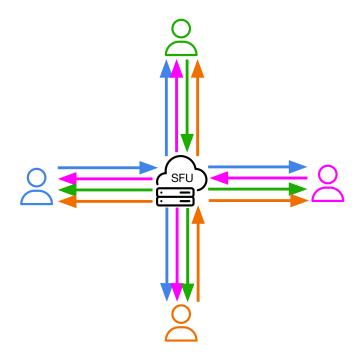




SFU operation is strikingly similar to traditional packet processing



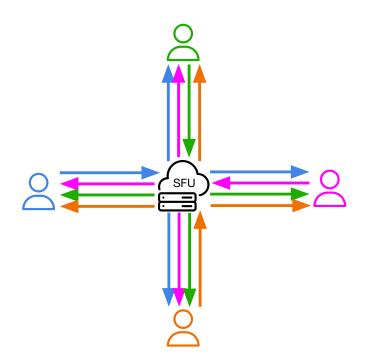
SFU operation is strikingly similar to traditional packet processing



(1) Relay audio and video streams



SFU operation is strikingly similar to traditional packet processing

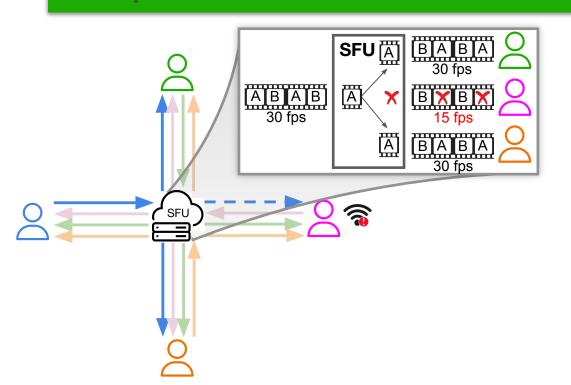


(1) Relay audio and video streams



Replicate traffic (Multicast)

SFU operation is strikingly similar to traditional packet processing

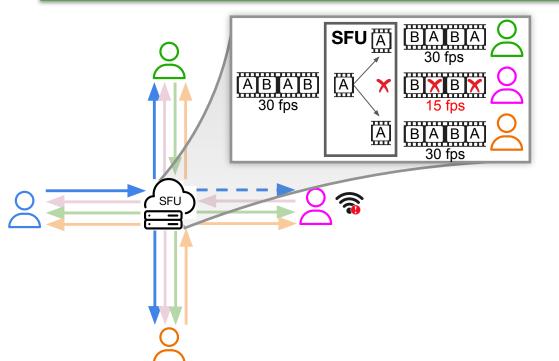


(1) Relay audio and video streams



(2) Monitor and adapt media signals

SFU operation is strikingly similar to traditional packet processing



(1) Relay audio and video streams



(2) Monitor and adapt media signals



Selectively forward traffic (Firewall)



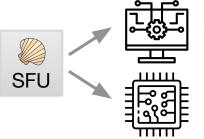




A novel hardware/software SFU co-design inspired by SDN



A novel hardware/software SFU co-design inspired by SDN



**Software:** Handle critical

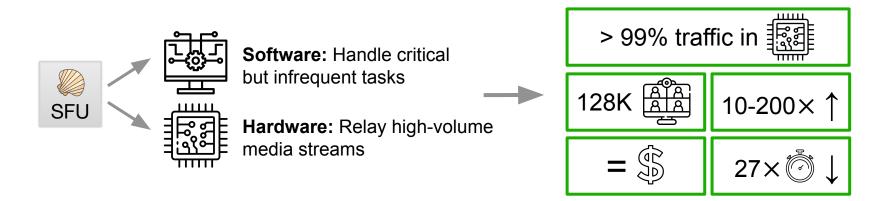
but infrequent tasks

Hardware: Relay high-volume

media streams



#### A novel hardware/software SFU co-design inspired by SDN



### Challenges

### Challenges

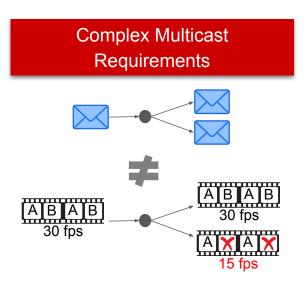
Monolithic Software Architecture



### Challenges

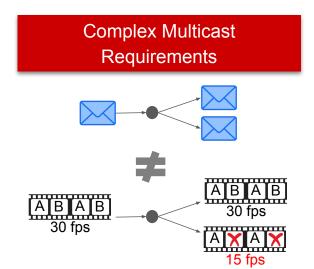
Monolithic Software Architecture





Monolithic
Software Architecture





Misaligned with Widely-Deployed Standard



Existing systems → split-proxy

Monolithic Software Architecture



Complex Multicast
Requirements

ABAB
30 fps

AXAX
15 fps

Misaligned with Widely-Deployed Standard

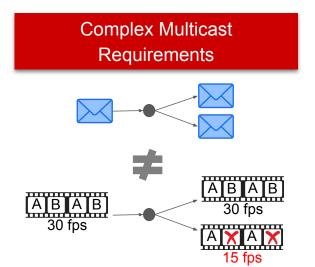


Existing systems → split-proxy

- Proxy architecture?
- Feedback semantics?
- Transparent rate adaptation?

Monolithic
Software Architecture





Misaligned with Widely-Deployed Standard



Existing systems → split-proxy

- Proxy architecture?
- Feedback semantics?
- Transparent rate adaptation?

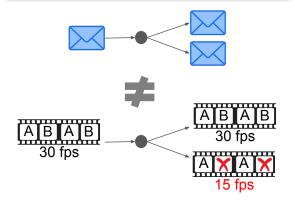
1

How to disaggregate into control and data planes?

Monolithic
Software Architecture



Complex Multicast Requirements



Misaligned with Widely-Deployed Standard



Existing systems → split-proxy

- Proxy architecture?
- Feedback semantics?
- Transparent rate adaptation?

1

How to disaggregate into control and data planes?

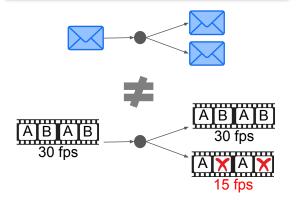
2

How to realize and scale application-layer multicast?

Monolithic
Software Architecture



Complex Multicast Requirements



Misaligned with Widely-Deployed Standard



Existing systems → split-proxy

- Proxy architecture?
- Feedback semantics?
- Transparent rate adaptation?

1

How to disaggregate into control and data planes?

2

How to realize and scale application-layer multicast?

3

How to make Scallop interoperable with WebRTC?



SFU workload is amenable to a control/data-plane split



SFU workload is amenable to a control/data-plane split

(1) >10ms latency, every few mins.

e.g., session management, signaling



SFU workload is amenable to a control/data-plane split

(1) >10ms latency, every few mins.

e.g., session management, signaling

(2) 1 < t < 10ms latency, 2-3 per sec.

e.g., handling feedback messages and connectivity checks



#### SFU workload is amenable to a control/data-plane split

(1) >10ms latency, every few mins.

e.g., session management, signaling

(2) 1 < t < 10ms latency, 2-3 per sec.

e.g., handling feedback messages and connectivity checks

(3) <1ms latency, 100s per sec.



SFU workload is amenable to a control/data-plane split

- (1) >10ms latency, every few mins.
- e.g., session management, signaling

- (2) 1 < t < 10ms latency, 2-3 per sec.
- e.g., handling feedback messages and connectivity checks
- (3) <1ms latency, 100s per sec.

replication and selective forwarding of media packets

requires lower latency



#### SFU workload is amenable to a control/data-plane split

- (1) >10ms latency, every few mins.
- e.g., session management, signaling

- (2) 1 < t < 10ms latency, 2-3 per sec.
- e.g., handling feedback messages and connectivity checks
- (3) <1ms latency, 100s per sec.

replication and selective forwarding of media packets

requires lower latency

higher event rate



#### SFU workload is amenable to a control/data-plane split

- (1) >10ms latency, every few mins.
- e.g., session management, signaling

- (2) 1 < t < 10ms latency, 2-3 per sec.
- e.g., handling feedback messages and connectivity checks
- (3) <1ms latency, 100s per sec.

replication and selective forwarding of media packets

requires lower latency

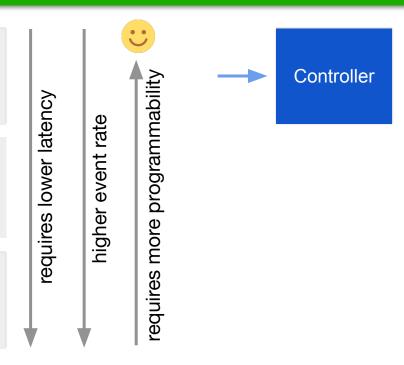
requires more programmability higher event rate



#### SFU workload is amenable to a control/data-plane split

- (1) >10ms latency, every few mins.
- e.g., session management, signaling

- (2) 1 < t < 10ms latency, 2-3 per sec.
- e.g., handling feedback messages and connectivity checks
- (3) <1ms latency, 100s per sec.

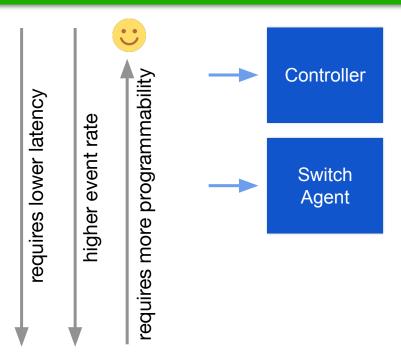




#### SFU workload is amenable to a control/data-plane split

- (1) >10ms latency, every few mins.
- e.g., session management, signaling

- (2) 1 < t < 10ms latency, 2-3 per sec.
- e.g., handling feedback messages and connectivity checks
- (3) <1ms latency, 100s per sec.





#### SFU workload is amenable to a control/data-plane split

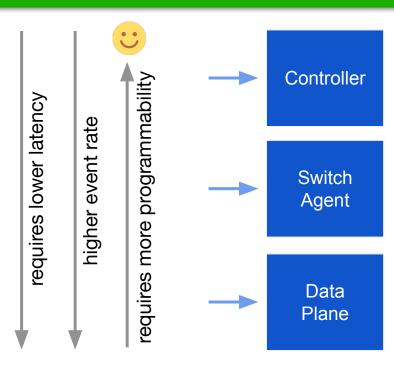
(1) >10ms latency, every few mins.

e.g., session management, signaling

(2) 1 < t < 10ms latency, 2-3 per sec.

e.g., handling feedback messages and connectivity checks

(3) <1ms latency, 100s per sec.





#### SFU workload is amenable to a control/data-plane split

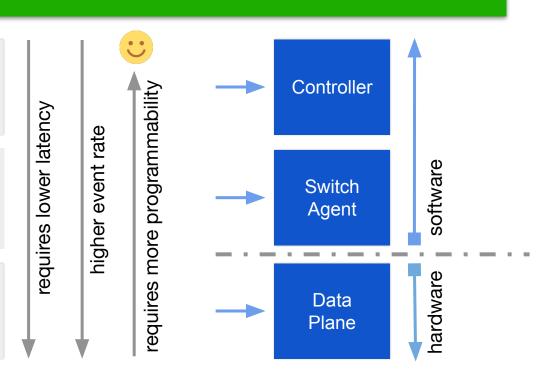
(1) >10ms latency, every few mins.

e.g., session management, signaling

(2) 1 < t < 10ms latency, 2-3 per sec.

e.g., handling feedback messages and connectivity checks

(3) <1ms latency, 100s per sec.





#### SFU workload is amenable to a control/data-plane split

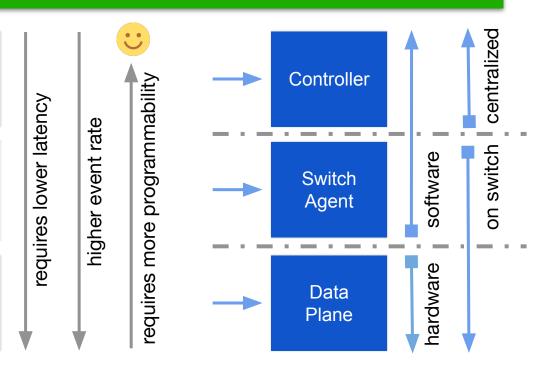
(1) >10ms latency, every few mins.

e.g., session management, signaling

(2) 1 < t < 10ms latency, 2-3 per sec.

e.g., handling feedback messages and connectivity checks

(3) <1ms latency, 100s per sec.





Hardware-native packet copying capabilities can be leveraged for SFU-style replication



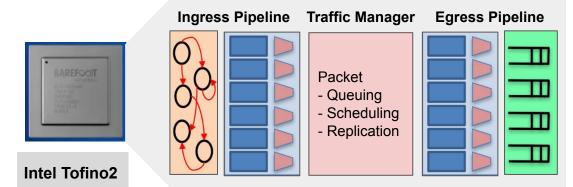
Hardware-native packet copying capabilities can be leveraged for SFU-style replication

Programmable Switches



Hardware-native packet copying capabilities can be leveraged for SFU-style replication

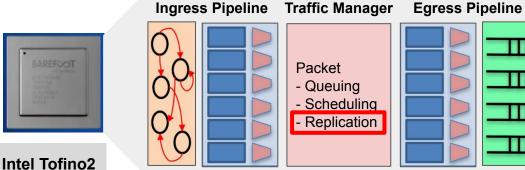
#### Programmable Switches





Hardware-native packet copying capabilities can be leveraged for SFU-style replication

#### Programmable Switches





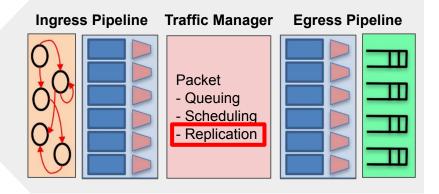
Hardware-native packet copying capabilities can be leveraged for SFU-style replication

#### Programmable Switches

SmartNICs/DPUs



**Intel Tofino2** 



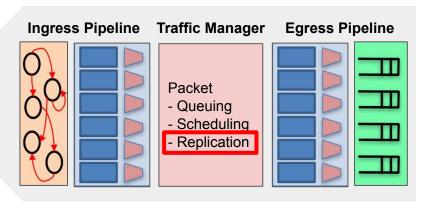


Hardware-native packet copying capabilities can be leveraged for SFU-style replication

#### Programmable Switches



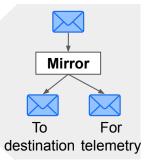
**Intel Tofino2** 



#### SmartNICs/DPUs







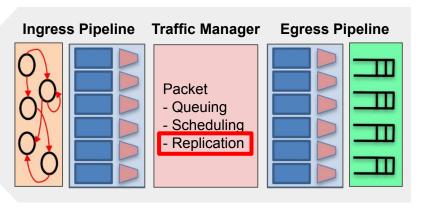


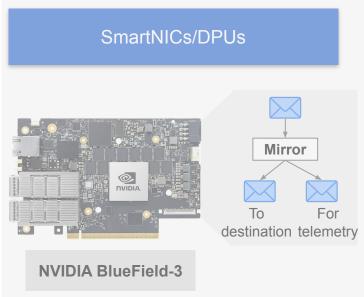
Hardware-native packet copying capabilities can be leveraged for SFU-style replication

#### Programmable Switches

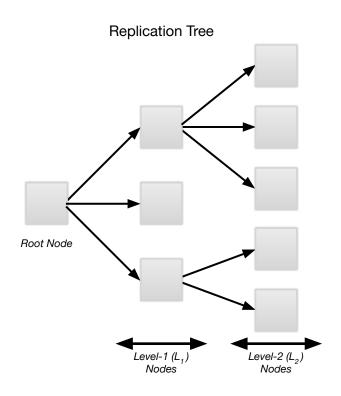


**Intel Tofino2** 



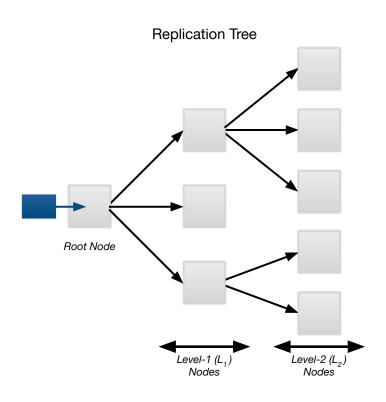


Background on Tofino's Packet Replication Engine (PRE)



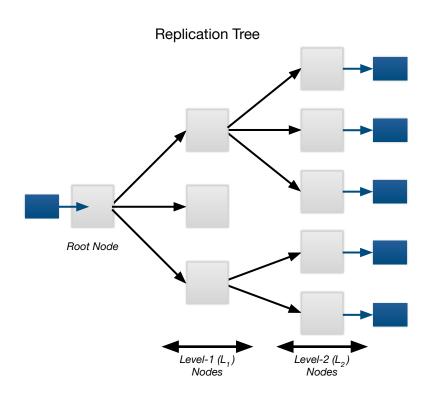
 Abstraction: replication tree with level-1 nodes and level-2 nodes

Background on Tofino's Packet Replication Engine (PRE)

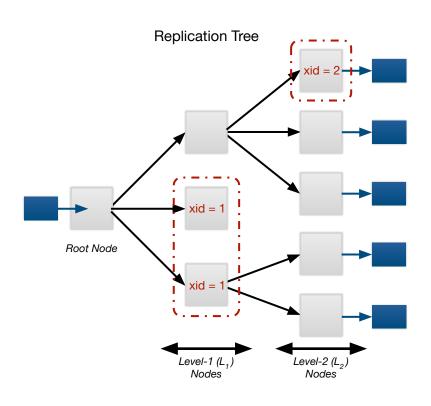


 Abstraction: replication tree with level-1 nodes and level-2 nodes

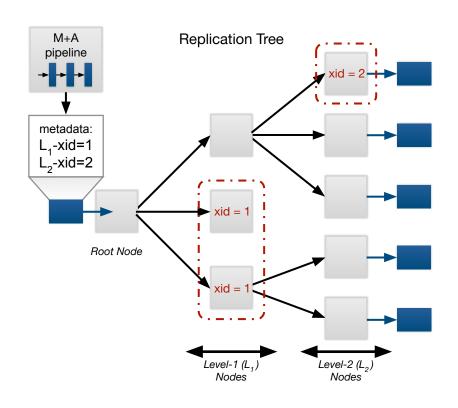
Background on Tofino's Packet Replication Engine (PRE)



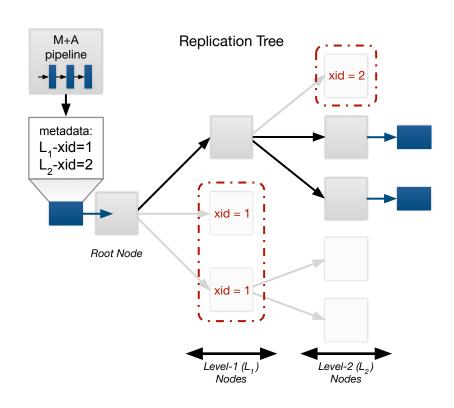
 Abstraction: replication tree with level-1 nodes and level-2 nodes



- Abstraction: replication tree with level-1 nodes and level-2 nodes
- Supports dynamic tree pruning
  - Each node can be associated with an exclusion ID (xid)



- Abstraction: replication tree with level-1 nodes and level-2 nodes
- Supports dynamic tree pruning
  - Each node can be associated with an exclusion ID (xid)
  - The ingress match+action pipeline can associate individual packets with xids



- Abstraction: replication tree with level-1 nodes and level-2 nodes
- Supports dynamic tree pruning
  - Each node can be associated with an exclusion ID (xid)
  - The ingress match+action pipeline can associate individual packets with xids
  - No replication along edges leading to excluded nodes

Packet Replication in Scallop: Challenges

Packet Replication in Scallop: Challenges

#### PRE-to-VC Mapping?

- PRE entities:
  - Root, L1 + L2 nodes
  - L1 + L2 xids
- VC entities:
  - Meetings, participants
  - Quality layers

Packet Replication in Scallop: Challenges

#### PRE-to-VC Mapping?

- PRE entities:
  - Root, L1 + L2 nodes
  - L1 + L2 xids
- VC entities:
  - Meetings, participants
  - Quality layers

#### **Different Meeting Configurations**

- Rate adaptation status
- Two-party vs. multiparty
- Can change dynamically

Packet Replication in Scallop: Challenges

#### PRE-to-VC Mapping?

- PRE entities:
  - Root, L1 + L2 nodes
  - L1 + L2 xids
- VC entities:
  - Meetings, participants
  - Quality layers

#### **Different Meeting Configurations**

- Rate adaptation status
- Two-party vs. multiparty
- Can change dynamically

#### Limited Resources

- 64,000 replication trees
- 2<sup>24</sup> L1 nodes

Packet Replication in Scallop: Challenges

#### PRE-to-VC Mapping?

- PRE entities:
  - Root, L1 + L2 nodes
  - L1 + L2 xids
- VC entities:
  - Meetings, participants
  - Quality layers

#### **Different Meeting Configurations**

- Rate adaptation status
- Two-party vs. multiparty
- Can change dynamically

#### Limited Resources

- 64,000 replication trees
- 2<sup>24</sup> L1 nodes



How to (i) correctly and (ii) efficiently map VC entities to PRE entities for each meeting configuration?

Packet Replication in Scallop: Solution

- Non-Rate-Adapted (NRA)
- Rate-Adapted (RA)
  - Receiver (RA-R)
  - Sender-Receiver (RA-SR)
- Two-Party (2P)

Packet Replication in Scallop: Solution

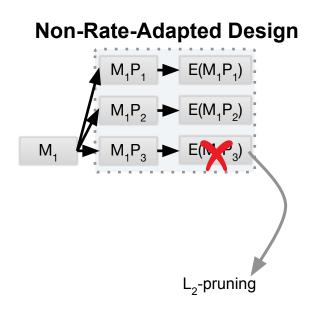
# Dynamic Migration

- Non-Rate-Adapted (NRA)
- Rate-Adapted (RA)
  - Receiver (RA-R)
  - Sender-Receiver (RA-SR)
- Two-Party (2P)

Packet Replication in Scallop: Solution

# Dynamic Migration

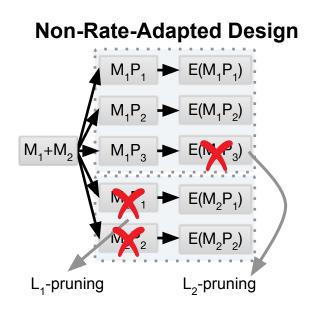
- Non-Rate-Adapted (NRA)
- Rate-Adapted (RA)
  - Receiver (RA-R)
  - Sender-Receiver (RA-SR)
- Two-Party (2P)



Packet Replication in Scallop: Solution

# Dynamic Migration

- Non-Rate-Adapted (NRA)
- Rate-Adapted (RA)
  - Receiver (RA-R)
  - Sender-Receiver (RA-SR)
- Two-Party (2P)



Packet Replication in Scallop: Solution

# Dynamic Migration

#### **Optimal Designs**

- Non-Rate-Adapted (NRA)
- Rate-Adapted (RA)
  - Receiver (RA-R)
  - Sender-Receiver (RA-SR)
- Two-Party (2P)

## Non-Rate-Adapted Design $M_1+M_2$

Supports up to 128K concurrent meetings and 2<sup>24</sup> concurrent participants

L<sub>2</sub>-pruning

L₁-pruning





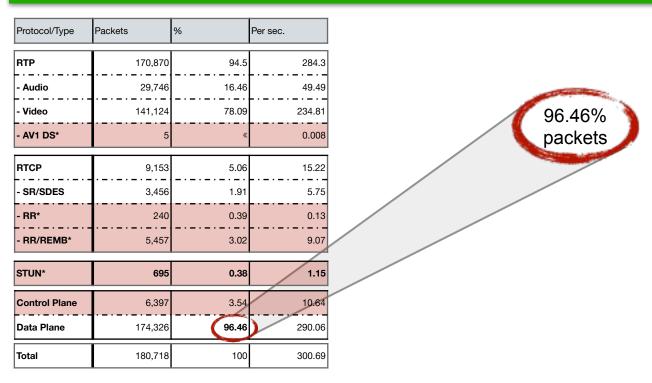


Protocol/Type	Packets	%	Per sec.
RTP	170,870	94.5	284.3
- Audio	29,746	16.46	49.49
- Video	141,124	78.09	234.81
- AV1 DS*	5	<u> </u>	0.008
RTCP	9,153	5.06	15.22
- SR/SDES	3,456	1.91	5.75
- RR*	240	0.39	0.13
- RR/REMB*	5,457	3.02	9.07
STUN*	695	0.38	1.15

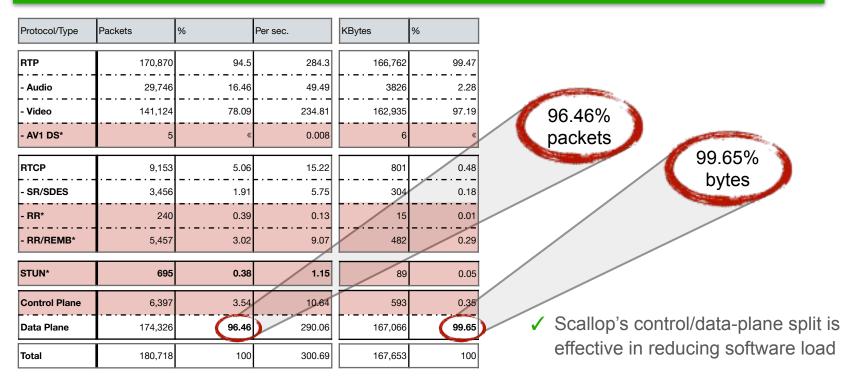


Protocol/Type	Packets	%	Per sec.
RTP	170,870	94.5	284.3
- Audio	29,746	16.46	49.49
- Video	141,124	78.09	234.81
- AV1 DS*	5	<b>«</b>	0.008
RTCP	9,153	5.06	15.22
- SR/SDES	3,456	1.91	5.75
- RR*	240	0.39	0.13
- RR/REMB*	5,457	3.02	9.07
STUN*	695	0.38	1.15



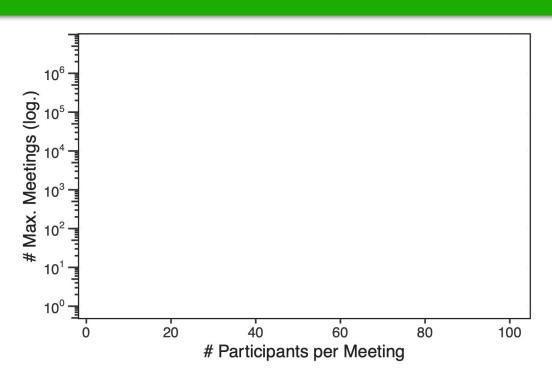




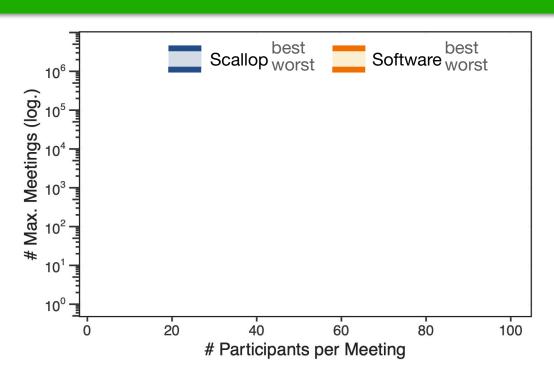




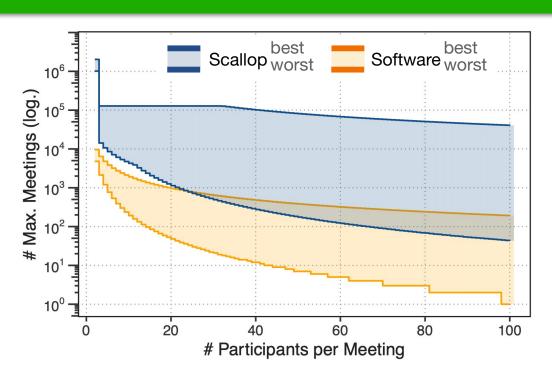




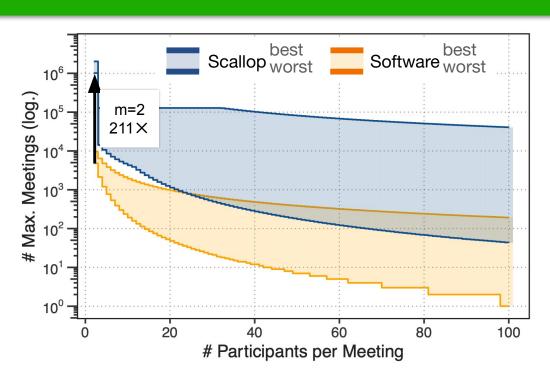




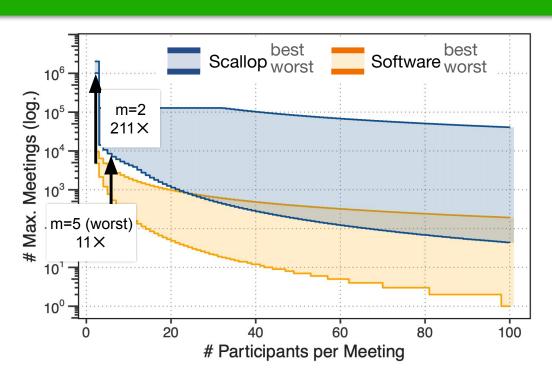




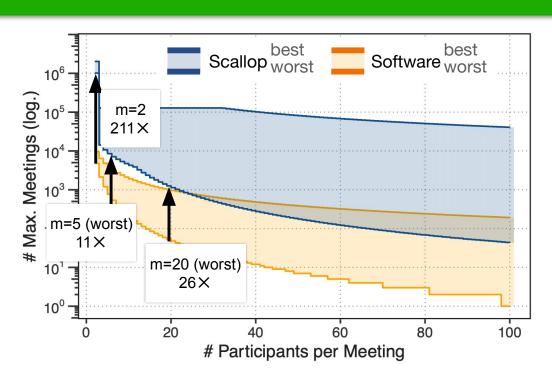




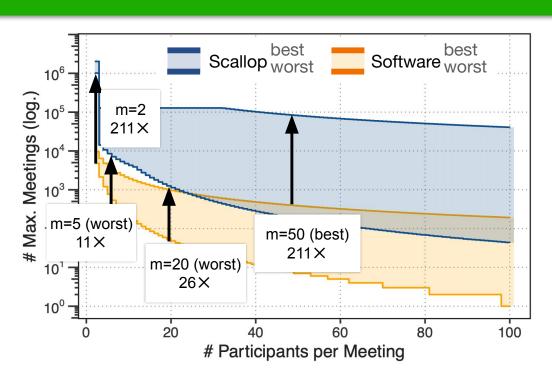






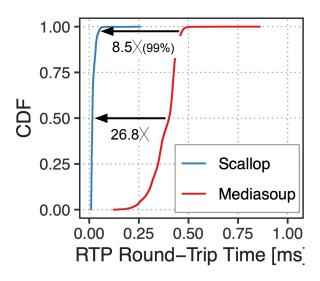








Scallop reduces SFU-induced latency: median by 27 × and tail by 8 ×



#### **Conclusion and Artifacts**

• Scallop: Novel, hardware-software co-designed SFU

#### Conclusion and Artifacts

Scallop: Novel, hardware-software co-designed SFU



#### Artifacts on &



- Control plane + software model of data plane
- Hardware prototypes of data plane:



Wireshark plugin

### Thank You!

- Code: <a href="https://github.com/Princeton-Cabernet/Scallop">https://github.com/Princeton-Cabernet/Scallop</a>
- Contact: <u>satadal.sengupta@princeton.edu</u>

I'm on the market for academic positions in the US, Canada, and Europe

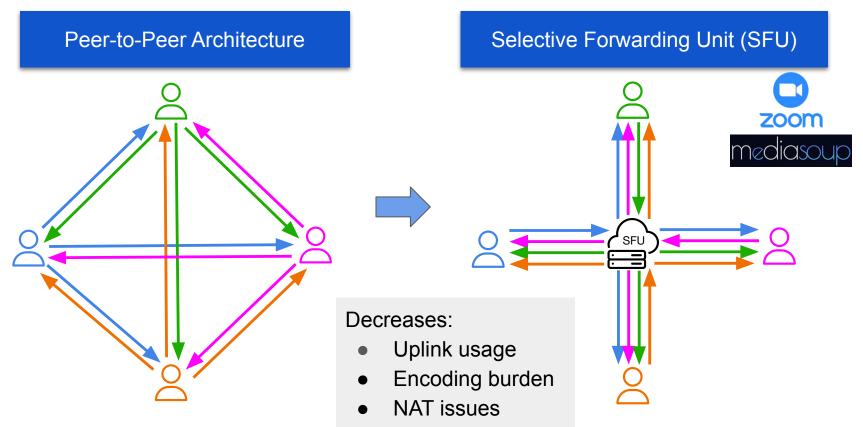






### Backup

#### Video-Conferencing Infrastructure



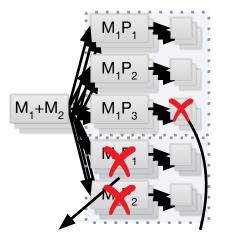
#### Scalable Semantics-Aware Replication

#### Packet Replication in Scallop

Straw Man one tree per quality laver Intel Tofino2: • 64,000 replication trees ( $\tau$ ) • 2<sup>24</sup> L<sub>4</sub> nodes (v) Scallop Prototype: • 3 quality layers (κ) Maximize: • N # of participants • M # of meetings

$$\begin{split} N &= \tau/\kappa = 64,000/3 \approx 21,333 \\ M &\leq N \end{split}$$

#### Optimization



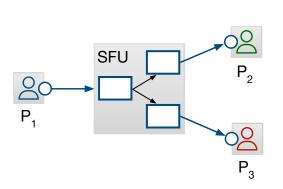
L<sub>1</sub>-pruning to diff. between meetings (max. 2 per tree) L<sub>2</sub>-pruning to ensure a participant does not receive their own packets

 $N \le v$  $M = 2\tau/\kappa = 2 \times 64,000/3 \approx 42,666$ 

#### Other Designs

- Non-Rate-Adapted (NRA)
- Rate-adapted/Receiver (RA-R)
- Rate-adapted/Sender-Receiver (RA-SR)
- 2-Party (2P)

#### Interoperability with WebRTC



#### **PC Endpoints**

- congestion control
- encryption/decryption
- packetization/de-packetization

  SFU
  P<sub>2</sub>

  WebRTC

  RTCPeerConnection (PC)

  P<sub>2</sub>

#### **Proxy SFU Architecture**



- hardware-friendly
- low overhead at SFU
- latency-friendly

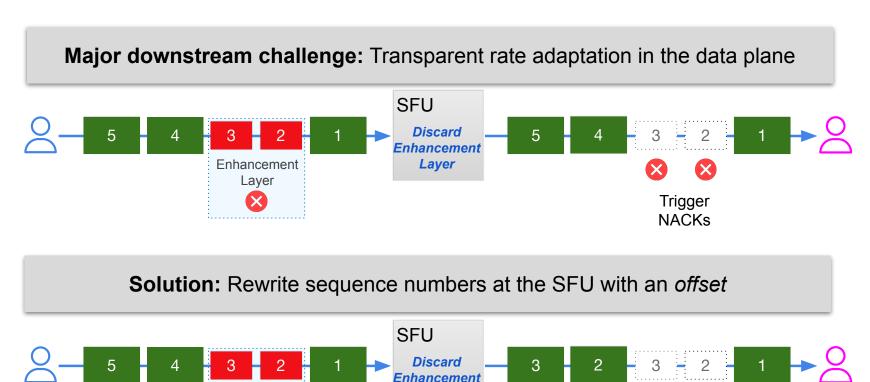
#### **Split-Proxy SFU Architecture**

- difficult in hardware
- lots of replicated logic at SFU
- introduces latency

### Interoperability with WebRTC

Enhancement

Layer



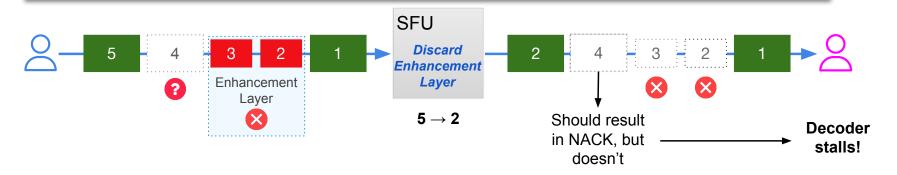
Layer

 $4 \rightarrow 2, 5 \rightarrow 3$ 

No NACKs!

### Interoperability with WebRTC

#### Challenge: Naive rewriting causes video freeze during network loss

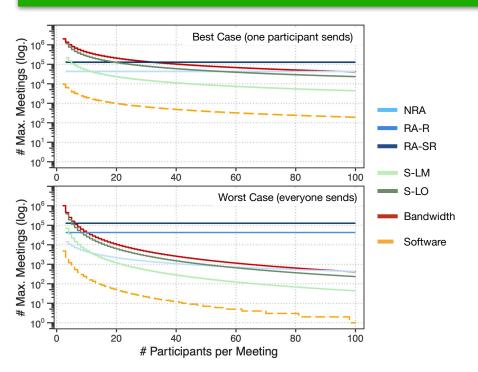


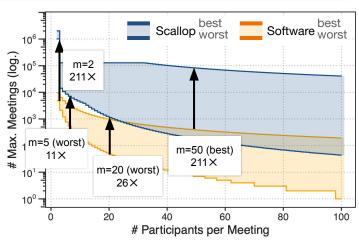
Observation: When unsure, leaving gap better than hiding one

Hardware-friendly heuristic that never hides loss at cost of possibly unnecessary retxs.

Two variations based on trade-off between unnecessary retxs. and switch memory (S-LO, S-LM)

### **(**





- Scale improvement depends on meeting composition and rate-adaptation characteristics
- Scallop improves scalability 7 x to 211 x over software and always performs better than software (for a given meeting composition)