

# Tunnel congestion Feedback

(draft-wei-tunnel-congestion-feedback-01)

Xinpeng Wei    Lei Zhu    **Lingli Deng**

Huawei            Huawei            China Mobile

IETF 89 London, UK

# Problem Statement

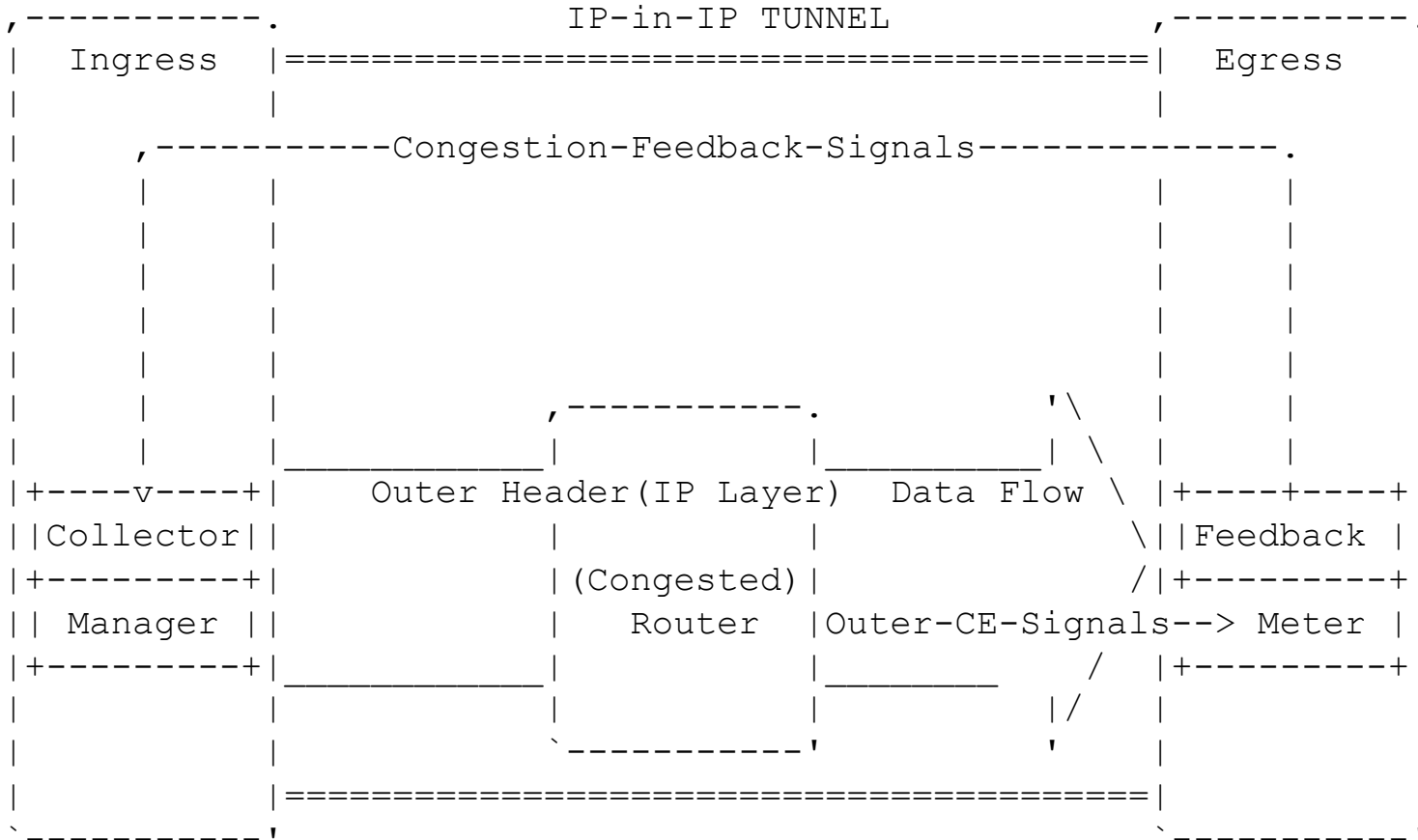
## ❖ Motivation

- ❖ There is significant variability in aggregated bandwidth demands and latency in managed networks such as mobile backhaul and DC network.
- ❖ Tunnels are widely deployed to carry end user flows in network (both backhaul network and DC network).
- ❖ Resource provision based on congestion status is expected to be helpful in promoting the overall resource utility of network resource.

## ❖ Problem

- ❖ While the congestion experienced in tunnels can be calculated according to RFC 6040, Appendix C.
- ❖ However, there is no standard feedback mechanism for congestion information from Egress to Ingress router of the tunnel.

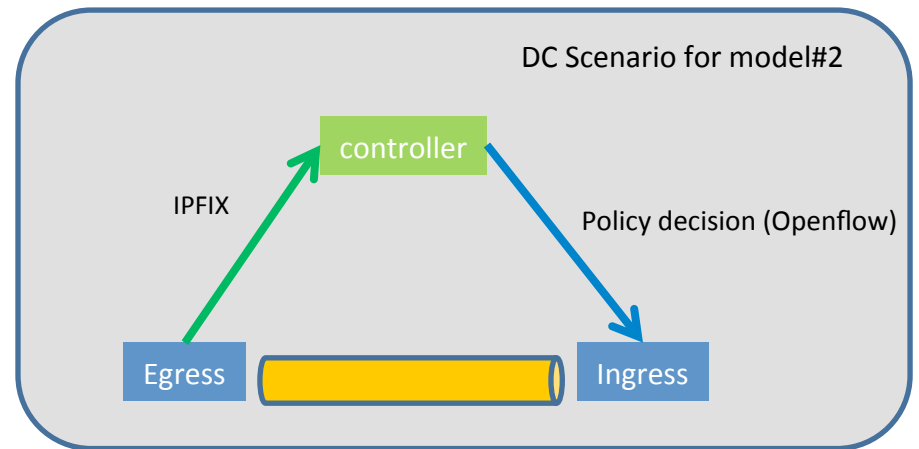
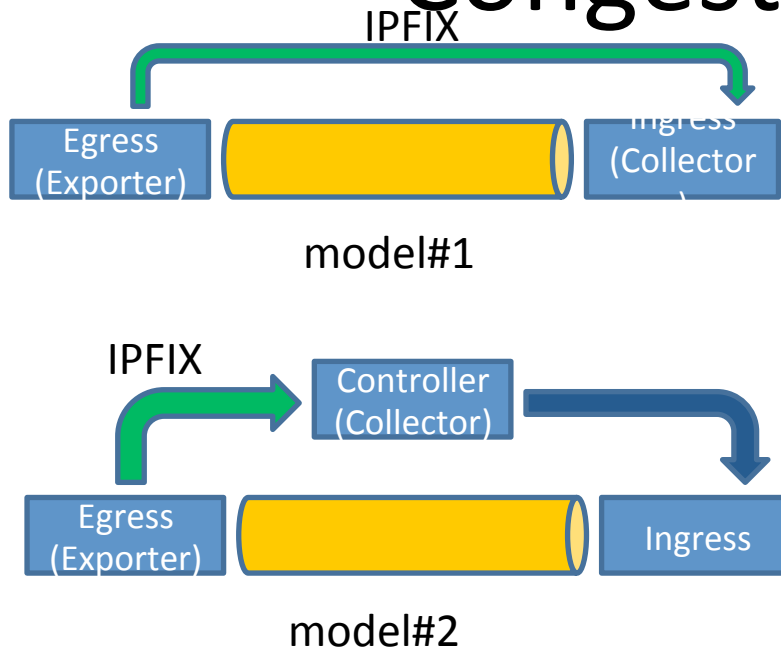
# Congestion Feedback Basic Model



# How to calculate congestion in tunnel

- Uses RFC6040 Appendix C.
- Egress calculates tunnel congestion in a statistical way.
  - the proportion of packets with CE marks within the tunnel
  - packets not marked in the inner header but have a CE marking in the outer header
  - based on moving average (MA) algorithm.

# IPFIX mapping for Tunnel Congestion Feedback



IPFIX is primarily used to convey Information about IP flows passing through a network element (extend to convey tunnel congestion):

- TCP (or SCTP) based (congestion friendly)
- Different choice of feedback triggering: on demand request, periodically.
- Extensible / flexible information export model.
- **Many existing information elements could be used to describe flow congestion information.**

# Relationships between tunnel congestion feedback and e2e CC

- Local feedback v.s. end2end feedback
  - Tunnel congestion control is a kind of local congestion control/ feedback within a specific administration domain on the path of the correspondent e2e flow.
  - It only responds to the congestion happened in the tunnel.
  - The tunnel congestion control is complementary with e2e ECN control.
- Network Management oriented v.s. end2end CC
  - The tunnel congestion feedback provides network administrator with network congestion level information that can be used as an input for resource provision management, not necessarily flow-based CC.
  - Example: differentiated traffic gating, if the tunnel is congested it will be a waste of resource to allow new low-priority traffic enter as they arrive spontaneously, as they may eventually get dropped in the tunnel.

# Next steps

- Is the capability negotiation needed between ingress and egress?
- What other congestion related information would be conveyed? And in what way?
  - e.g. the major contributor to the congestion.
  - Is the information aggregation among different flows happens at the meter or collector?

# Backup Slides



# How to calculate congestion in tunnel

- The algorithm to calculate congestion.

The basic idea of calculating congestion statistically in tunnel is :

Calculating the congestion level of a subset of traffic flows in the tunnel, and take the result as congestion level of the whole tunnel.

Rationale: all the traffics are treated equally by router according to RED, and when congestion occurs in the router, the router randomly selects the packets to mark.

Here we take the traffic that is ECN capable and not congestion-marked before tunnel to calculate congestion.

When ingress is conformant to RFC6040, the packets collected by egress can be divided into to 4 categories, shown as the figure.

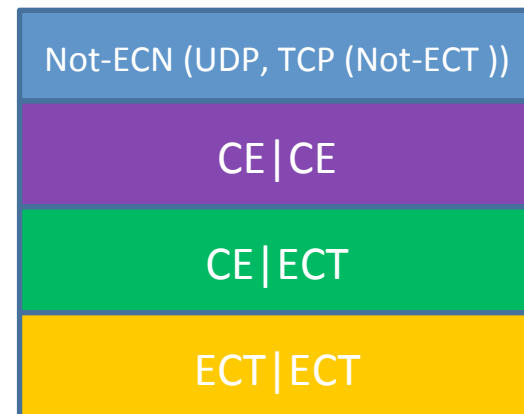
The tag before “|” stands for ECN field in outer header; and the tag after “|” stands for ECN field in inner header.

“Not-ECN” means the traffic that don’t support ECN, such as UDP and Not-ECT marked TCP;

“CE|CE” means the ECN capable packets that have CE-marked before entering tunnel;

“CE|ECT” means ECN capable packets that CE-marked in tunnel;

“ECT|ECT” means ECN capable packets that have not congested in tunnel.

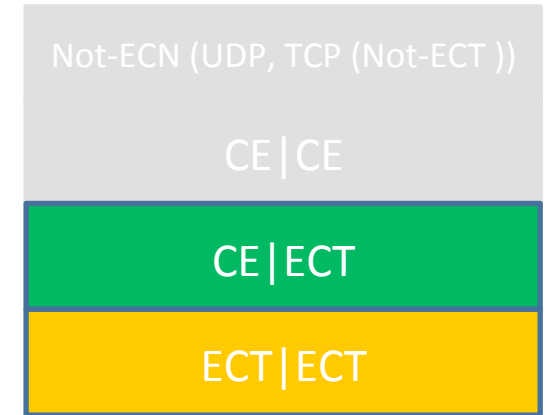


# How to calculate congestion in tunnel

Assuming the quantity of CE|ECT packets is A, the quantity of ECT|ECT packets is B, then the congestion level (C) can be calculate as following:

$$C=A/(A+B)$$

Here as an example, we take 100 packets to calculate the moving average. As analyzed above, we just need to take CE|ECT and ECT|ECT packets into consideration. Every time we calculate the congestion, we use the current packet (CE|ECT or ECT|ECT) and the last 99 packets (CE|ECT or ECT|ECT) to get the moving average result which stands for current congestion.



## NOTES:

- (1) There only shows a simple method of moving average, some other method, such as weight moving average may also be used .
- (2) The calculation is based on the assumption that all the packets are treated equally by routers, but the existence of DSCP may has some impacts for this.
- (3) According to the calculation process the UDP traffic may don't have too much impact.
- (4) In 3GPP scenario, the congestion may be calculated by ECN field, e.g. when the congestion occurs in RAN.

# The basic procedure of congestion feedback

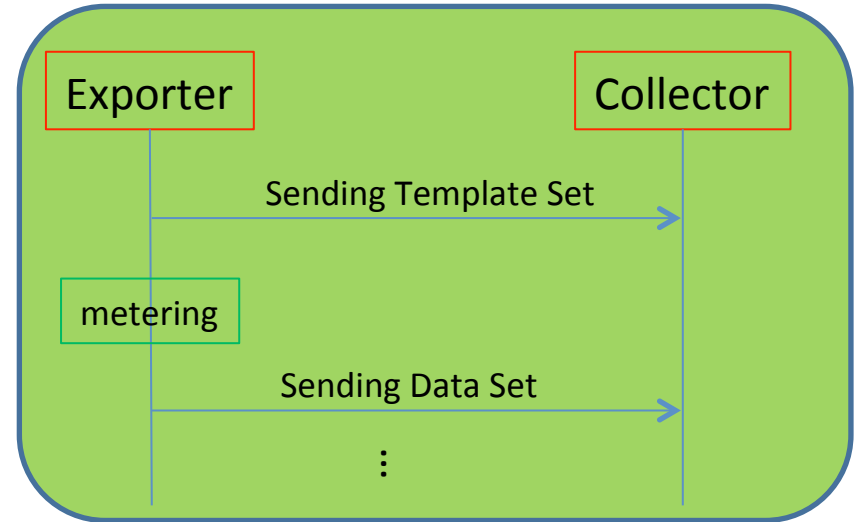
Here an example is shown to illustrate how IPFIX can be used for congestion feedback, the information conveyed here may be incomplete. **The exact information to be conveyed from exporter to collector needs further discussion.**

## (1) Sending Template Set

The exporter use Template Set to inform the collector how to interpret the IEs in the following Data Set.

Set ID=2	Length = octets
Template ID= 257	Field Count =
exporterIPv4Address = 130	Field Length = 4
collectorIPv4Address = 211	Field Length = 4
Congestion Level = TBD1	Field Length = 2
Enterprise Number = TBD2	

Exporter sends Data Set periodically or by trigger.



## (2) Sending Data Set

The exporter meters the traffic and sends the congestion information to collector by Data Set.

Set ID = 257	Length = octets
	192.0.2.12
	192.0.2.34
	15

# The basic procedure of congestion feedback

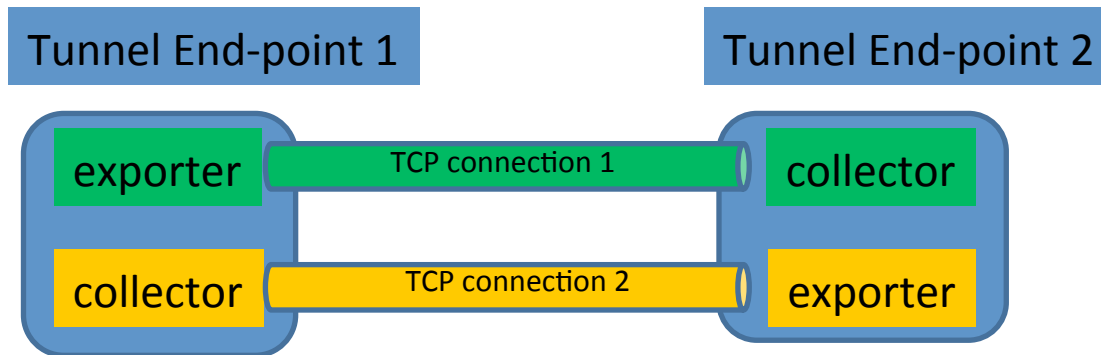
Congestion information to be reported:

Congestion volume	mandatory
Egress IP address	mandatory
Ingress IP address	mandatory
.....	.....

# Evaluation of IPFIX

The message flow of IPFIX is **unidirectional**, which means the information is only from exporter to collector. But in the tunnel scenario, the ingress/egress can host both exporting process and collecting process, so for a pair of ingress and egress there should be two TCP connections to convey IPFIX message bidirectionally.

*[NOTES: The need of two TCP connections between ingress and egress may be a shortcoming here. But I am wondering if information in two direction can be transported through one TCP connection.]*



The scenario that two TCP connections are established between two tunnel end-point, each connection for one direction. There will be an exporter and a collector process on each tunnel end-point.