

A circular logo with a blue border and a light blue background. Inside the circle, the text "P4 Expert Roundtable Series" is written in blue. The "P4" is in a larger, stylized font with a purple 'P' and a green '4'. Below "P4" is the word "Expert" in a smaller font. Below "Expert" is the word "Roundtable Series" in a larger font. To the right of "Roundtable Series" is a small white cartoon animal. Below the text is the date "April 28-29, 2020". Below the date is the text "Hosted by:" followed by the ONF logo, which consists of a red ampersand and the letters "ONF" in blue.

P4
Expert
Roundtable Series

April 28-29, 2020

Hosted by:



Offloading Media Traffic to P4 Programmable Data Plane Switches

Elie Kfoury¹ (Ph.D. Student), Jorge Crichigno¹, Elias Bou-Harb²

¹ University of South Carolina, Columbia, SC

² University of Texas, San Antonio, TX

Supported by the NSF awards 1925484 and 1829698

Special thanks to Vladimir Gurevich (Barefoot Networks, an Intel Company) for his insightful feedback on various technical issues.

Agenda

- Introduction
- Background Information
 - Session Initiation Protocol (SIP) and Real Time Protocol (RTP)
 - Network Address Translation (NAT) traversal problem
 - P4 switches
- Proposed solution
- Evaluation
- Lessons learned

Introduction

- According to estimations, media traffic represents approximately 80% of the total traffic over the Internet¹
- Much media traffic is generated by end users communicating with each other
- Media services (voice, video) running alongside the data network in campuses are becoming standard

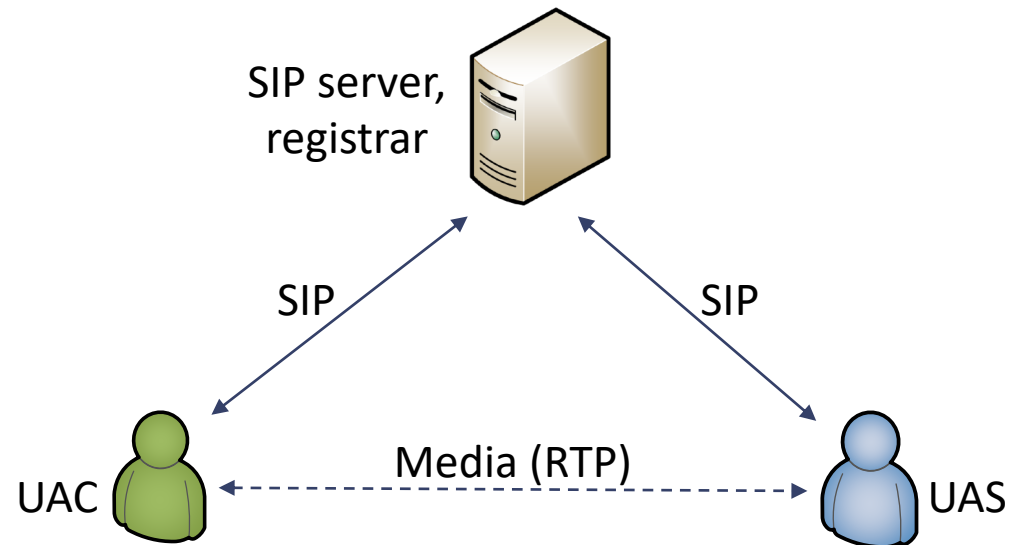
¹ H. W. Barz and G. A. Bassett, Multimedia networks: protocols, design and applications, John Wiley and Sons, 2016.

Voice and Video

- Conversational Voice- and Video-over-IP are widely used today
 - Open and proprietary solutions
- Supporting protocols are divided into two main categories
 - Session control protocols (signaling): establish and manage the session
 - E.g., Session Initiation Protocol (SIP)
- Media protocols (media)
 - Transfer audio and video streams between the end-users
 - E.g., Real Time Protocol (RTP)
- Desirable Quality-of-Service (QoS) characteristics
 - Delay- and jitter-sensitive, low values
 - Occasional losses are tolerated

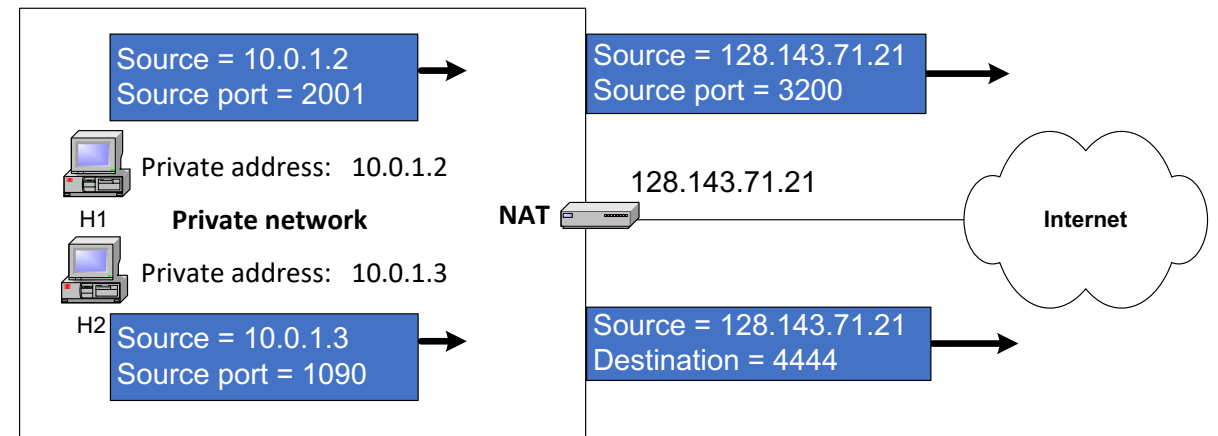
Signaling and Media Protocols

- SIP initiates, maintains, and terminates multimedia sessions between endpoints
 - User agent client (UAC)
 - User agent server (UAS)
- RTP transports real-time data, such as audio and video



Network Address Translation (NAT)

- NAT maps ports, private IP addresses to public IP addresses
 - Used in campus / enterprise networks, operators¹
- NAT introduces various issues
 - Violation of the end-to-end principle
 - Traversal of end-to-end sessions



¹I. Livadariu et al., "Inferring carrier-grade NAT deployment in the wild," in IEEE 2018 INFOCOM, 2018.

Network Address Translation (NAT)

- NAT prevents a user from outside from initiating a session
- If both users have NATs, then neither can accept a call
 - IP translation is recorded by a SIP registrar server
- SIP carries the IP addresses and ports to be used by RTP to send/receive media
- Ports are unknown until RTP traffic starts
- Several solutions proposed for NAT traversal
 - STUN - RFC 53891, TURN - RFC 75662, ICE - RFC 84453

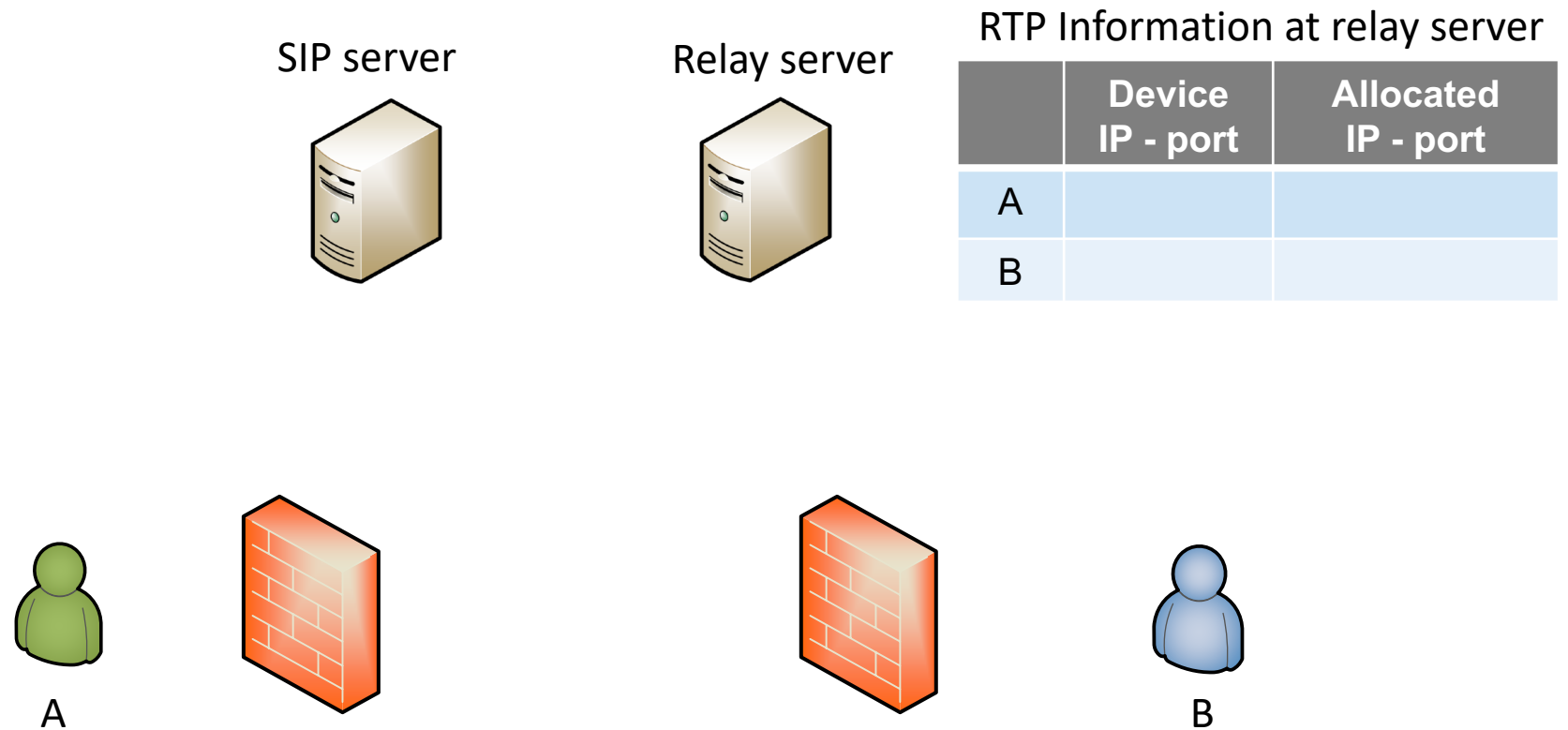
¹D. Wing, P. Matthews, R. Mahy, and J. Rosenberg, "RFC 5389 - STUN: Session traversal utilities for NAT," 2008.

²M. Petit-Huguenin, S. Nandakumar, G. Salgueiro, and P. Jones, "RFC 7566 - TURN: Traversal using relays around NAT (TURN) uniform resource identifiers," 2013.

³J. Rosenberg and C. Holmberg, "RFC 8445 - ICE: Interactive connectivity establishment: a protocol for Network Address Translator (NAT) traversal," 2018.

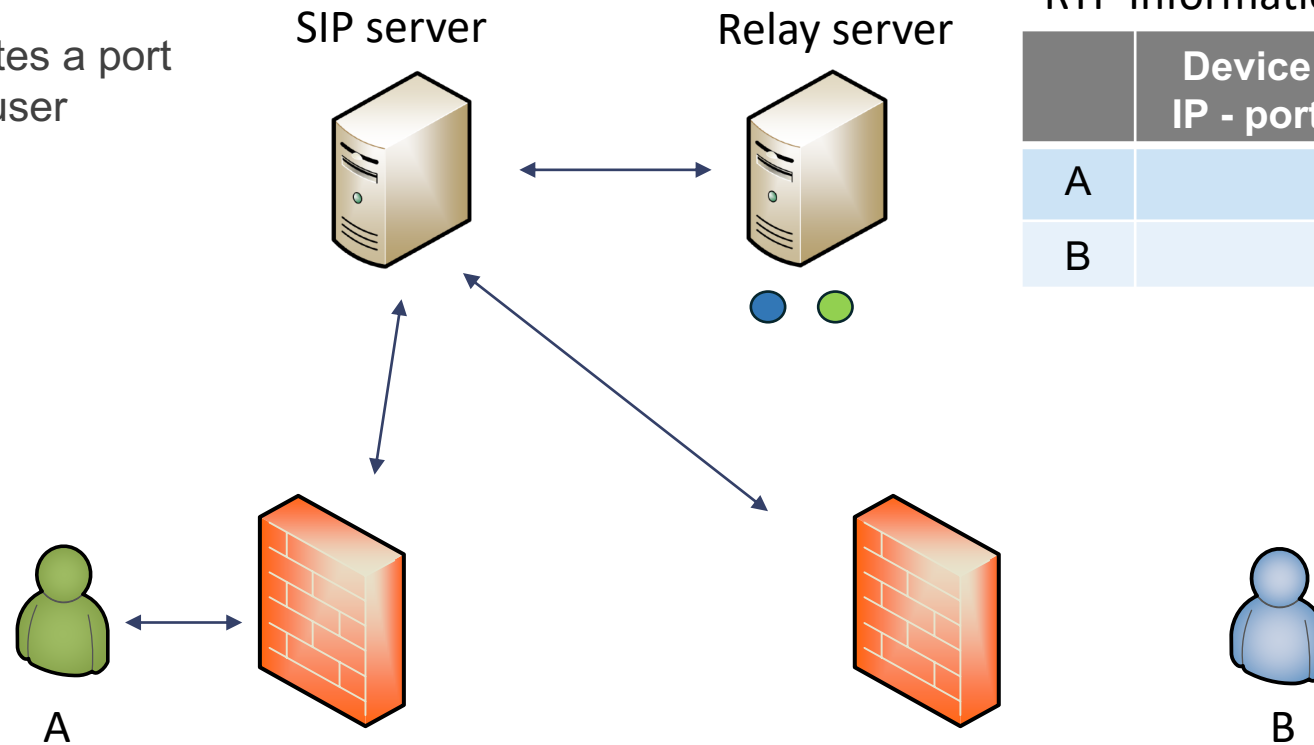
Relay Server for Media Traffic

- Intermediary device



Relay Server for Media Traffic

- Intermediary device
- SIP establishes the session
 - RTP ports are unknown
 - The relay server allocates a port on behalf of each end user

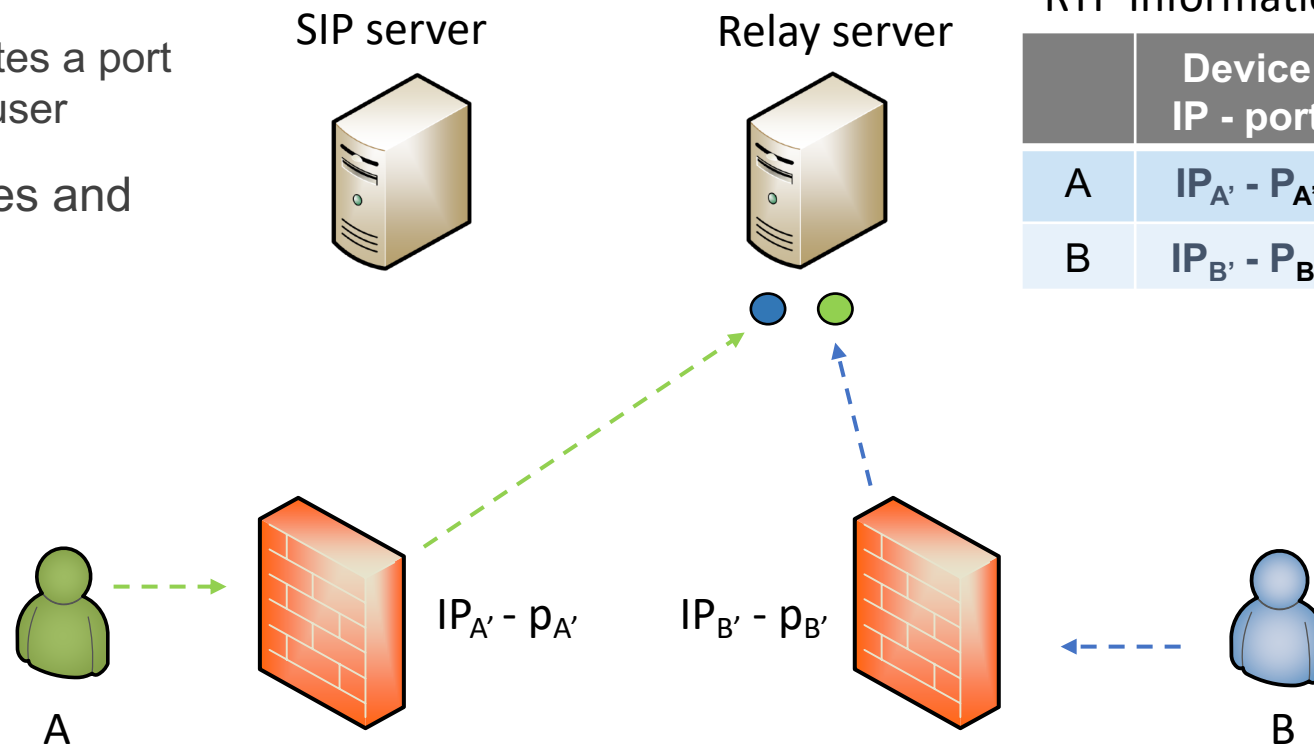


RTP Information at relay server

	Device IP - port	Allocated IP - port
A		$IP_R - P_{RA}$
B		$IP_R - P_{RB}$

Relay Server for Media Traffic

- Intermediary device
- SIP establishes the session
 - RTP ports are unknown
 - The relay server allocates a port on behalf of each end user
- The relay server receives and relays the RTP traffic

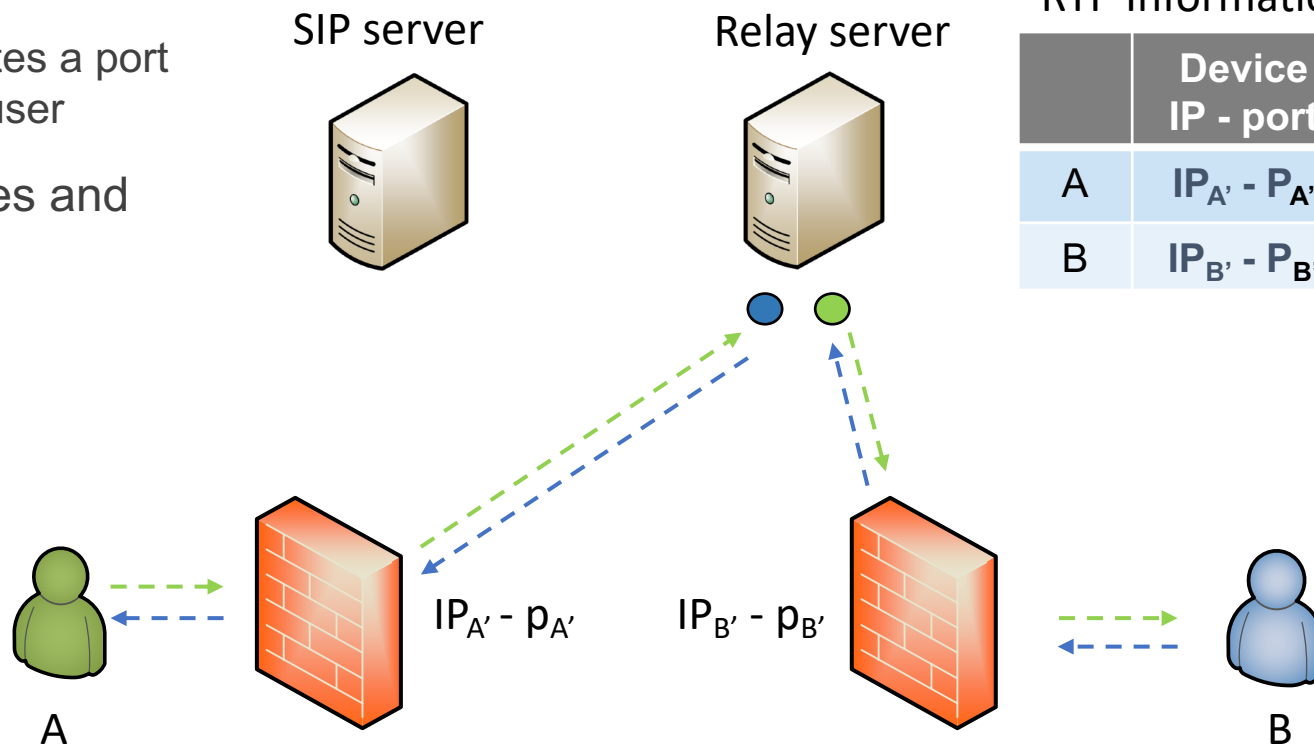


RTP Information at relay server

	Device IP - port	Allocated IP - port
A	IP _{A'} - P _{A'}	IP _R - P _{RA}
B	IP _{B'} - P _{B'}	IP _R - P _{RB}

Relay Server for Media Traffic

- Intermediary device
- SIP establishes the session
 - RTP ports are unknown
 - The relay server allocates a port on behalf of each end user
- The relay server receives and relays the RTP traffic



RTP Information at relay server

	Device IP - port	Allocated IP - port
A	IP _{A'} - P _{A'}	IP _R - P _{RA}
B	IP _{B'} - P _{B'}	IP _R - P _{RB}

Overview P4 Switches

- P4 switches permit programmer to program the data plane
- Add proprietary features; e.g., emulate RTP relay server
 - Parse packet headers, including UDP packets carrying RTP traffic
 - Header inspection, identifying media sessions using the 5-tuple
 - Modify fields, IP addresses and ports
- If the P4 program compiles, it runs on the chip at line rate

```
136 /*****  
137 *****/  
138 # PARSER  
139 /*****  
140 #  
141 # state parse_ethernet {  
142 #     packet.extract(hdr.ethernet);  
143 #     transition select(hdr.ethernet.etherType) {  
144 #         TYPE_IPV4: parse_ipv4;  
145 #         default: accept;  
146 #     }  
147 # }  
148 # state parse_ipv4 {  
149 #     packet.extract(hdr.ipv4);  
150 #     verify(hdr.ipv4.ihl >= 5, error.IPHeaderTooShort);  
151 #     transition select(hdr.ipv4.ihl) {  
152 #         5 : accept;  
153 #         default : parse_ipv4_option;  
154 #     }  
155 # }
```

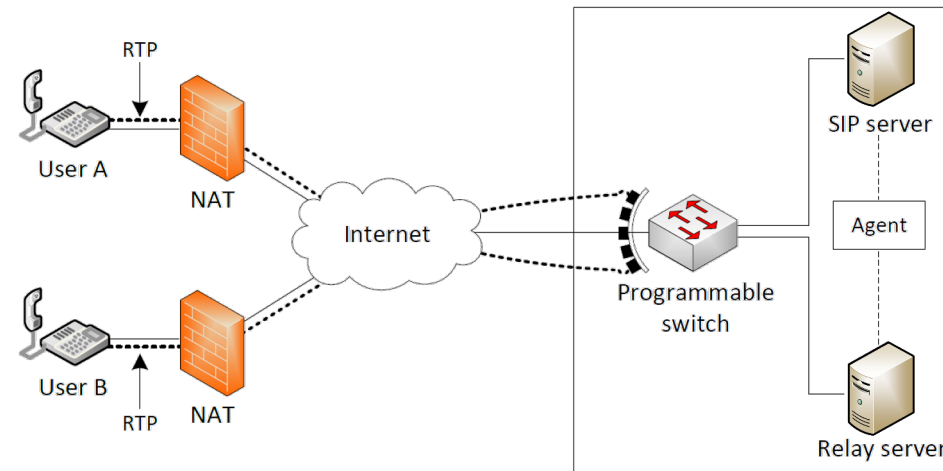
P4 code



Programmable chip

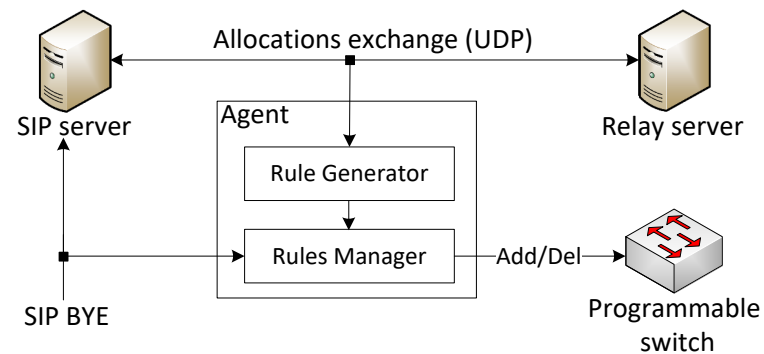
Proposed System

- Emulate the behavior of the relay server using programmable switch:
 1. Parse the incoming packet carrying media traffic from the first party, say user A
 2. Identify the session this packet belongs to by using the 5-tuple
 3. Replace the source IP with that of the relay server, and the source port with that used by the relay server to receive traffic from user A
 4. Replace the destination IP and the destination port with those of user B
 5. Recalculate both IPv4 and UDP checksums
 6. Forward the packet to user B



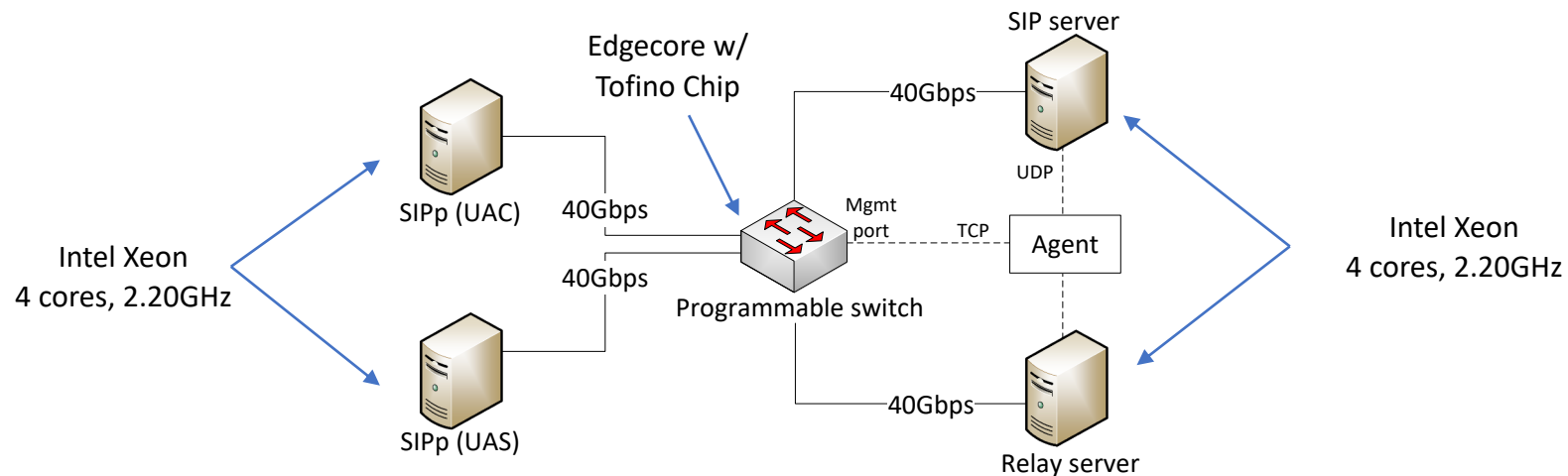
Proposed System

- A custom software (agent) learns the ports allocated to a media session by the relay server
- The Rule Generator uses the 5-tuple allocated to the media session to construct a unique session identifier
- It stores identifiers of the media sessions and the new header' values in the switch
- It also clears media sessions allocated in the switch when a call is teared down



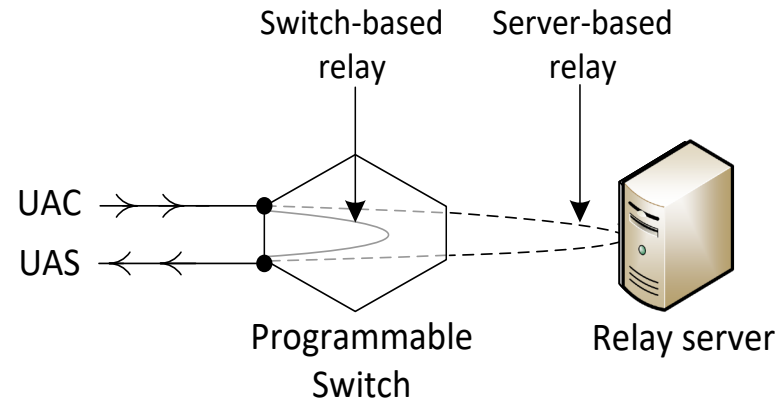
Implementation and Evaluation

- OpenSIPS, an open source implementation of a SIP server
- RTPProxy, a high-performance relay server for RTP streams
- SIPp: an open source SIP traffic generator that can establish multiple concurrent sessions and generate media (RTP) traffic
- Iperf3: traffic generator used to generate background UDP traffic
- Edgecore Wedge100BF-32X: programmable switch



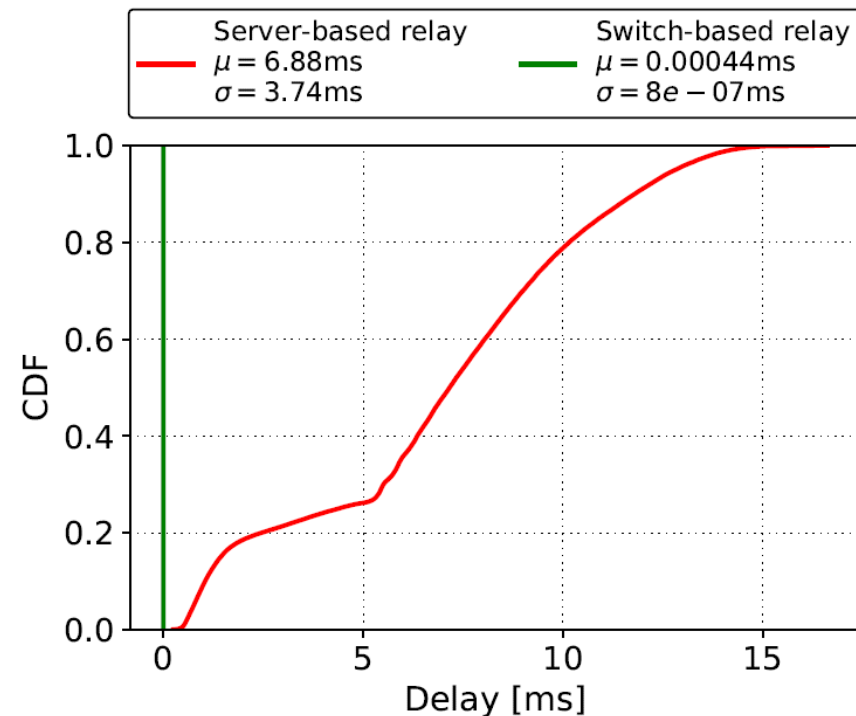
Implementation and Evaluation

- Two scenarios are considered:
 - “Server-based relay”: relay server is used to relay media between end devices, without the intervention of the switch
 - “Switch-based relay”: the switch is used to relay media
- UAC (SIPp) generates 900 media sessions, 30 per second
- The test lasts for 300 seconds
- G.711 media encoding codec (160 bytes every 20ms)



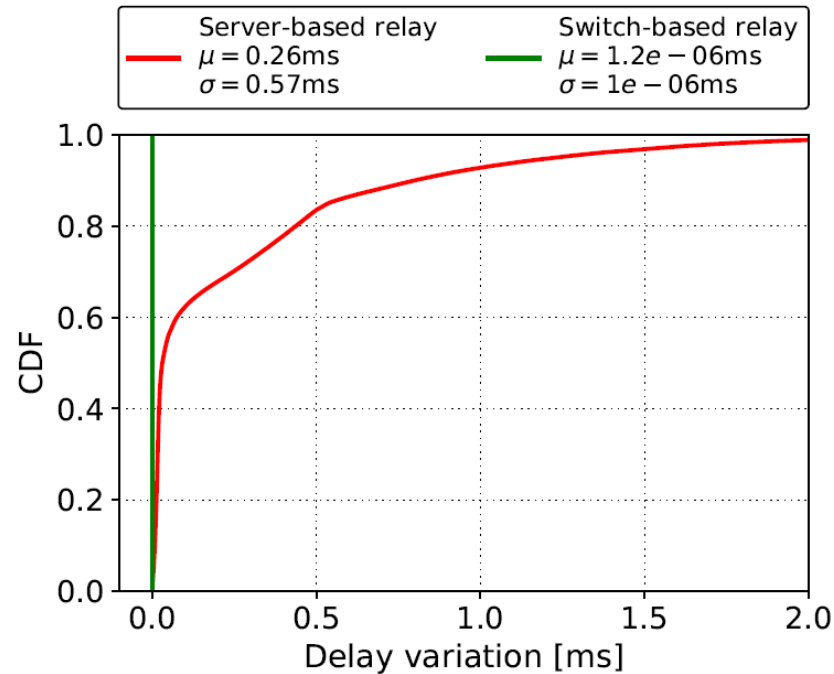
Results

- Delay: time interval starting when a packet is received from the UAC by the switch's ingress port and ending when the packet is forwarded by the switch's egress port to the UAS
 - Delay contributions of the switch and the relay server



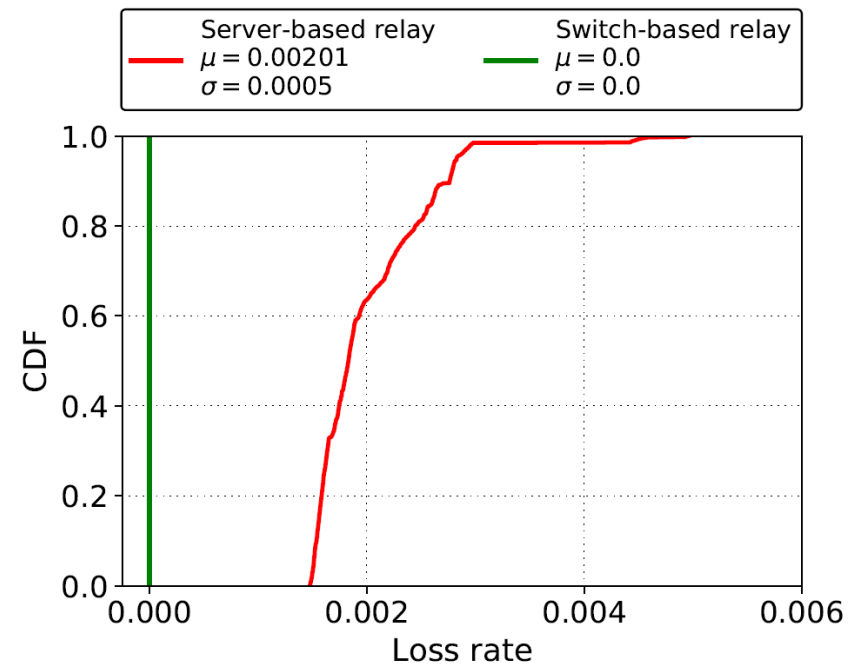
Results

- Delay variation: the absolute value of the difference between the delay of two consecutive packets
 - Analogous to jitter, as defined by RFC 4689



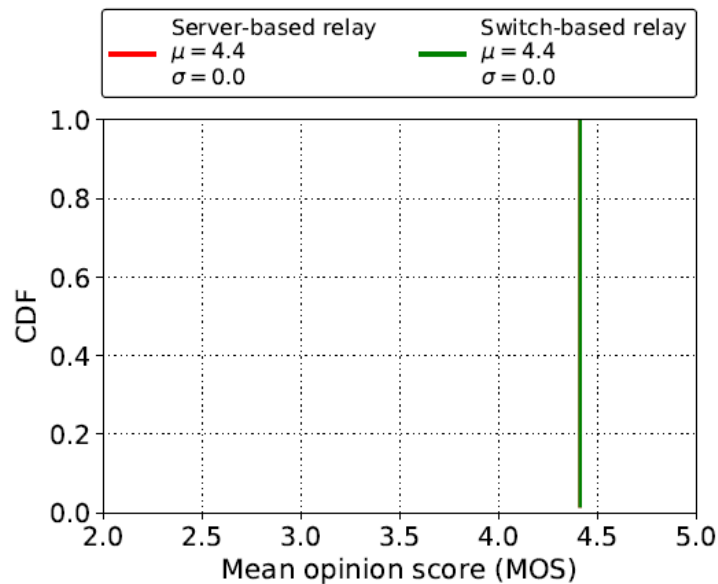
Results

- Loss rate: number of packets that fail to reach the destination
 - Calculation is based on the sequence number of the RTP header

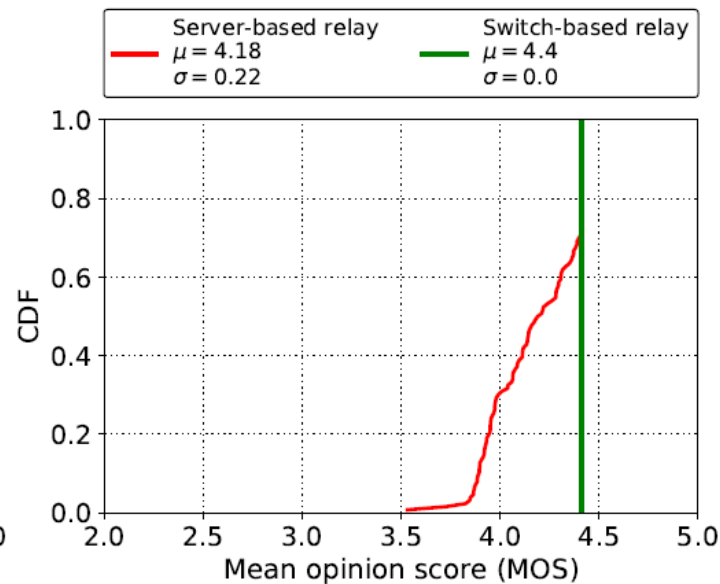


Results

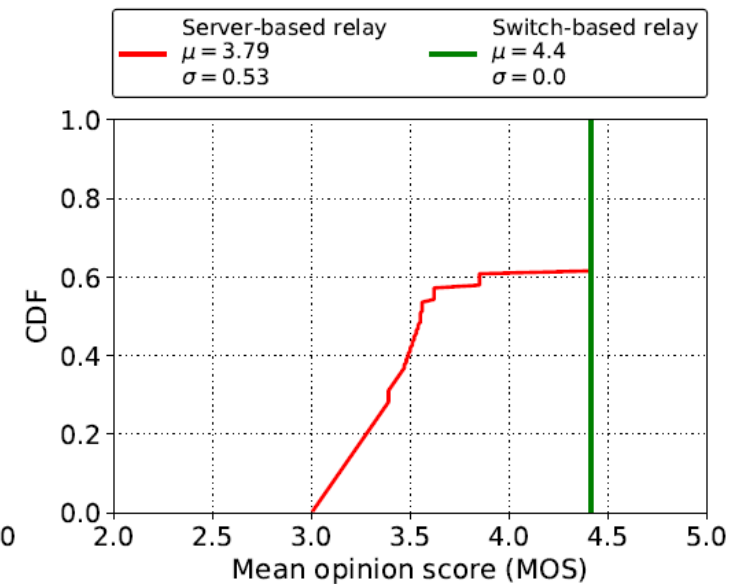
- Mean Opinion Score (MOS): estimation of the quality of the media session
 - A reference quality indicator standardized by ITU-T
 - Maximum for G.711 is ~4.4



(a) 750 simultaneous sessions.



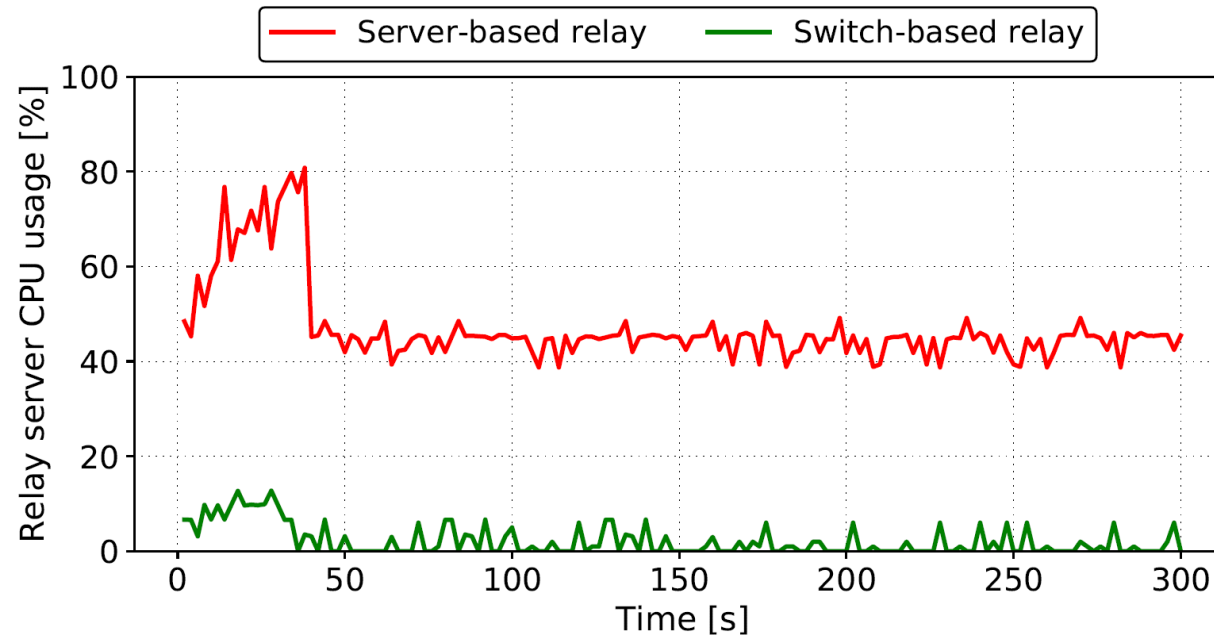
(b) 1500 simultaneous sessions.



(c) 1800 simultaneous sessions.

Results

- CPU usage: the percentage of the CPU's capacity used by the relay server



Resource Consumption

- The prototype is implemented in two different scenarios:
 - On top of the baseline switch program (switch.p4): implements various features including Layer 2/3 functionalities, ACL, QoS, etc.
 - Standalone implementation

On top of switch.p4			
Table size	SRAM	Hash Bits	TCAM
32,000	+8.45%	+2.7%	+0%
64,000	+16.2%	+4.6%	+0%

Standalone program			
Table size	SRAM	Hash Bits	TCAM
500,000	-----	-----	-----
1,000,000	+97.84%	+86.4%	+0%
1,050,000	+107.5%	+89.8%	+0%

Additional hardware resources used when the solution is deployed on top of switch.p4 and as a standalone program

Lessons Learned

- Advantages of using a switch-based relay:
 - Performance: ~1,000,000 sessions vs ~1,000 sessions per core
 - Optimal QoS parameters: delay, delay variation, packet loss rate
 - Flexibility: switch permits to modify / forward packets using non-standard fields
 - Precise timing information: measuring delay and its variation on the P4 switch results in precise high-resolution timing information
 - Programmer can free unused resources and customize program: accommodate additional sessions
- Limited resources
- Avoid complex application logic

Acknowledgement

- Thanks to the National Science Foundation (NSF)!
- Activities in the CI Lab at the UofSC are supported by NSF, Office of Advanced Cyberinfrastructure (OAC), awards 1925484 and 1829698





P4
Expert
Roundtable Series

April 28-29, 2020

Hosted by:



Thank You

Contact info for further questions

ekfoury@email.sc.edu

jcrichigno@cec.sc.edu

Full text

<https://tinyurl.com/wab7yej>

CI Lab website

<http://ce.sc.edu/cyberinfra/>