



# Building E2E Service Experience Assured Network (SEAN)

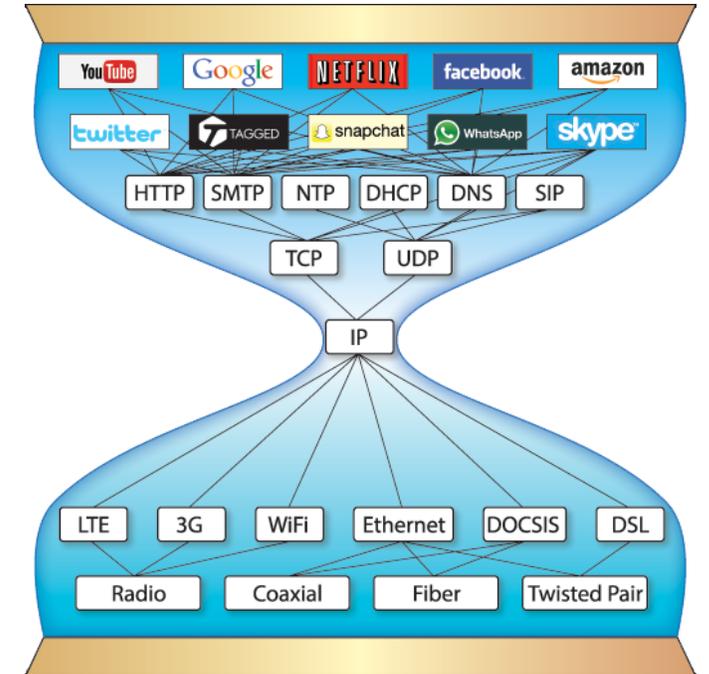
Andy Malis  
[andrew.malis@huawei.com](mailto:andrew.malis@huawei.com)

# Outline

- Overview of the problem space
- Requirements and challenges from new applications
- Current work to address these requirements (at IETF and elsewhere)
- What else may be needed, and proposals for the future
  - Spoiler alert – I'm not here to propose any particular solutions, but to provide food for thought and further discussion

# Challenges Facing the Internet

- IPv4/IPv6 are networking's “narrow waist” and support countless applications over the Internet
  - Some say the narrow waist is HTML, but the majority of applications don't run in a browser
- However, as we know, today's Internet has a number of deficiencies that still need to be addressed, including a lack of widely-available end-to-end service guarantees for latency, jitter, and packet loss
- Acts to limit the availability of some applications on wide-area
- Focus of this talk is on the wide-area multi-domain open Internet, not private/campus/data center networks or single-domain, well-managed, traffic-engineered enterprise services networks



# Some of the Architectural Issues Preventing E2E Internet Service Guarantees

- **Coordination across multiple layers**
  - Link Layer
  - Network layer
  - Upper layers
  - Control Plane
- **Inter-SDO Coordination and Cooperation**
  - Each is looking at their slice of the problem space, without taking a wider system view
  - Difficult to have E2E interoperable solutions
  - Liaisons don't always work effectively for such complicated multiple layer & domains solutions
- **Service Provider Coordination and Cooperation**
  - Inconsistent implementation of Diff-Serv in the public Internet
  - Inability to aggregate like flows that require similar treatment through the backbone
  - Even if these are implemented within a particular AS, lack of coordination at cross-AS boundaries

# Challenges to IP Networks from Wireless 5G

## Low latency is an important requirement

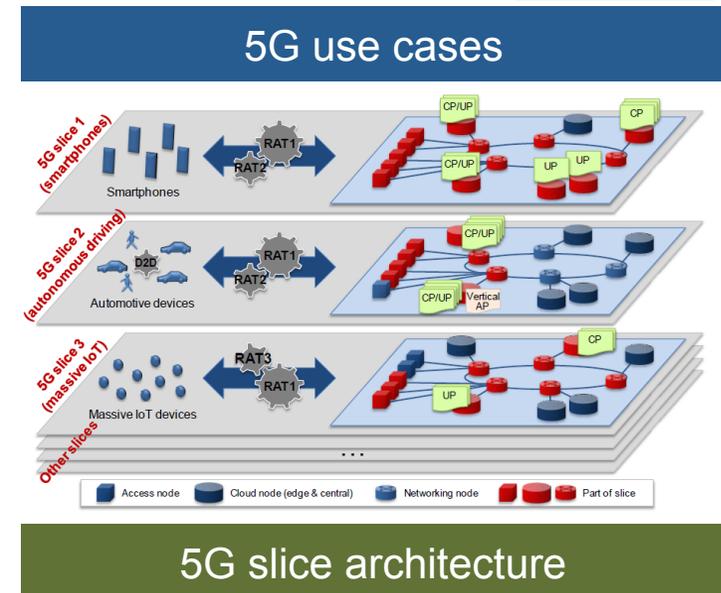
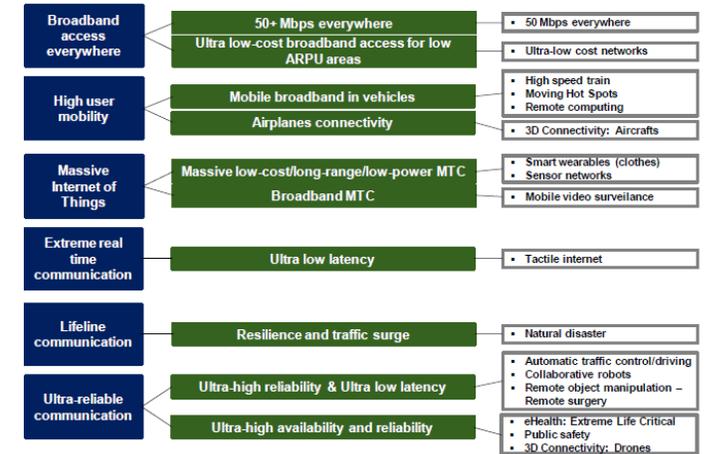
Use case category	User Experienced Data Rate	E2E Latency	Mobility
Broadband access in dense areas	DL: 300 Mbps UL: 50 Mbps	10 ms	On demand, 0-100 km/h
Indoor ultra-high broadband access	DL: 1 Gbps, UL: 500 Mbps	10 ms	Pedestrian
Broadband access in a crowd	DL: 25 Mbps UL: 50 Mbps	10 ms	Pedestrian
50+ Mbps everywhere	DL: 50 Mbps UL: 25 Mbps	10 ms	0-120 km/h
Ultra-low cost broadband access for low ARPU areas	DL: 10 Mbps UL: 10 Mbps	50 ms	on demand: 0-50 km/h
Mobile broadband in vehicles (cars, trains)	DL: 50 Mbps UL: 25 Mbps	10 ms	On demand, up to 500 km/h
Airplanes connectivity	DL: 15 Mbps per user UL: 7.5 Mbps per user	10 ms	Up to 1000 km/h
Massive low-cost/long-range/low-power MTC	Low (typically 1-100 kbps)	Seconds to hours	on demand: 0-500 km/h
Broadband MTC	See the requirements for the Broadband access in dense areas and 50+Mbps everywhere categories		
Ultra-low latency	DL: 50 Mbps UL: 25 Mbps	<1 ms	Pedestrian
Resilience and traffic surge	DL: 0.1-1 Mbps UL: 0.1-1 Mbps	Regular communication: not critical	0-120 km/h
Ultra-high reliability & Ultra-low latency	DL: From 50 kbps to 10 Mbps; UL: From a few bps to 10 Mbps	1 ms	on demand: 0-500 km/h
Ultra-high availability & reliability	DL: 10 Mbps UL: 10 Mbps	10 ms	On demand, 0-500 km/h
Broadcast like services	DL: Up to 200 Mbps UL: Modest (e.g. 500 kbps)	<100 ms	on demand: 0-500 km/h

From NGNM 5G white paper

- In 5G, a network slice supports the communication services for a particular connection type or application
- Each type has its specific properties and requirements
- E2E latency, in particular, becomes a strong requirement for some network slices (applications)

# 5G Requirements Require QoS Innovations

- 5G slice architecture is the main motivation
- Current QoS mechanisms cannot meet the new requirement for future 5G networks
- New opportunities and challenges
  - Isolation of resources
    - mix of exclusive and shared resources
  - Heterogeneous service requirements in the same forwarder
    - Each slice can be defined to support any combination of network requirement such as high availability, low latency, no packet drop, etc.



# Challenges from Emerging VR/AR



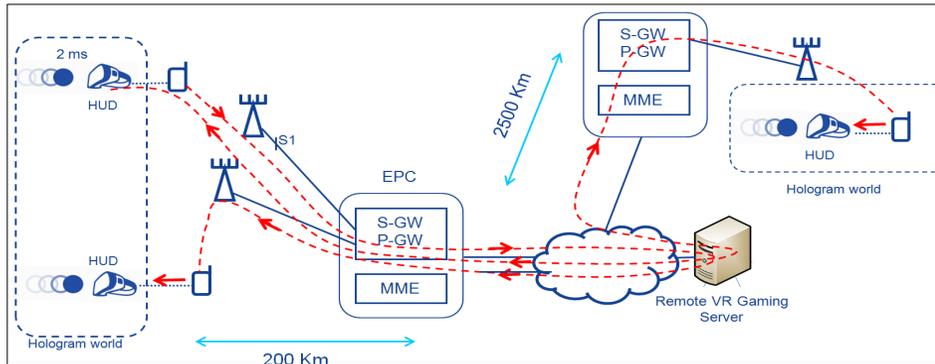
VR

Also known as immersive multimedia or computer-simulated reality, is a computer technology that replicates an environment, real or imagined, and simulates a user's physical presence and environment to allow for user interaction.



AR

AR is a live direct or indirect view of a physical, real-world environment whose elements are augmented (or supplemented) by computer-generated sensory input such as sound, video, graphics or GPS data.



## Implement Process:

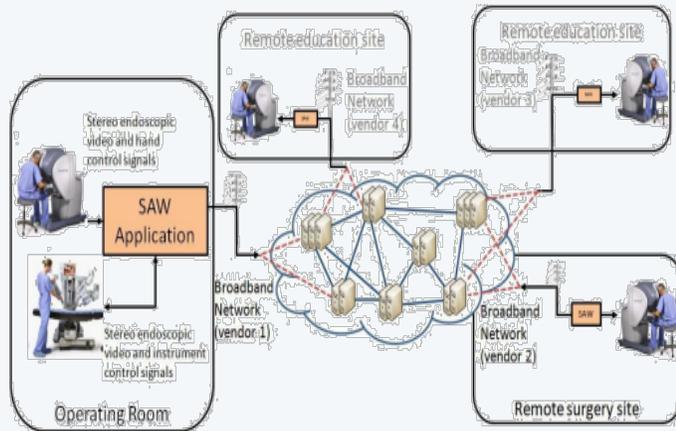
1. Collect input data from sensor
2. Transmit data from sensor to processing node in network
3. Process data and render view
4. Send rendered view to display devices
5. User experiences view from the screen

## Requirements for VR/AR:

- **Ultra-low latency**  
End-to-end latency **less than 20ms** including network transmission and processing time. (If latency exceeds 20ms, users will feel dizzy since users' motion perception doesn't match the images)
- **Bandwidth**  
A perfect experience requires bandwidth greater than **several Gb/s** depending on image resolution, frame rate, etc.
- **Screen resolution**  
The resolution should be larger than 4K to prevent screen door effect or latticing
- **Screen refresh rate**  
The refresh rate should be at least 60Hz

# Requirements From Other Emerging Applications

## Remote Healthcare



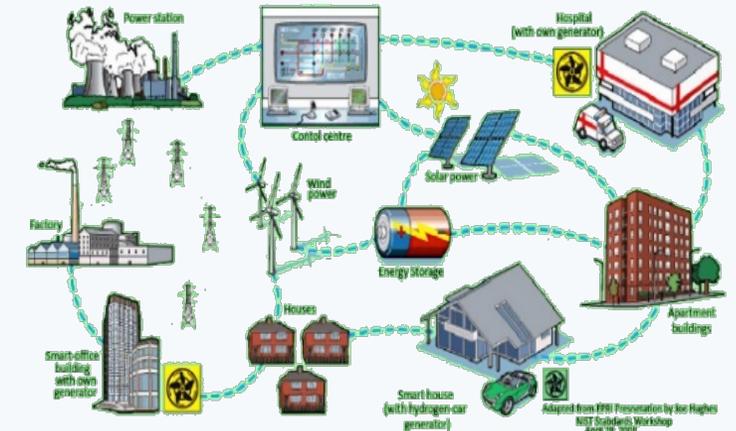
- High-fidelity interaction is fundamental to the safe deployment of tele-medical technologies
- To achieve the fidelity required for tele-medical applications, it is necessary to achieve end-to-end latencies of 1-10 ms and highly reliable data transmission.

## Factory Automation



- Factory automation can enable many new possibilities for discrete manufacturing and help producers achieve more efficient production.
- The sensitivity of control circuits when controlling devices moving rapidly (such as industrial robots) requires an end-to-end latency significantly below 1 ms per sensor.

## Smart Grid



- Optimize consumers' power supply and reduce associated costs.
- A synchronous co-phasing of power suppliers requires an end-to-end latency in the order of 1 ms. This 1 ms latency results in a phase shift of  $18^\circ$  (50 Hertz AC network) or  $21.6^\circ$  (60 Hertz AC network).

# 2013 ISOC Workshop on Latency

- **In 2013, Internet Society hosted a Reducing Internet Latency workshop**
- **Major sources of latency:**
  - Processing: Computational translation, forwarding, encap/decap, NAT, encrypt, auth., compression, error coding, signal translation
  - Multiplexing: Delays needed to support sharing, shared channel acquisition, output queuing, connection establishment
  - Grouping: Reduced frequency of control information and processing, packetization, message aggregation
- **The workshop's report (<http://www.bobbriscoe.net/projects/latency/latws-ccr.pdf>) concluded:**
  - There are fundamental limits to the extent to which latency can be reduced, but there is considerable capacity for improvement throughout the system, making Internet latency a multifaceted challenge.
  - “How to standardize a definition of access network latency such that latency could be used (like bandwidth) as a unit of commerce. It turns out that's a hard problem.”

# What's Changed since 2013?

- **Content is more distributed**
  - Low latency is so crucial to many applications that many of those apps' management systems need visibility to network latency characteristics to intelligently distribute content in places with the minimal latency.
  - Easier with stored content, much harder with dynamically rendered content (AR/VR)
- **Advancements and innovations have been made at various layer in exploring technologies to meet low latency requirements**
  - Upper Layer: QUIC, L4S
  - Network layer: IP/MPLS Hardened Pipe (RFC 7625), latency-optimized router design, and BBF's Broadband Assured Services (BAS).
  - Link layer: IETF DETNET, IEEE 802.1 TSN (Time Sensitive Networking), Flex Ethernet (OIF).
  - 3GPP has started multiple projects related to reducing latency in RAN, core, and in backhaul.
- **With the latest technologies advancement in packet networks, it is becoming feasible to partition access network resources to dedicate resources to flows that need deterministic latency**
- **It is no longer impossible to have solutions to the “hard problem of the access network”**

# BBF BAS addresses performance-assured IP services

The Broadband Forum has started work on a BAS (Broadband Assured IP Services) project in its Innovation and Architecture work areas.



New requirements for for dynamic, high speed on-demand services such as performance-assured interactive videoconferencing, interactive gaming on new platforms, high quality 4K/8K content delivery, 5G mobile, secure vertical market applications (finance, healthcare, government).



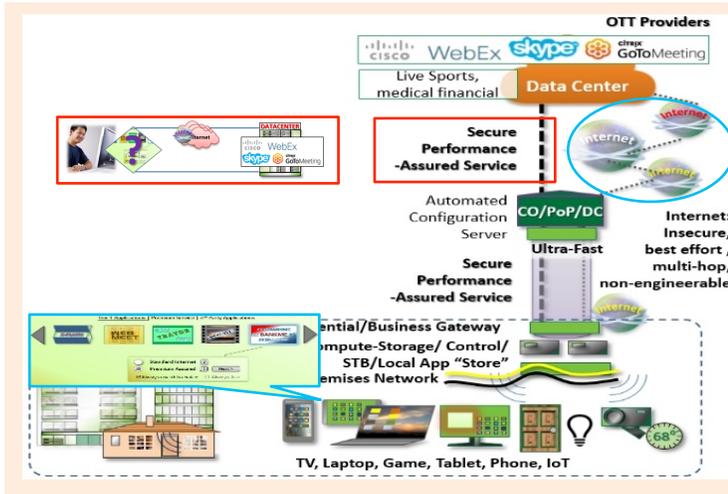
Emergence of virtualization, SDN, G.fast, NGPON, mobile 4G/5G, and other new access technologies that bring down infrastructure costs and enable new service offerings.



Strong support of access network service providers and their customers for supporting new services and applications over and above current IP service offerings.

BBF exploring on-demand performance-assured IP services, with cloud services, data center interconnect, and fixed/mobile convergence the initial focus areas for end-customer applications.

# BAS Covers Scenarios, Terminology , Requirements, and Architecture



## Typical Use Case Scenarios:

- End-user Assured Cloud Service delivery
- Data Center Interconnect
- Converged, hybrid wireless/wireline
- IoT secure access
- Teleconferencing
- Assured mobile services at Wi-Fi hotspots
- .....

BAS will provide a new E2E performance assured network service:

- Bandwidth
- Latency
- Packet loss
- Security
- .....

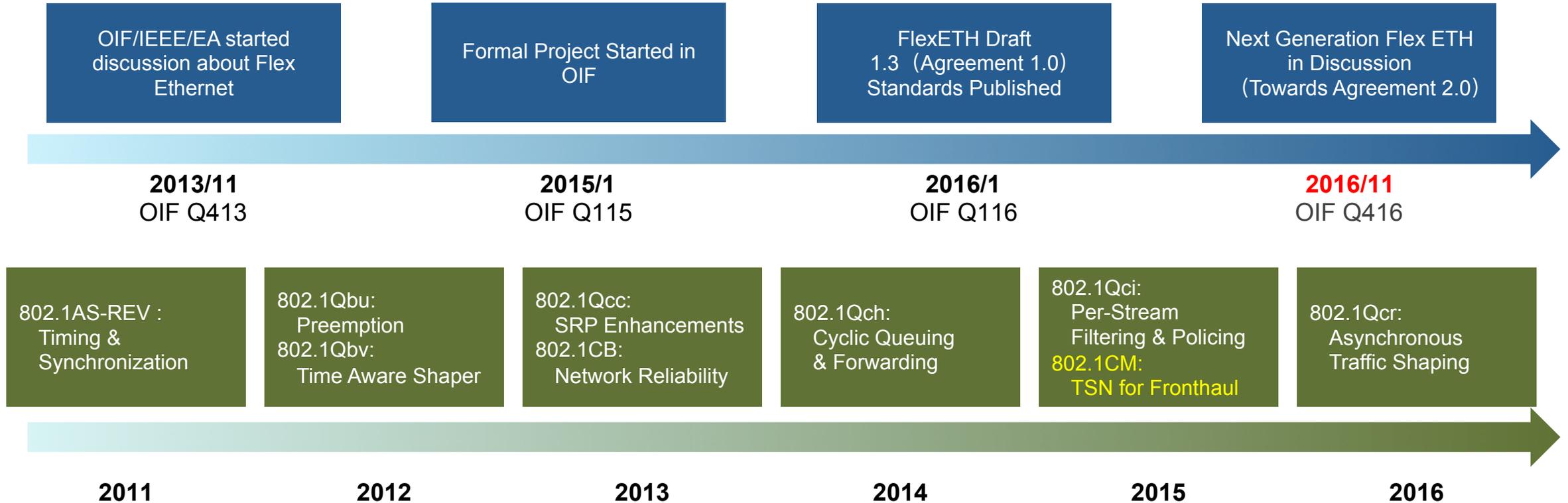
## Project Scope:

- Terminology
- BAS Architecture
- BAS Use Cases and Requirements
- BAS Information, Data Models
- Management Model
- Service Definitions and Performance Objectives
- Performance Monitoring, Assurance
- Business and Marketing Documents

## Out of Scope:

- Changes to underlying technology owned by another SDO
- If changes necessary, will partner with other SDOs (primarily expected to be IETF and IEEE) as required to jointly improve the technology

# Flex Ethernet: Related Standards Work in OIF and IEEE



- Flex Ethernet provides strong assurance for low latency and deterministic jitter at link level
- IEEE 802.1 TSN targets low latency, low jitter Ethernet bridges, to enable time sensitive applications like 5G fronthaul and in-vehicle networking, etc.

# IETF Exploring Mechanisms to Reduce Network Latency

- L4S BOF: Low Latency Low Loss Scalable throughput

- Latency (queuing delay) is the factor limiting application performance
- L4S will work on the fine saw-teeth congestion control to enable the low latency and low loss TCP

- QUIC: Quick UDP Internet Connection

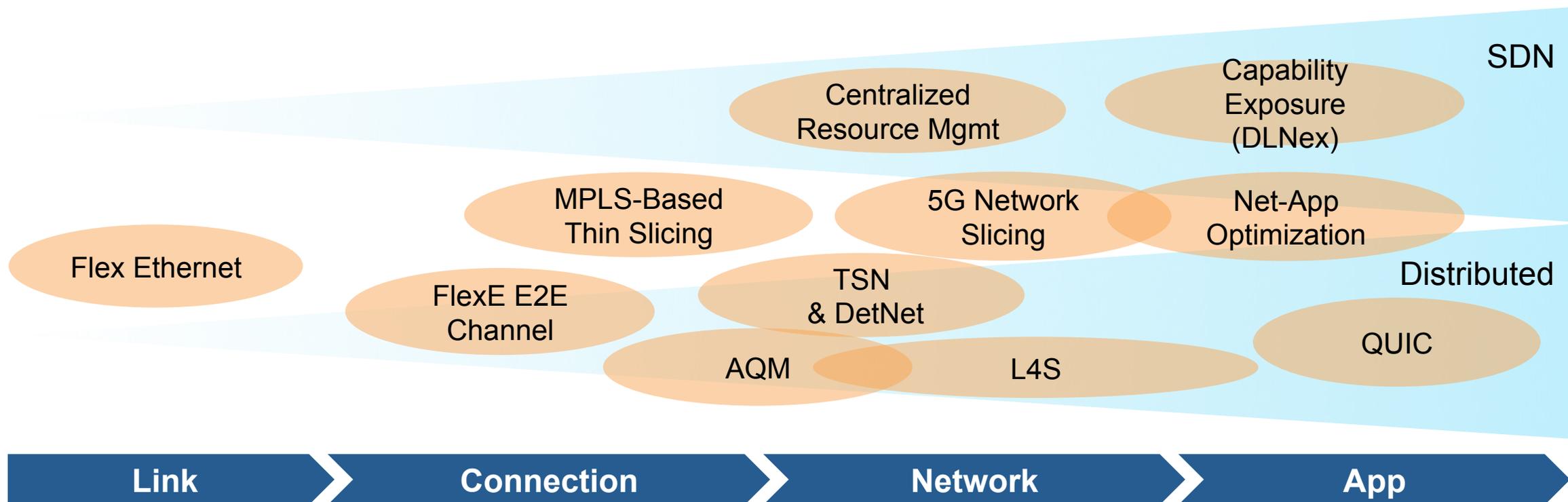
- Define a new standards track IETF transport protocol on top of UDP based on deployment experience from Google, etc.
- Reduces latency for HTTP when compared to TCP

- DetNet: Deterministic Networking

- Focuses on deterministic data paths that operate over Layer 2 bridged and Layer 3 routed segments
- Provides bounds on latency, loss, and packet delay variation (jitter), and high reliability.
- Will not spend energy on solutions for large groups of domains such as the Internet.



# Point Solutions in Multiple SDOs



Key to successful standardization:

1. Innovating new mechanisms to support new capability
2. Reusing protocols to minimize standards & products development effort
3. Collaboration: SDOs, vendors, operators, academic research, etc.

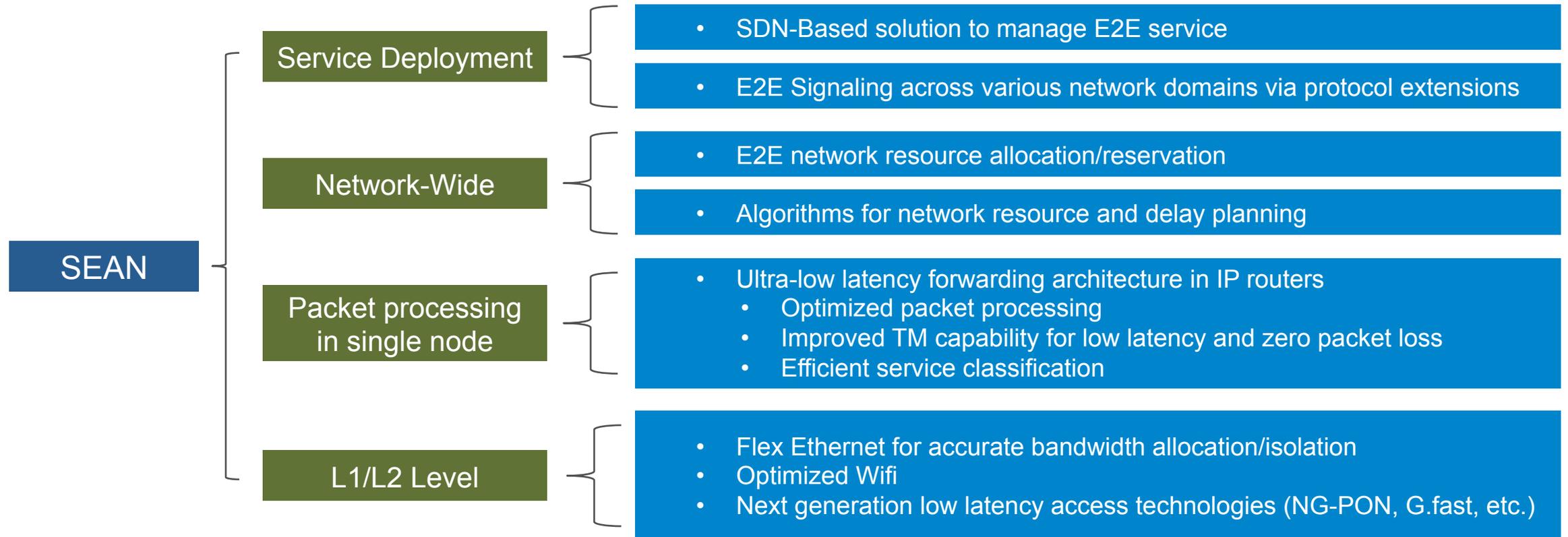
# SEAN: Service Experience Assured Network

## Key Characteristics about SEAN:

- Ultra low latency and jitter
- Zero packet loss
- Constant bandwidth assurance
- On demand “zero wait” deployment

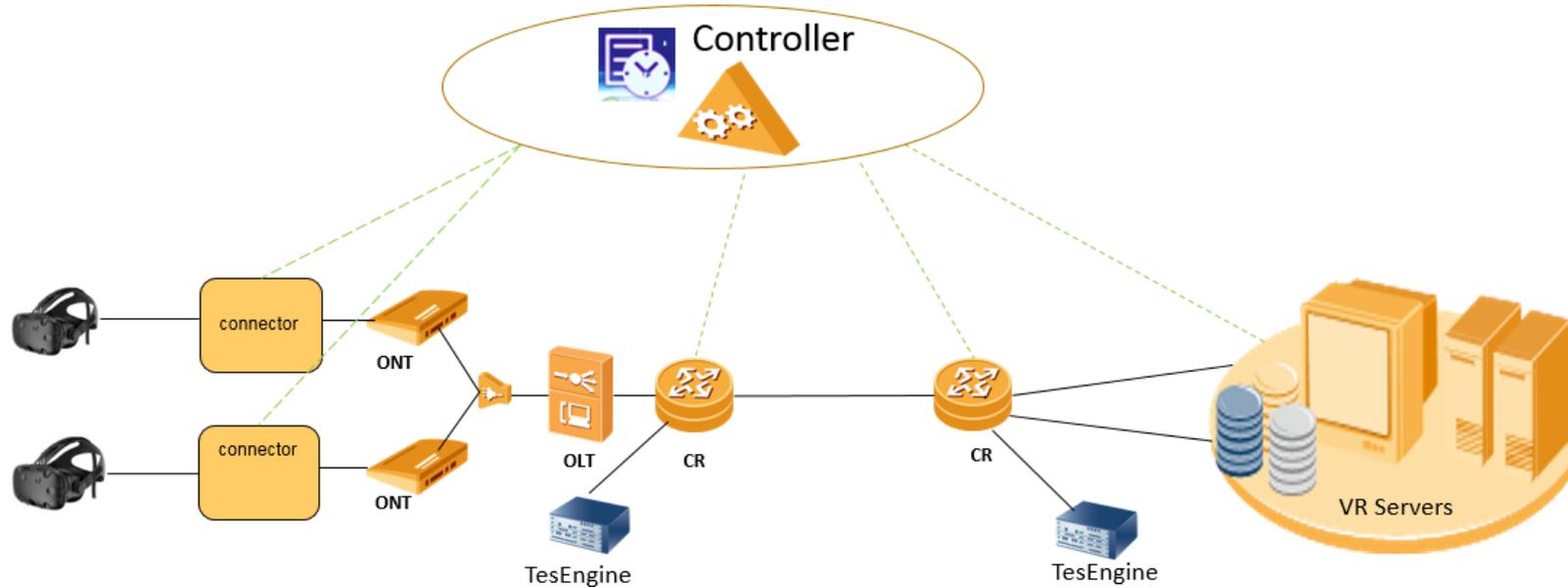


# Key Technologies to make SEAN Possible



- To enable SEAN we need a comprehensive approach encompassing multiple technologies
- Standards will be an important part to make SEAN possible

# SEAN Experimentation in Development



Cloud real-time VR service are assured with SEAN even in a congested network

Key technologies adopted

1. SDN based solution with E2E network resource reservation
2. Flex Ethernet with strict bandwidth guarantee
3. Innovative ultra-low latency forwarding architecture
4. Optimized access technologies in OTN and OLT

VR requires latency less than 20 ms, we expect network total RTT less than 5 ms (Metro)

# Where Do We Go From Here?

- **To further the development of available and new technologies to achieve low end-to-end latency and loss over wide-area packet networks**
  - Enable more latency sensitive applications (even the ones we haven't imaged yet) to go through the Internet
- **To better utilize the technologies on low latency initiatives developed by other SDOs**
  - A possible workshop to spur new developments
  - Focus on cross-layer interaction issues and system-level issues
  - Perhaps leading to new standardization based on workshop results
- **To answer questions like:**
  - Are there advantages to be gained from taking a broader view of these different technologies at different layers and how they might interact?
  - What are the effective Interaction/coordination between upper and lower layers so that efficient optimization can be achieved for latency- and/or loss-sensitive services? For example, some applications would prefer the network to drop their packets instead of transient routers' large buffers causing jitter & latency?
  - Given the larger universe of issues that prevent E2E service guarantees in the Internet, where can we most intelligently attack the problem to get the "biggest bang for the buck"?

# Questions?