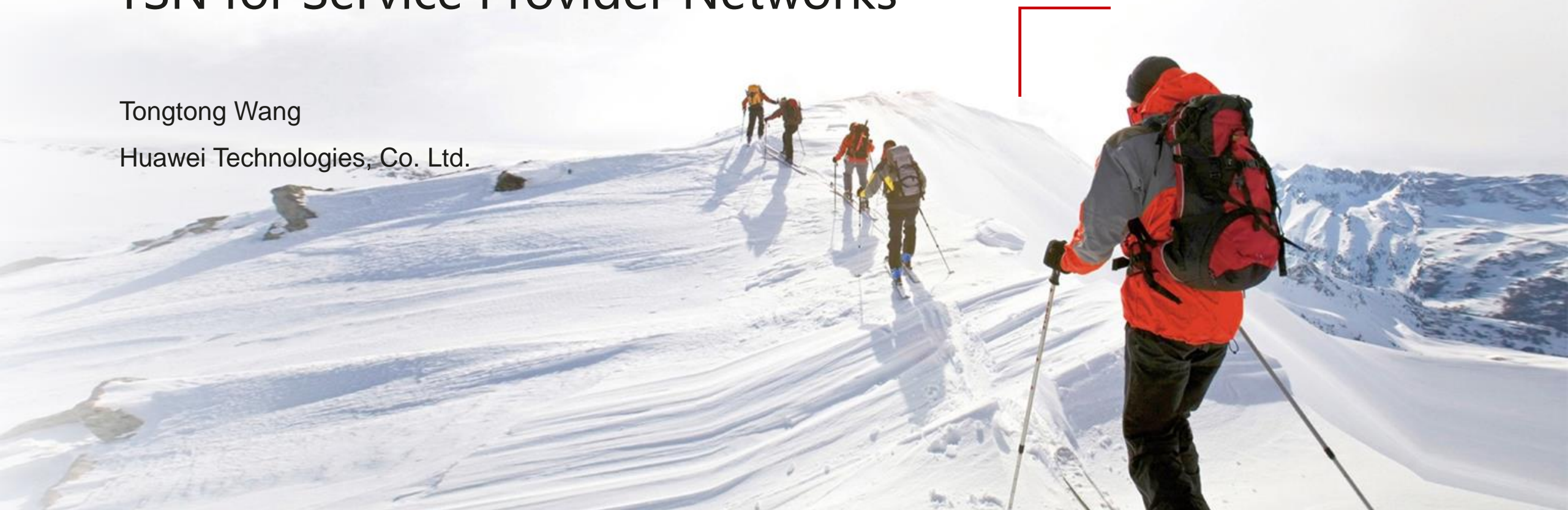


Latency Analysis in TSN for Service Provider Networks

Tongtong Wang

Huawei Technologies, Co. Ltd.



TSN for Service Provider Networks

- Use cases / Requirement

5QI Value	Resource Type	Default Priority Level	Packet Delay Budget	Packet Error Rate	Default Maximum Data Burst Volume(NOTE 2)	Default Averaging Window	Example Services
1	GBR NOTE 1	20	100 ms	10^{-2}	N/A	2000 ms	Conversational Voice
2		40	150 ms	10^{-3}	N/A	2000 ms	Conversational Video (Live Streaming)
3		30	50 ms	10^{-3}	N/A	2000 ms	Real Time Gaming, V2X messages Electricity distribution - medium voltage, Process automation - monitoring
4		50	300 ms	10^{-6}	N/A	2000 ms	Non-Conversational Video (Buffered Streaming)
65		Bandwidth Sensitive Services				2000 ms	Mission Critical user plane Push To Talk voice (e.g., MCPTT)
66						2000 ms	Non-Mission-Critical user plane Push To Talk voice
67		15	100 ms	10^{-2}	N/A	2000 ms	Mission Critical Video user plane
75		25	50 ms	10^{-2}	N/A	2000 ms	V2X messages
5		10	100 ms	10^{-6}	N/A	N/A	IMS Signalling
6	Non-GBR NOTE 1	60	300 ms	10^{-6}	N/A	N/A	Video (Buffered Streaming) TCP-based (e.g., www, e-mail, chat, ftp, p2p file sharing, progressive video, etc.)
7		70	100 ms	10^{-3}	N/A	N/A	Voice, Video (Live Streaming) Interactive Gaming
8		80	---	---	N/A	N/A	Video (Buffered Streaming) TCP-based (e.g., www, e-mail, chat, ftp, p2p file sharing, progressive video, etc.)
9		Connection Services				N/A	Mission Critical delay sensitive signalling (e.g., MC-PTT signalling)
69		55	200 ms	10^{-6}	N/A	N/A	Mission Critical Data (e.g. example services are the same as QCI 6/8/9)
70		65	50 ms	10^{-2}	N/A	N/A	V2X messages
79		68	10 ms	10^{-6}	N/A	N/A	Low Latency eMBB applications Augmented Reality
80		11	5 ms	10^{-5}	160 B	2000 ms	Remote control (see TS 22 261 [2])
82	Delay Critical GBR	12	10 ms NOTE 1	10^{-5}	320 B	2000 ms	Intelligent transport systems
83		URLLC Latency Sensitive Services				10 ms	Intelligent Transport Systems
84		17	10 ms	10^{-5}	2352 B	2000 ms	Discrete Automation
85		22	10 ms	10^{-4}	1358 B NOTE 3	2000 ms	Discrete Automation

3GPP SA2 Table 5.7.4-1: Standardized 5QI to QoS characteristics mapping

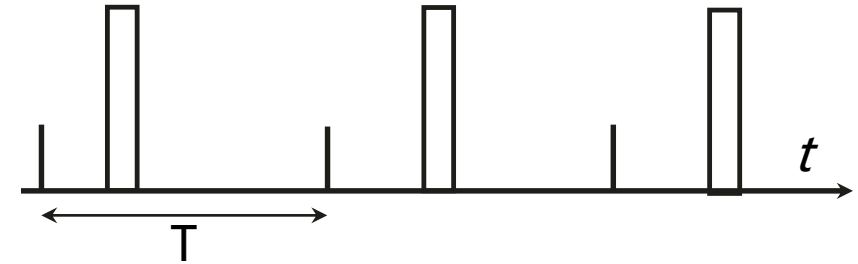
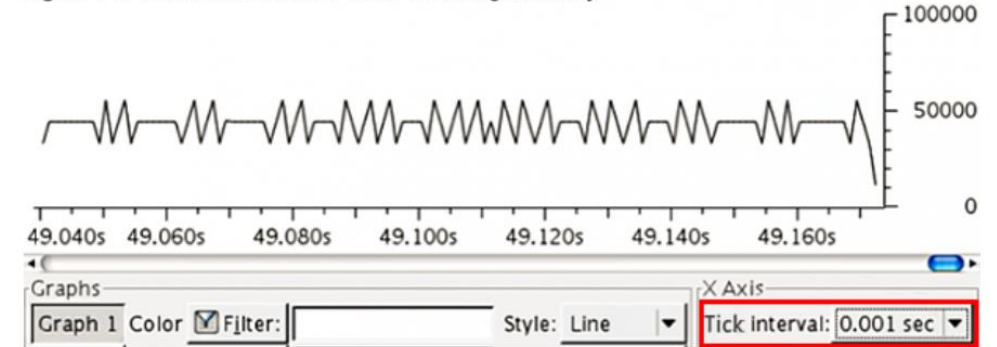
Three main types for services co-exist in 5G carrier networks, especially in backhaul networks or metro networks, according to 3GPP requirement document.

- *Bandwidth sensitive services* care more about throughput as long as its short term bursting can be buffered and transit later.
- *Latency sensitive services* are new applications appear with new era carrier networks. TSN techniques are most useful on this aspect.
- *Connection services* are the legacy services on 802 or IP networks.

Differentiated SLA in Service Provider Networks

- Type 1 Bandwidth sensitive application
 - Care more about average rate and peak rate, usually have CIR/PIR parameters to specify user requirements. Not so much on latency, have tolerance on microburst and congestion, use buffering to solve bursting and gaps in data stream.
 - Use cases : TV/Sport Videos, Surveillance Video, etc.
- Type 2 Latency sensitive application
 - Care more about bounded latency, have T-SPEC parameter to specify traffic model (TrafficInterval, BurstSize), and latency bound requirement.
 - Use cases : Smart Grid Teleprotection, Cloud VR.
- Type 3 Connection services
 - Messages

Figure 1-2 Traffic statistics at a lower level of granularity



TSN Toolbox on Forwarding Plane

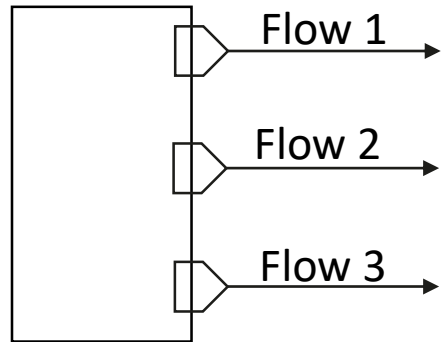
Both scheduling and shaping functions affect latency bound. For some specific TSN techniques, they have shaping and scheduling function combined together.

- TAS/Async TAS; (Refer to IEEE Std 802.1Qbv)
- Generic CQF; (Refer to IEEE Std 802.1Qch)
- Strict Priority/WRR; (Refer to IEEE Std 802.1Q-2018)
- CBS/BLS; (Refer to IEEE Std 802.1Qav)
- Dedicated physical lines

Some techniques are more for bounded latency, some others are more about to get low jitter; while Strict Priority/WRR algorithms generally have lower average delay and may benefit on bandwidth utilizations.

Within Smart Grid scenarios, we compared legacy QoS algorithms(Strict Priority/Deficit Round Robin/etc.) and dedicated physical lines to ensure that with fundamental 802.1 TSN techniques , bounded latency services can be achieved in service provider networks; nevertheless, performance on latency and jitter varies.

TSN Toolbox on Forwarding Schedulers and Shapers



Dedicate Ports/Links

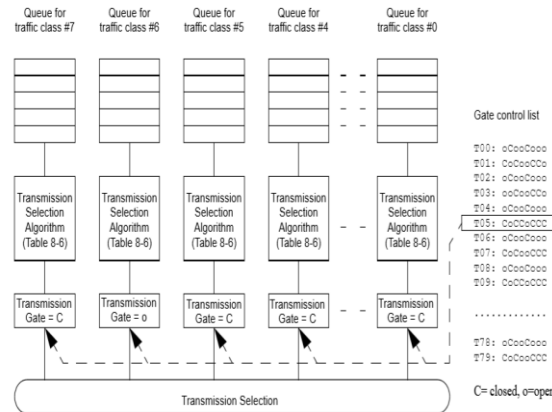
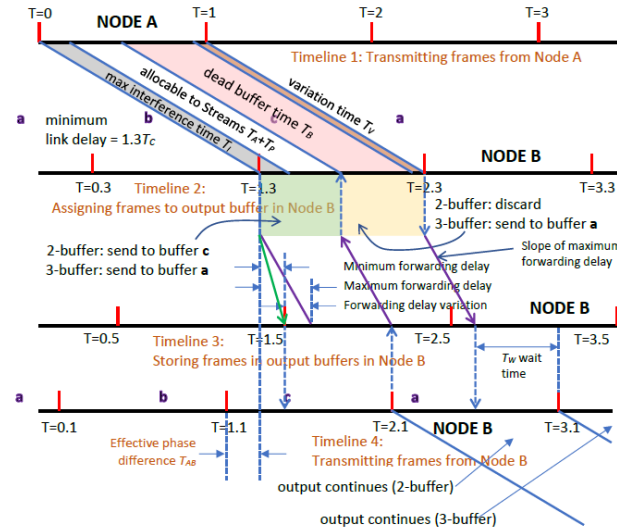
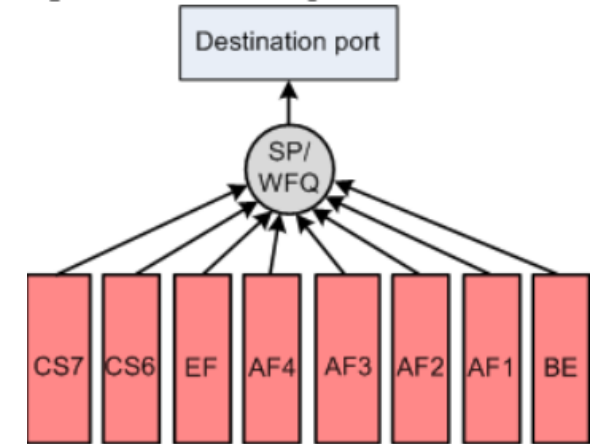


Figure 8-14—Transmission selection with gates

Time Aware Shaper



*Generic CQF
(2 buffer/3buffer/etc.)*



Strict Priority/Weight Fair Queueing

By using different schedulers and shapers, multiple types of traffic are transmitted with differentiated service levels (SLA) on shared network resources;

This is similar to network slicing concept, to divide network up and share among users/applications;

TSN techniques are capable to support network slicing with multiple levels of service guarantee;

Delay Decomposition and Comparison

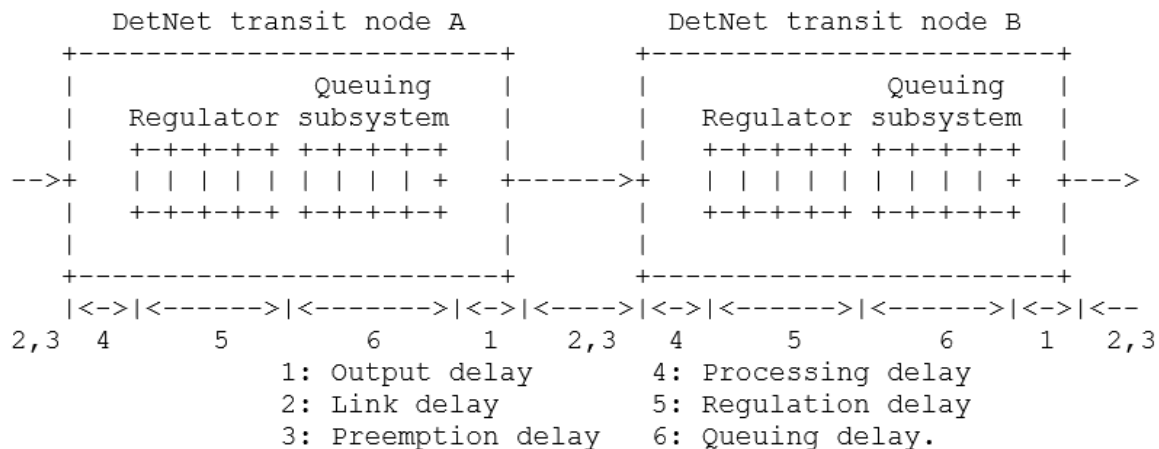
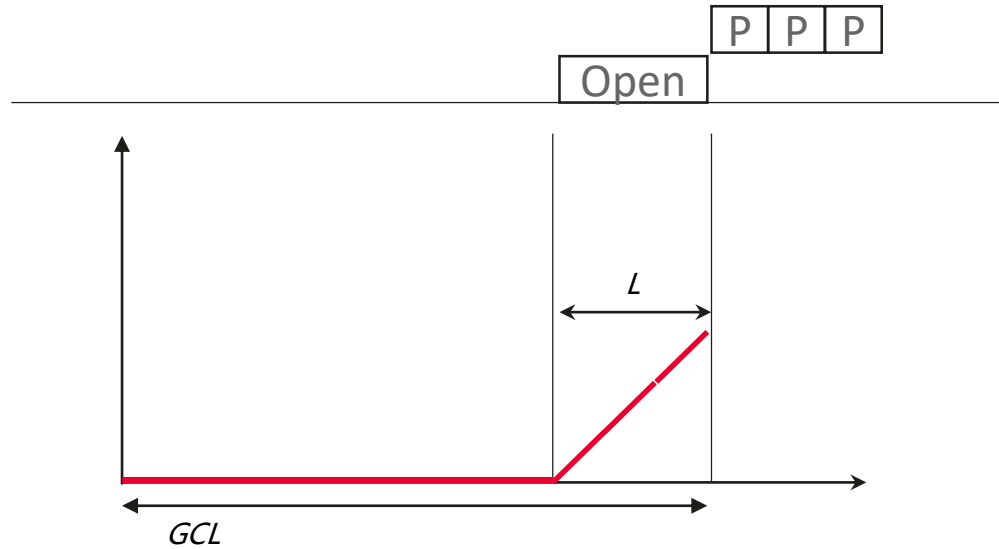


Figure 1: Timing model for DetNet or TSN

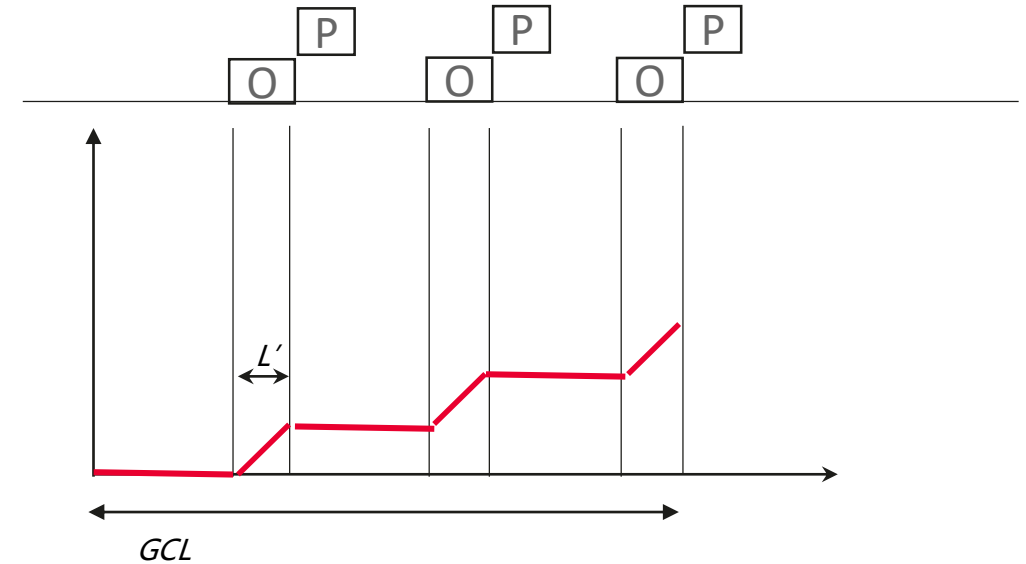
<https://tools.ietf.org/html/draft-finn-detnet-bounded-latency-04>

	Queueing Delay
Dedicate Port	$D_{Port} = \frac{b}{R}$; //R is output port data rate
TAS	Depends on GCL configuration and traffic timing; //need more assumption to compare
CQF	$2 * T_c$; //more analysis from Norman's whitepaper
SP/WFQ	$\frac{b_i}{R \frac{Q_i}{\sum Q}} + \frac{b}{R}$ // Q_i is the share of the bandwidth for flow i
CBS	$T^A = \frac{1}{R} \left(L^A + \frac{\bar{L}}{R} \right)$ //for Class A traffic, assuming no CDT

TAS Latency Analysis



$$\text{Max Delay} = \text{GCL} - L + \frac{L}{R}$$



$$\text{Max Delay} = \text{GCL}/3 - L/3 + \frac{L}{3R}$$

- Max Delay on TAS scheme depends on whether gate control configurations suits with traffic bursts;
- Complex TAS configuration schemes may cause more complicated latency formula;

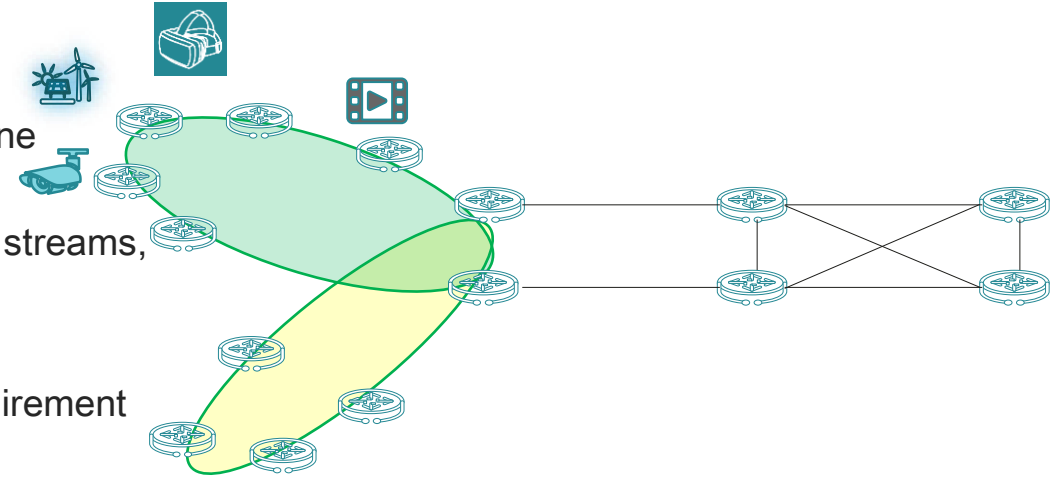
L. Zhao, P. Pop and S. S. Craciunas, "Worst-Case Latency Analysis for IEEE 802.1Qbv Time Sensitive Networks Using Network Calculus," in IEEE Access, vol. 6, pp. 41803-41815, 2018.

IEEE 802.1 TSN Geneva, Jan 2020

Delay Analysis in TSN for Service Provider Networks

- Comparison Between Schedulers

- Typical topology on backhaul network, with access ring, aggregate layer and backbone layer.
- Assume 8 DTU(digital transmitting unit) on access ring、16 VR stream and 40 video streams, with periodic burst traffic model.
 - Smart grid : in every 5ms , 14*380Byte. //Tight latency requirement ~2ms.
 - VR Game : in every 15ms , 60*1522Byte (~50Mbps) //medium latency requirement ~10ms
 - Video : in every 70ms , 168*1522Byte (~30Mbps) //loose latency requirement ~10ms
- Different ways of bandwidth sharing(scheduling and shaping) cause performance varying.



Worst Case Latency Calculation (ms)	Dedicate Link	FIFO @10G	Strict Priority @10G	Strict Priority @10G + Edge Shaping②	CQF @10G③
Smart Grid	0.267 @2G	12.883	0.077 (H) ①	1.081 (H)	TBA
Teleprotection					
VR Game	3.796 @6G	13.347	2.565 (M)	8.138 (M)	
Video	73.438 @2G	14.778	23.019(L)	19.121(L)	
Worst case Latency in Simulation (ms)	Dedicate Link	FIFO @10G	Strict Priority @10G	Strict Priority @10G + Edge Shaping	CQF @10G
Smart Grid	0.244 @2G	7.991	0.050	0.728	TBA
Teleprotection					
VR Game	3.219 @6G	8.797	1.896	7.845	
Video	50.257 @2G	10.378	11.456	11.432	

TSN for Service Provider Networks Annex :

- Network Calculus

- Briefly introduce Network calculus as necessary to explain delay analysis for different schedulers and shapers. //This work may update CBS latency analysis and bring up amendment request.

- Arrival Curves :

- Current research on Network calculus use token bucket attributes on traffic models to setup arrive curve $\alpha(t)$, such as

$$\alpha_{\sigma,\rho}(t) = \begin{cases} \rho t + \sigma, & t \geq 0 \\ 0, & t < 0, \end{cases}$$

while each flow has a burstsize σ , and data rate ρ ;

//Note: TSN standard usually use T-SPEC to describe input traffic model, need conversion here.

//Aggregating behavior also changes traffic model

- Service Curves:

- The service offer by the scheduler on an outgoing port can be characterized by a minimum service curve, denoted by $\beta(t)$.
 - Service Curves for different schedulers(Strict Priority/WRR/CBS/TAS) will be sorted out and then evaluate the queueing latency

Thank you.

把数字世界带入每个人、每个家庭、
每个组织，构建万物互联的智能世界。

Bring digital to every person, home and
organization for a fully connected,
intelligent world.

