

Contract Number: IST-1999-20841
Project Title: SEQUIN



Deliverable D2.1 - Addendum 1

Implementation architecture specification for the Premium IP service

Contractual Date:
Actual Date:
Work Package: WP2
Work Item: WI8.2
Nature of Deliverable: R - Report
Dissemination Level: PU - Public

Authors: Mauro Campanella INFN/GARR Mauro.Campanella@garr.it

ABSTRACT

This addendum specifies the implementation architecture for the Premium IP service described in Deliverable D 9.1, which aims at offering the equivalent of an end to end virtual leased line service at the IP layer across multiple domains. The architecture is targeted at the GÉANT network and is applicable to each connected NRENs and local Diffserv domains.

The architecture leverages the scalability features of Differentiated Services and takes a pragmatic approach to balance configuration complexity and benefits.

1	EXECUTIVE SUMMARY	3
2	ARCHITECTURE SUMMARY	4
	BASIC PRINCIPLES.....	5
3	DETAILED ARCHITECTURE SPECIFICATION.....	7
3.1	INTRODUCTION.....	7
3.2	SERVICE COMPONENTS SPECIFICATION.....	7
3.2.1	<i>Shaping</i>	7
3.2.2	<i>Fair sharing between Elastic and Anelastic flows</i>	8
3.2.3	<i>Policing</i>	8
3.2.4	<i>Choice of token bucket depth and MTU size</i>	8
3.2.5	<i>Admission control and Classification</i>	9
3.2.6	<i>Marking</i>	9
3.2.7	<i>Scheduling</i>	10
3.2.8	<i>Premium IP link capacity</i>	10
3.2.9	<i>Monitoring and accounting</i>	10
3.3	SPECIFICATION OF FUNCTION PER NODE	10
3.3.1	<i>Source node</i>	11
3.3.2	<i>Domain L1</i>	11
3.3.3	<i>Domain N1</i>	11
3.3.4	<i>Domain CORE</i>	12
3.3.5	<i>Domain N2</i>	12
3.3.6	<i>Domain L2</i>	13
4	PRACTICAL CONSIDERATIONS.....	14
5	RISK ANALYSYS AND LIMITS.....	14
6	SECURITY CONSIDERATIONS	15
7	ACKNOWLEDGEMENTS.....	15
8	REFERENCES	16
9	ACRONYMS.....	17

1 EXECUTIVE SUMMARY

This addendum specifies the implementation architecture for the Premium IP service defined in Deliverable D 2.1 [D2.1]. The Premium IP service aims at offering the equivalent of an end to end virtual leased line service at the IP layer across multiple domains. Although the architecture is targeted at the GÉANT [GÉANT] network, it is general enough to be applicable to any similar topology of communicating Diffserv domains like each connected NRENs and their local user Diffserv domains.

According to Deliverable D2.1 the architecture is based on Differentiated Services and the Expedited Forwarding (EF) Per Hop Behaviour. The network is decomposed in Diffserv domains and rules for interconnection, which require the interconnection to behave as an EF hop.

The architecture detailed here minimises the number of actions to be performed on every packet at each node and builds an initial configuration, which does not use a signaling protocol. In particular shaping is a requirement for the incoming user flow and will not be performed by the network. Policing is required only for ingress traffic and strictly only at the ingress to the first domain.

The aim is at delivering a QoS service, based on current knowledge and availability of QoS technologies, which could be implemented in a short time scale, compatible with the start of the GÉANT network.

According to the experience gained through the use of the service and new set of experimental evidence and theoretical developments, the structure of the architecture is such that it can be easily improved and adapted to new requirements and scaling needs.

2 ARCHITECTURE SUMMARY

The GÉANT [GÉANT] network will be based on very high speed links, equal or greater than 2.5 Gigabit per second, with a limited amount of meshing. GÉANT will connect a large number of National Research and Educational Networks (NREN), more than 25, which have or will have a core backbone engineered along the same guidelines.

Figure 1 depicts a simplified version of the overall network structure which will be used to illustrate and analyse the Premium IP implementation architecture.

The sample network is decomposed into multiple, communicating Diffserv domains, named L1, L2, N1, N2 and CORE, and no particular hypotheses are required about internal topology, physical structure or transmission technology of each of them.

The architecture will detail the behaviour of a Diffserv domain and the rules for their interconnection. If the model is applied to the GÉANT network, additional simplifications can be assumed, like very high speed in the core.

According to the EF PHB specifications [EFPHB] the rate at which EF traffic is served at a given output interface should be at least the configured rate R , over a suitably defined interval, independent of the offered load of non-EF traffic to that interface.

Flow shaping can be the foundation of a correct behaviour of the whole service. It ensures that the flow does not incur in packet losses due to policing and minimises creation of burstiness due to aggregation between different flows. Last but not least, shaping each flow ensures a fair sharing of the services between elastic and anelastic transport protocols like TCP and UDP.

The architecture requires shaping at the source, or a second choice, by the network as close as possible to the source. The network will then not apply any additional shaping before the delivery to final destination.

The flows must also be strictly policed as near to the source as possible and packets violating the contract are discarded. Policing is performed according to at least three mandatory parameters: IP source, IP destination prefixes and agreed sending rate.

At the initial policing point, packets successfully admitted to the service, are marked with an appropriate DSCP or IP Precedence value and queued in the highest priority queue for delivery.

In addition, it is suggested that at the border between different domains, for example at the core accesses points, an additional policing action can be performed based on aggregate bandwidth specifications for each (ingress, egress) pair of the domain.

The capacity, here intended as aggregate sending or receiving rate, limit for border policing is suggested to be configured at a value larger than sum of the agreed sending rates for the Premium IP service flows crossing the border. This policing action is performed as a safety measure, to limit service degradation to only a part of the network in case of incorrect configuration or denial of service attacks. Policing is never performed at egress interfaces and it is not mandatory when exiting a core, or trusted, core domain toward a user site.

High priority queuing, according to the Premium IP tag only, is enabled at all Premium IP participating nodes of each domain. The Premium IP enabled node set can be a sub-set of all the nodes of the domain.

The decision of not enabling additional shaping and policing is a consequence of the highest priority scheduling provided to Premium IP packets, of the very high speed and over-provisioned characteristics of the GÉANT and NRENs core backbone. The experimental validation of the assumption for domains, which deploy much lower link speeds, is in progress.

Monitoring of the service performance is enabled from the start of the service and it is performed both reading SNMP counters and by in-band active measurements of the performance of basic QoS parameters (delay, its variation, capacity and packet loss).

According to the experience gained through the use of the service and new set of experimental evidence and theoretical developments, the architecture structure is such that can be easily improved and adapted to new requirements. In particular the specification for shaping and policing location and the techniques for policing can be fine-tuned to improve the performance of the service and its scaling capabilities.

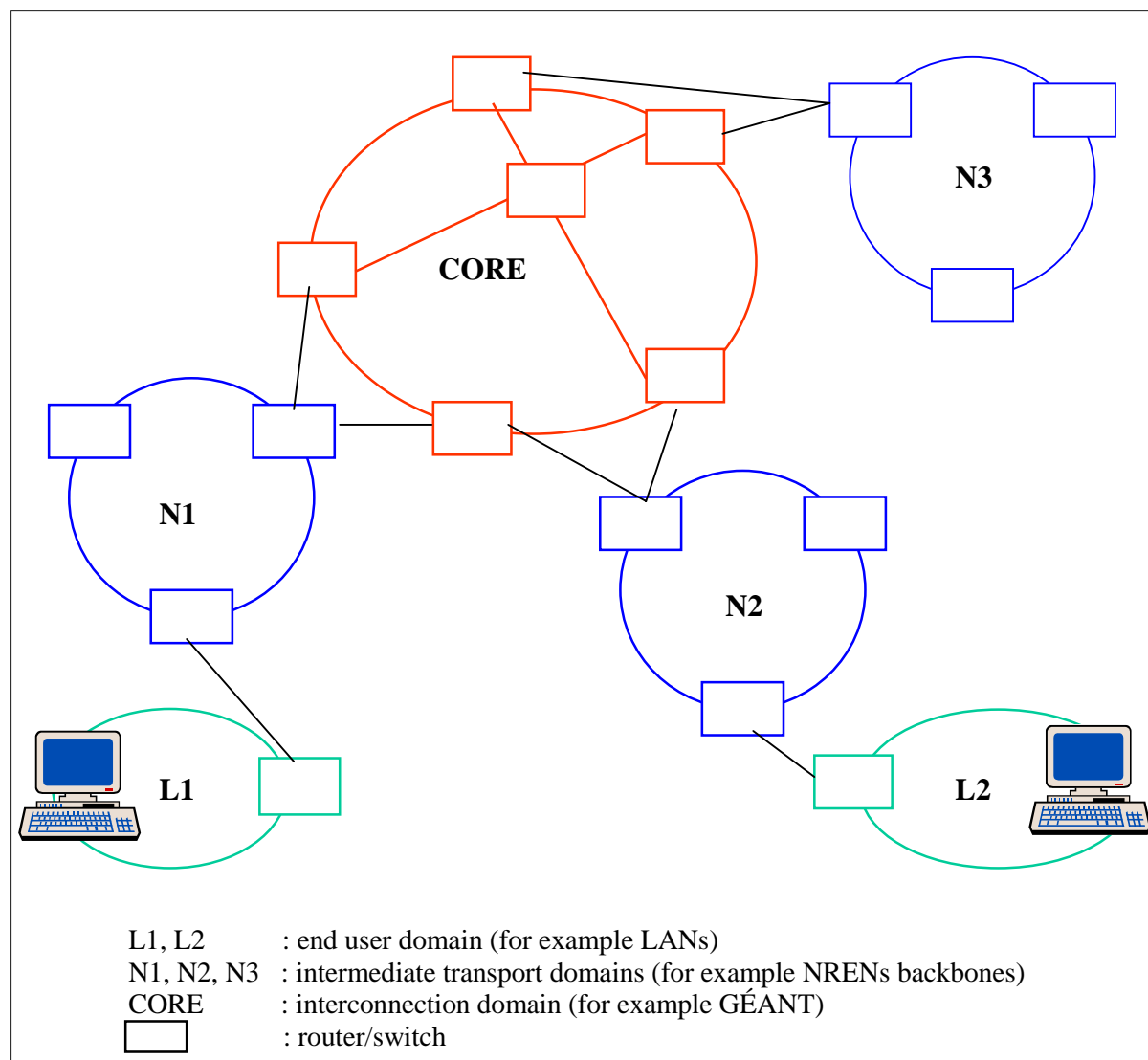


Figure 1: simplified multidomain network

2.1 BASIC PRINCIPLES

Enabling Quality of Service mandates a set of actions to be performed on every packet, even if not each action will be performed at each network node. The architecture detailed here minimises the number of actions at each node and builds an initial configuration, which does not use a signaling protocol.

The aim is at enabling a scalable and flexible QoS Premium IP service in the same timeframe of the start of the GÉANT network, based on current knowledge and availability of QoS technologies.

A set of assumptions and decisions is taken:

- The service will not police or shape per flow, being based on an aggregation model. Nonetheless it may police aggregates according to destination domains as a safety measure. An initial policing stage to ensure compliance to agreed sending rate is mandatory.
- The foreseen speed of the core links and the highest priority scheduling for Premium IP packets make delay variation small even at aggregation points. At 2.5 Gb/s the transmission time of a 1500 bytes packet is about 5 microseconds. The consideration suggests starting the service without enabling shaping in the core and it hints to have shaping optional also at the border, provided the sources produce well shaped flows.
- The sending host or source is required to shape flows it sends according to its allowed sending rate, to avoid initial packet loss due to policing when entering the Diffserv domain, and to ensure a fair share of the aggregate Premium IP capacity amongst all its simultaneous flows.
- The network is not responsible for fair sharing of premium capacity between microflows.
- Only a small fraction of the total bandwidth of the high speed links will be made available for the Premium IP service, 5% as an initial estimate or less. The choice minimises the probability of instantaneous burstiness at aggregation nodes, which leads to packet loss, and avoids any possibility of starving Best Effort traffic on lower speed links.
- Policing will be performed by means of a token bucket. Token bucket depth will be chosen larger than one MTU in the core and about two closest to the source. This choice is made to avoid, as much as possible, any packet loss, at the price of a small increase of delay variation and it is supported by experimental evidence [QTP-D6.2].
- Admission to the Premium IP service will be based, at the border nearest to the source, on both IP source and destination prefixes and packets will be policed according to the agreed sending rate. Packets exceeding the agreed sending rate will be discarded. In the core packet will be served according to the QoS tag (DSCP or IP Precedence), "trusting" the ingress domain and performing a less stringent policing for safety reason only. The admission control can be based also on other parameters, as defined case by case. A particular case is that the source is capable of tagging the packets and admission is then granted only when the tag is present.
- Packet admitted to the Premium IP service will be marked with a DSCP or IP Precedence value which is strongly suggested to be equal in all involved domains.
- There will be no policing and shaping applied at egress from a domain. The above-described choices will ensure that egress Premium IP traffic will not exceed the total agreed capacity.
- The link between Diffserv domains is required to behave according to the EF PHB.
- The Premium IP service is aimed at providing end to end QoS. To fulfil this goal the establishment of a particular service instance, for example between a node in L1 and a node in L2, must be made known to all domains involved. The service must be defined both as an end to end service level agreement and be accepted as a modification in the chain of service level agreements between all involved domains. For example the capacity requested between node in L1 and a node in L2 will be seen by domains N1, N2 and CORE as an increase of the aggregated premium capacity agreed between them.

The architecture is considered scalable up to one or two hundredth border links for each domain. Scaling to higher number of links or very large number of simultaneous premium flows in the same nodes of the order of many thousands requires additional investigation and possibly dedicated hardware not yet available.

Figure 2 contains a pictorial view of a link that carries both Premium IP and Best Effort traffic in a way to show the architecture key features.

The effect of the highest priority queuing can be thought as filling the link with EF traffic from top, while Best Effort traffic fills the link from bottom. Due to the scheduling mechanism, the EF traffic has always the precedence on BE traffic, even in case of bursts. The choice of assigning a small

percentage of total link to EF traffic ensures that the BE traffic is never starved and that EF traffic can use spare capacity in case of burstiness, avoiding packet loss.

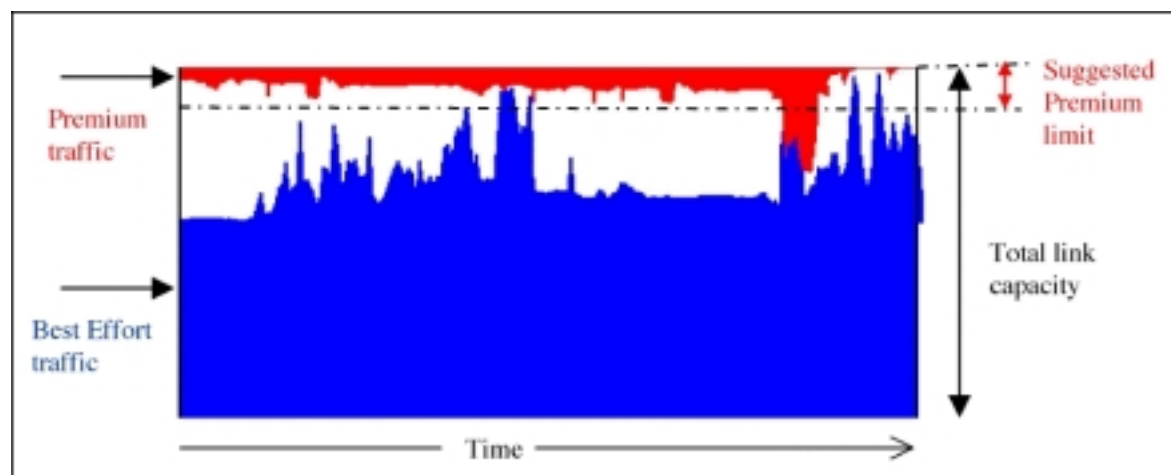


Fig. 2 Pictorial description of a link loaded with EF and BE traffic

3 DETAILED ARCHITECTURE SPECIFICATION.

3.1 INTRODUCTION

A detailed specification will be provided for each Premium IP service component that clarifies the minimum set of features required from hardware. The actions needed at each node in the network will be listed.

3.2 SERVICE COMPONENTS SPECIFICATION

3.2.1 Shaping

The initial shaping of each flow is mandatory for the successfulness of the service. Shaping is intended here as limiting the rate of transmission of data to a specific value. As the Premium IP service performs a strict policing at ingress on the agreed service sending rate, the burst size should be ideally null or equal to one MTU at most.

The sending source is required by this architecture to produce a traffic destined to Premium IP service that conforms to agreed service specification for sending rate.

The preferred way of producing a well-shaped flow from a host is by enabling shaping inside the application and/or the operating system itself. This procedure allows to avoid packet loss, as the internal feedback loop, prevents the application to send more data when the internal buffer space is exhausted.

Shaping can be performed externally to the traffic source. But there is a high probability of packet loss due to policing, unless the application intrinsically shapes traffic or the requested Premium IP sending rate is much larger than the mean send rate or the physical capacity link of the sender is lower than the agreed premium rate.

A host send for example UDP packets at maximum link rate, as soon as they are ready for transmission and TCP traffic is intrinsically bursty to probe for congestion and hence leads to packet loss if policing is performed outside the sending host.

In case the sending host is not capable of shaping, the knowledge of the mean rate and burstiness profile of the traffic generated by the application and the operating system is required to agree on the best value of the Premium IP service sending rate.

3.2.2 Fair sharing between Elastic and Anelastic flows

The fairness of the sharing of agreed sending rate between flows transported using elastic and anelastic protocols, like a mixture of TCP and UDP flows, is not considered a responsibility of the network, exactly like the responsibility of the shaping of the initial flow. If shaping is correctly implemented in the source nodes the issue is of minimal relevance.

3.2.3 Policing

Microflow policing should be done as close as possible to the source of the flow, in the first Diffserv domain, using the agreed sending rate.

Policing will be done the first time using a token bucket of minimal depth of two MTU. In case already at the initial policer multiple flows are expected, a depth of three or more MTU is suggested. Packet exceeding the allowed sending rate will be discarded.

It is suggested that policing be performed only on ingress traffic and never on egress traffic. When entering successive Diffserv domain the policing function based on aggregate EF capacity can be relaxed, according to the relevant agreements.

In case of the GÉANT domain, a possible implementation is to enable at its border a policing based on Premium IP DSCP value and the aggregate EF capacity per each pair of source and destination NRENs. This implies a simpler set of rules and higher scalability. The core domain, in this case, does not need to know the addresses of participating Premium IP end nodes, but needs to maintain a global table of agreed aggregated EF rates between each pair of NRENs which use the service. The table does not need to be symmetric. It is suggested that the rule enforce a rate limit between each pair of NRENs, which is greater than the sum of all the contracted values between each pair, as computed from the table. If no traffic rate is agreed between a particular NREN pair, but packets are encountered, the packets should be dropped, as the virtual line has an equivalent capacity of zero.

An additional extension to this simplification is that policing may be avoided on traffic coming from a trusted Diffserv domain. For example for the premium traffic which flows from GÉANT to an NREN and from a NREN backbone to an end user domain.

In case of two domains connected by more than one link, as shown in figure 1 between domain N1 and CORE, the rules have to be applied at every ingress. It is suggested that the rule sets are identical, so that, in case of routing failure, the Premium IP traffic is not affected. If the rules are based only on the IP Premium DSCP value, the sum of the allowed EF rates on the two links will be twice the agreed value if the rules set are identical. Appropriate values for the rates have to be investigated case by case, according to routing patterns.

3.2.4 Choice of token bucket depth and MTU size

In case of a network device which has multiple interfaces, each carrying Premium IP flows, it exists the possibility of a collision of Premium IP packets coming from different interfaces on the same egress bucket, even if the links are unloaded. Moreover, if the interfaces do not have the same speed,

for example when a packet flows from a higher speed link to a lower speed, a small burstiness in the high speed part might cause packet discard in the lower speed link.

This implementation requires policing only on ingress flows in selected nodes, usually only at the border, and never at egress. Hence packet loss due to egress policing is automatically avoided, although the decision to avoid both shaping and egress policing might increase burstiness.

For these reasons a depth of one full MTU is considered insufficient and experimental tests supports the conclusion [QTP-D6.2]. It is suggested that the depth be progressively increased when moving farther away from the source. Initial values can be set at 2 MTU near the source and 5 MTU at the ingress to GÉANT; larger values can be configured when the number of hops in the path becomes large.

It has to be underlined that the total depth is used only when needed and, provided a correct, limited configuration on the amount of Premium IP capacity it should be completely used only in very rare cases.

Although the most common MTU is the Ethernet one, equal to 1500 bytes, the wide area network and Gigabit Ethernet interfaces support a larger MTU value. A common choice is to have an MTU of 4470 bytes.

In any case, it is suggested that the token bucket depth is not less than 4470 bytes. If the use of this larger MTU is considered also for end nodes, the minimum granularity should be this larger value.

3.2.5 Admission control and Classification

Admission to the Premium IP service will be done as close as possible to the source host and will be based on a minimum mandatory set of parameters, which are the IP source and destination prefixes. Source and destination IP addresses can be specified as address prefixes or specific host addresses. Requiring the IP addresses is a measure to protect the service against incorrect configuration, denial of service attack or malicious users and it has to be considered as a protection for entitled users.

In the first admission control stage, packet exhibiting a DSCP or IP Precedence value equal to the Premium IP service value for that domain and not matching the rule for IP prefixes will have the tag reset to a value of zero, which is considered here as the default value and it is assigned to Best Effort traffic.

Between domains packet will be classified according to the QoS tag (DSCP or IP Precedence), "trusting" the ingress domain.

The admission control can be based also on other parameters, as defined case by case. A particular case is that the source is capable of tagging the packets and admission is then granted when the tag is present, which is discouraged on a LAN due to security concern.

Admission control and classification must be enabled on all border routers in the form of a general deny unless explicitly allowed by a rule. The general "deny" rule must be active before the service is started.

3.2.6 Marking

Although a single DSCP value for all domains is not mandatory according to Diffserv specification, an identical value on all domains is strongly suggested to avoid additional complexity.

As an example, in the case of neighbor domains, which choose different values for the Premium IP, tag there is the need of detecting and remarking crossing traffic for both values. Moreover there is the

need to remark, to tag equal to zero, any ingress traffic from the other domain with the “wrong” tag, that is the tag value chosen for Premium IP by the domain the traffic is entering.

Packets undergoing classification for the first time, which exhibit a correct tag, but wrong IP prefixes will have the tag reset to zero and will be treated as best effort.

All other packets will have the QoS bits in the IP header left untouched.

In the LAN environment 802.1p may be used instead of Diffserv to provide QoS assurances to selected flows. In this case there is the need that LAN edge routers translate between DSCP and 802.1p markings. A VLAN, coupled to 802.1p, in which only Premium IP allowed nodes are members, can provide an effective access control to the service itself

3.2.7 Scheduling

For the Premium IP service the scheduling must use the highest priority queuing algorithm available, for example Priority Queuing or Weighted Round Robin with the maximum weight on Premium IP queue.

Priority queuing has to be enabled at all nodes of involved Diffserv domains, or at least along all the relevant paths.

3.2.8 Premium IP link capacity

There is limit on the amount of capacity to devote to Premium IP, due to

- the type of service, which does not tolerate loss after initial policing, being the equivalent of a leased line;
- the choice of never starving the Best Effort traffic.

Moreover it has been shown by A.Charny and J. Le Boudec [Charny] that in a network with aggregate FIFO scheduling, for sufficiently low enough utilisation factors, deterministic delay bounds can be obtained as a function of the bound on utilisation's of any link and the maximum hop count of any flow.

It is thus suggested that the amount of Premium IP capacity subscribed does not exceed 5% of the speed link. The computation should take into account the link speed between domains and total EF rate may vary between each link. The smallness of the percentage ensures also that in case of re-routing, the service will continue to work without packet loss, albeit the delay and delay variation will be different from base values.

The admission control based on IP source and destination prefixes allows computing in each node of the domain the maximum amount of Premium IP traffic that may flow through it. It is suggested that each domain builds a matrix to compute and account the rate subscribed between each pair of its border links.

For GÉANT, for example, it will be the matrix of NREN to NREN aggregate Premium IP traffic. The matrix can in principle be asymmetrical.

3.2.9 Monitoring and accounting

Monitoring will be based on measurements of performance variables from routers and switches and active measurements of in-band Premium traffic. Measurement will concentrate on QoS parameters (delay, delay variation and packet loss) and report statistic also for usage percentages, traffic matrixes. Threshold and alarms will be set to act proactively before service degradation occurs.

The value ranges for QoS parameters may be different for each service level agreement.

Deliverable D9.4 [D9.4] [QOS-MON] of GÉANT investigates the subject.

3.3 SPECIFICATION OF FUNCTION PER NODE

The specification will be provided for a unidirectional flow from the source to the destination node.

Referring to figure 1, the picture can represent, as an example two LANs (L1 and L2) in two different countries, each LAN is connected to a different NREN backbone (N1 and N2), which in turn are connected together by the CORE GÉANT network.

Focus is given to the mandatory actions and to some of the most common possibilities. Not all the possibilities are listed. For example a LAN environment may be implemented as an Integrated Services domain, using RSVP as dynamic signaling protocol for both admission and policy propagation, but it must comply with the listed mandatory tasks: admit, mark, police, queue and in particular propagate according to the EF PHB.

3.3.1 Source node

The source node in domain L1 SHOULD perform shaping of outgoing traffic and MUST be responsible for sharing fairness of the premium capacity it is allowed to use between elastic and inelastic protocols.

The source node MAY tag the premium packets with the correct Premium IP tag value (for the domain it is in).

3.3.2 Domain L1

This is the first domain the packet encounters and usually contains the sending host and it is a LAN.

The first domain MUST perform

- as near as possible to the source
 - admission control based on IP source and destination prefixes,
 - marking of valid premium packets with agreed DSCP or IP Precedence value
 - remarking of invalid packet to best effort
 - policing according to a token bucket of depth of 2 MTU to the agreed sending rate
- enable queuing using PQ or WRR or similar queuing mechanism, with premium packets being assigned to the highest priority queue on all its border and internal routers/switches
- propagate packets inside the domain according to the EF PHB along all hops of its path
- propagate packets on links to a different domain according to the EF PHB

The domain SHOULD (in case it has multiple links to external Diffserv domains):

- police at each ingress to the domain according to a series of policers defined for each domain ingress-egress pair (aggregate policing). It is suggested to use for the policers a sending rate value greater than the contracted value, between 1.2 and two times larger and a token bucket with a depth of at least 5 MTU or more. If no traffic rate is agreed between a particular pair, but packets are encountered, the packets should be dropped, as the virtual line has an equivalent capacity of zero.

The domain MAY

- propagate the rules for aggregated traffic using signaling techniques like QoS policy propagation over BGP.
- shape the ingress traffic
- shape in selected or all transport nodes inside the domain
- shape at egress from the domain
- police at egress from the domain

3.3.3 Domain N1

This is the second domain the packet encounters.

The domain MUST perform:

- admission control based on DSCP or IP Precedence value at its border
- enable queuing using PQ or WRR or similar queuing mechanism, with premium packets being assigned to the highest priority queue on all its border and internal routers/switches
- propagate packets inside the domain according to the EF PHB along all hops of its path

- propagate packets on links to a different domain according to the EF PHB

The domain SHOULD

- police at each ingress to the domain according to a series of policers defined for each domain ingress-egress pair (aggregate policing). It is suggested to use for the policers a sending rate value greater than the contracted value, between 1.2 and two times larger and a token bucket with a depth of at least 5 MTU or more. If no traffic rate is agreed between a particular pair, but packets are encountered, the packets should be dropped, as the virtual line has an equivalent capacity of zero.

The domain MAY

- propagate the rules for aggregated traffic using signaling techniques like QoS policy propagation over BGP.
- shape the ingress traffic
- shape in selected or all transport nodes inside the domain
- shape at egress from the domain
- police at egress from the domain

3.3.4 Domain CORE

The domain tasks are identical to domain N1.

The domain MUST perform:

- admission control based on DSCP or IP Precedence value at its border
- enable queuing using PQ or WRR or similar queuing mechanism, with premium packets being assigned to the highest priority queue on all its border and internal routers/switches
- propagate packets inside the domain according to the EF PHB along all hops of its path
- propagate packets on links to a different domain according to the EF PHB

The domain SHOULD

- police at each ingress to the domain according to a series of policers defined for each domain ingress-egress pair (aggregate policing). It is suggested to use for the policers a sending rate value greater than the contracted value, between 1.2 and two times larger and a token bucket with a depth of at least 5 MTU or more. If no traffic rate is agreed between a particular pair, but packets are encountered, the packets should be dropped, as the virtual line has an equivalent capacity of zero.

The domain MAY

- propagate the rules for aggregated traffic using signaling techniques like QoS policy propagation over BGP.
- shape the ingress traffic
- shape in selected or all transport nodes inside the domain
- shape at egress from the domain
- police at egress from the domain

3.3.5 Domain N2

This domain receives packet from a core "trusted" domain and the traffic may have accumulated a small amount of burstiness, for these reasons it is suggest that the domain performs a very limited amount of checks on the ingress premium traffic.

The domain MUST perform:

- admission control based on DSCP or IP Precedence value at its border
- enable queuing using PQ or WRR or similar queuing mechanism, with premium packets being assigned to the highest priority queue on all its border and internal routers/switches

- propagate packets inside the domain using the EF PHB along all hops of its path
- propagate packets on links with a different domain according to the EF PHB

The domain **MAY**

- police at each ingress to the domain according to a series of policers defined for each domain ingress-egress pair (aggregate policing). It is suggested to use for the policers a sending rate value greater than the contracted value, between 1.5 and two times larger and a token bucket with a depth of at least 7 MTU or more. If no traffic rate is agreed between a particular pair, but packets are encountered, the packets should be dropped, as the virtual line has an equivalent capacity of zero.
- propagate the rules for aggregated traffic using signaling techniques like QoS policy propagation over BGP.
- also require a valid IP source and destination prefixes pair at ingress to the domain

The domain **SHOULD AVOID**, unless required by experimental evidence to:

- shape the ingress traffic at the border and enable shaping inside the domain or at its egresses
- police at egress from the domain

3.3.6 Domain L2

The domain **MUST** perform:

- admission control based on DSCP or IP Precedence value at its border
- enable queuing using PQ or WRR or similar queuing mechanism, with premium packets being assigned to the highest priority queue on all its border and internal routers/switches
- propagate packets inside the domain using the EF PHB along all hops of its path
- propagate packets on links with a different domain according to the EF PHB

The domain **MAY**

- police at each ingress to the domain according to a series of policers defined for each domain ingress-egress pair (aggregate policing). It is suggested to use for the policers a sending rate value greater than the contracted value, between 1.5 and two times larger and a token bucket with a depth of at least 9 MTU or more. If no traffic rate is agreed between a particular pair, but packets are encountered, the packets should be dropped, as the virtual line has an equivalent capacity of zero.
- propagate the rules for aggregated traffic using signaling techniques like QoS policy propagation over BGP.

The domain **SHOULD AVOID**, unless required by experimental evidence to:

- shape the ingress traffic at the border and enable shaping inside the domain or at its egresses
- police at egress from the domain

4 PRACTICAL CONSIDERATIONS

This section contains some practical suggestion to simplify the implementation, ensure compatibility between domains or to signal alternatives:

- The basic unit of measure is bits per seconds.
- The actions at the same node in the network can be performed by more than one hardware box, for example an access router and a core router, the first to classify, mark and police, the second to switch/route packets. The link between the boxes has to behave according to EF specification.
- The LAN environment can implement the EF behaviour in many ways. For example, the EF host can use dedicated wiring to connected to the Diffserv border router or be a member of a VLAN with appropriate QoS guarantees, for example using 802.1p.
- Hosts that provide key services at the application layer to other hosts that use the Premium IP service should also benefit of Premium IP. For example an MCU for videoconferences or a DNS server.
- The hardware capabilities should match the link speed. On low speed links (less than one Gigabit) inexpensive hardware, in the form of a router or a switch, can provide an EF PHB using standard mechanism, like committed access rate or absolute priority queueing.
- A simple way to provide a flow with a maximum sending rate of 10 Mb/s (or 100 Mb/s) is to send the traffic through a physical network interface in the host at that speed. It is suggested to use the full duplex version to avoid link contention.
- The deployment of the service can be planned in steps, which increase the number of participating nodes in a domain; there is not the need to configure all nodes at the same time, but rather just the involved in the EF path.
- For every border node that is part of a Diffserv domain, access control **MUST** be configured at ALL interfaces, to block the use of the service by unauthorised nodes.
- During the start up period of the service, it might be easier to substitute remarking invalid packet to Best Effort with dropping, using ACL. A study is under way to understand the current presence of non zero IP Precedence or DSCP values in the traffic.

5 RISK ANALYSYS AND LIMITS

Technically speaking there is a set of risks in the path to the implementation, in particular:

- Actual capabilities of routing and switching hardware. The hardware must perform flawlessly in the core at Gigabit link speed and be capable of performing at least all of the basic set of actions.
- For links at speed lower than one Gigabit, the experimental results are not yet completed. However hints from ongoing measurement do not anticipate problems, provided that the hardware capabilities are matched to the connected links speed.

The only known limitations a scaling problem due to the hardware bounds on the number of rules to be applied to ingress traffic. The limitation is currently around few hundreds or thousands of rules per border router.

The limitation is negligible both near the user, where rules based on IP prefixes and the foreseen number of IP prefixes per site is well below the limit, and between domains, where flow aggregates are policed, but the foreseen number of domain is again below the quoted limit.

It is worth noting the architecture does not intrinsically place bounds on the number of flows or allowed hosts.

The probability of packet loss due collision of Premium IP packets that enter simultaneously the same node heading to the same token bucket of an egress interface is considered quite low due to the limit on total EF capacity per link. In the worst case, in a node with N interfaces, all handling premium

traffic, the maximum number of simultaneous EF packets is $N-1$ and this value can be used to dimension the token bucket depth.

The choice of not shaping may increase the burstiness of the traffic aggregates in the core, but should not significantly modify the shape of the single flow. Measurements are planned to quantify the effect.

The service, to be considered useful, should be implemented in a large percentage of connected domains in a reasonable time.

Experimental evidence on the usefulness of the service and safety of the implementation should be available at the start time of GÉANT.

6 SECURITY CONSIDERATIONS

Classification and policing according to IP prefixes at the first stage and aggregate policing at later stages, greatly reduces the risk of Denial of Service attacks or unauthorised use. Misuse in part of a domain should also have a limited effect on other parts or domains.

Careful configuration of access control must be applied to every interface of border routers, in the default form of a “deny” for traffic that carries the DSCP tag (with a remarking of the traffic to BE), unless explicitly allowed.

The architecture can apply more stringent checks if needed, according to hardware support and performance.

7 ACKNOWLEDGEMENTS

This work has received key input from the effort of the TF-NGN [TF-NGN] and [SEQUIN] projects. The author particularly thanks Larry Dunn of Cisco System and Simon Leinen of Switch for the fruitful discussion and in depth comments.

8 REFERENCES

- [D9.1] GN1 (GÉANT) Deliverable D9.1 - "Specification and implementation plan for a Premium IP service" 9 April 2001 - <http://www.dante.net/tf-ngn/GEA-01-032.pdf>
- [D9.4] GN1 (GÉANT) Deliverable D9.4 GEA-01-113 - "Testing of traffic measurement tools", 22-Oct-2001 - <http://www.dante.net/tf-ngn/D9.4.pdf>
- [Charny] A. Charny and J.Y. Le Boudec, "Delay bounds in a network with aggregate scheduling," in Proc. First International Workshop of Quality of future Internet Services (QofIS'2000), Sept. 25--26, 2000, Berlin, Germany
- [CIT-ITU] Citkusev L., "ITU update: IP Performance and Availability Objectives and Allocations", December 2000
- [Diffserv-WG] <http://www.ietf.org/html.charters/Diffserv-charter.html>
- [EFPHB] An Expedited Forwarding PHB, Bruce Davie, Editor, Anna Charny, Fred Baker, <http://www.ietf.org/internet-drafts/draft-ietf-Diffserv-rfc2598bis-02.txt>
- [EFSUPP] Anna Charny, ed., "Supplemental Information for the New Definition of the EF PHB" draft-ietf-Diffserv-ef-supplemental-01.txt
- [GÉANT] <http://www.dante.net/geant/index.html>
- [INT-SRV] RFC 1633 Integrated Services in the Internet Architecture: an Overview. R. Braden, D. Clark, S. Shenker. June 1994.
- [Mezger95] Mezger, K. and D. W. Petr, "Bounded Delay for Weighted Round Robin", University of Kansas, Technical Report TISL-10230-07, May 1995. <http://www.tisl.ukans.edu/lite/publications/techreports/tr-tisl-10230-08.ps>
- [QOS-MON] QoS monitoring and SLS auditing, Victor Reijs, http://www.heanet.ie/heanet/projects/nat_infrastruct/qosmonitoringtf-ngn.html
- [RFC-2205] Resource ReSerVation Protocol (RSVP) -- Version 1 Functional Specification. R. Braden, Ed. , L. Zhang, S. Berson, S. Herzog, S. Jamin. September 1997. (Updated by RFC2750)
- [RFC-2208] Resource ReSerVation Protocol (RSVP) -- Version 1 Applicability Statement Some Guidelines on Deployment. A. Mankin, Ed. , F. Baker, B. Braden, S. Bradner, M. O'Dell, A. Romanow, A. Weinrib, L. Zhang. September 1997.
- [RFC-2212] Specification of Guaranteed Quality of Service. S. Shenker, C. Partridge, R. Guerin. September 1997.
- [RFC-2474] RFC-2474: Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers. K. Nichols, S. Blake, F. Baker, D. Black. December 1998.
- [RFC-2475] RFC2475 An Architecture for Differentiated Service. S. Blake, D. Black, M. Carlson, E. Davies, Z. Wang, W. Weiss. December 1998.
- [RFC3086] K. Nichols and B. Carpenter, "Definition of Differentiated Services Per Domain Behaviors and Rules for their Specification", April 2001.
- [SEQUIN] <http://www.dante.net/sequin>
- [QTP-D6.2] Deliverable D6.2 - Report on Results of Quantum test Programme. QUA-00-015 23 June 2000. <http://www.dante.net/quantum/qtp/final-report.pdf>
- [Y-1541] ITU Study Group 13, "Revised draft Recommendation Y. 1541 'Internet protocol communication service - IP Performance and Availability Objectives and Allocations'", November 2000

9 ACRONYMS

ATM	Asynchronous Transfer Mode
CU	Currently Unused
BGP	Border Gateway Protocol
DSCP	Differentiated Services Code Point
DoS	Denial of Service
EF PHB	Expedited Forwarding Per Hop Behaviour
IETF	Internet Engineering Task Force
FIFO	First In First Out
IP	Internet Protocol
IPv4	Internet protocol version 4
IPv6	Internet Protocol version 6
ITU	International Telecommunications Unit
IPDV	IP Packet Delay Variation
IPPM	IP Performance Measurement
LAN	Local Area Network
MAC	Medium Access Control
MDRR	Modified Deficit Round Robin
MPLS	Multi Protocol Label Switching
MTU	Maximum Transfer Unit
NREN	National Research and Educational Network
PDB	Per Domain Behaviour
PHB	Per Hop behaviour
PQ	Priority Queuing
QoS	Quality of Service
RSVP	Resource Reservation Protocol
SLA	Service Level Agreement
SLS	Service Level Specification
SNMP	Simple Network Management Protocol
TCP	Transmission Control protocol
ToS	Type of Service
UDP	User Datagram Protocol
WFQ	Weighted Fair Queuing
WRED	Weighted Random Early Detection
WRR	Weighted Round Robin