**Contract Number:  IST-2000-26417**
**Project Title:         GN1 (GÉANT)**

**Deliverable D9.1 - Addendum 1**

**Implementation architecture specification for the Premium IP service**

Contractual Date:
Actual Date:
Work Package:              WP9
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 architectures leverages the scalability features of Differentiated Services and takes a pragmatic approach to balance configuration complexity and benefits.

## 1 EXECUTIVE SUMMARY

This addendum specifies the implementation architecture for the Premium IP service defined in Deliverable D 9.1. [D9.1] 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 its local user Diffserv domains.

According to Deliverable 9.1 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 interconnections, which require the interconnection to behave as an EF hop .

The architecture detailed here minimises the number of action to be performed on every packet at each node and builds an initial configuration which does not use a signalling protocol.
In particular shaping is a requirement for the incoming user flow and 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 in a short time scale, which is based on current knowledge and availability of QoS technologies.
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 it is such that 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 high speed in the core (validity of the "onion" model).

According EF PHB specifications [EFPHB], shaping is required for each  flow. Shaping is 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 addresses and agreed capacity. 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 of a different domain, 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 limit for border policing is suggested to be configured at a value larger than sum of the agreed capacities for the Premium IP service flows crossing the border . This policing is performed as a safety measure, to limit service degradation to a part of the network in case of incorrect configuration or denial of service attacks. Policing never performed at egress interfaces and is not mandatory when exiting a core backbone 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 decision of not enabling additional shaping and policing is a consequence of the highest priority given to Premium IP packets and the very high speed and over provisioned characteristics of the GÉANT and NRENs core backbone and experimental validation for domains which deploy much lower link speeds is in progress.

Monitoring of the service performance is enabled from the beginning 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 losses)

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.
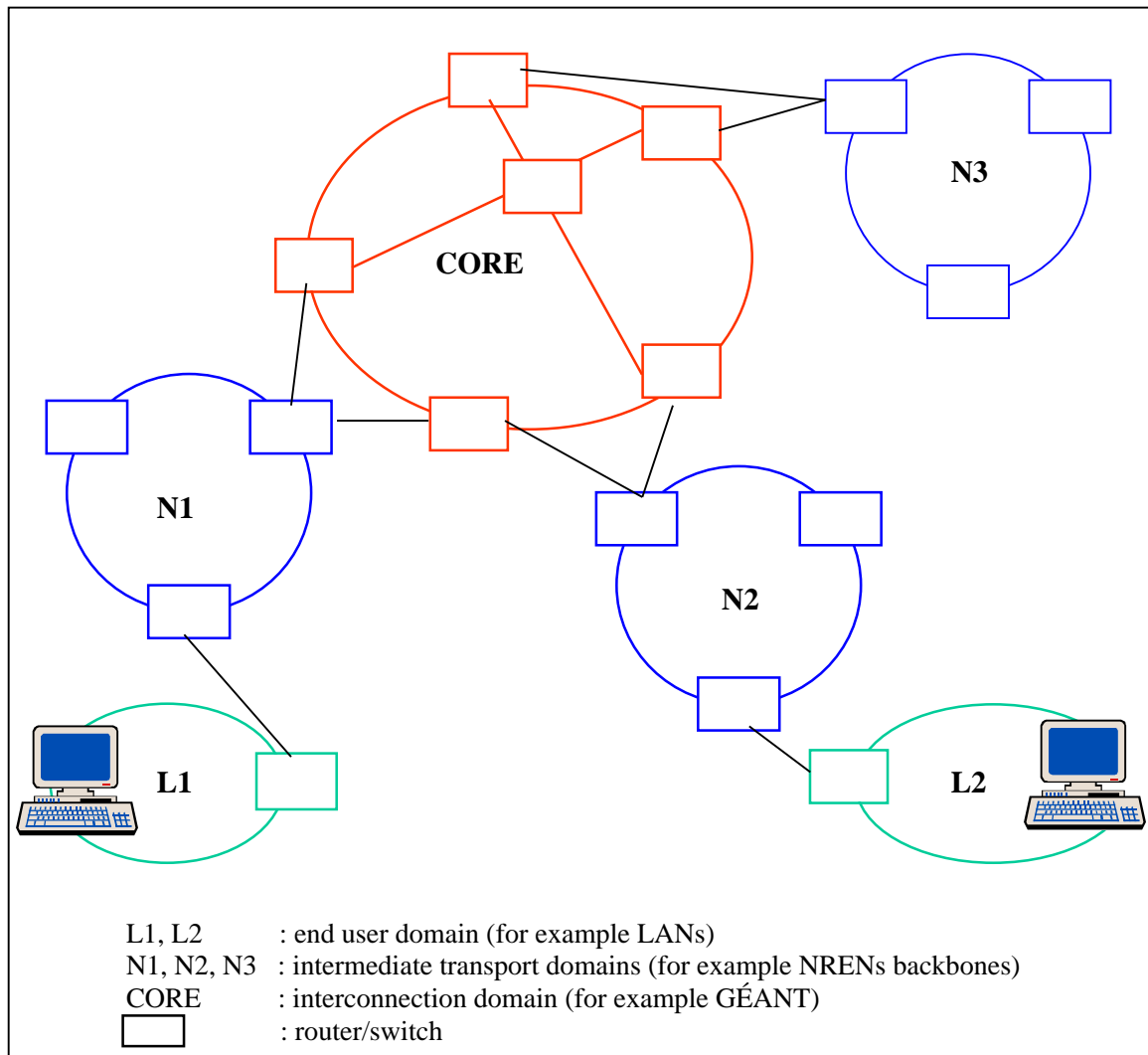


L1, L2          : end user domain (for example LANs)
N1, N2, N3   : intermediate transport domains (for example NRENs backbones)
CORE          : interconnection domain (for example GÉANT)
▭                 : router/switch

**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 signalling protocol.
The aim is at enabling a QoS Premium IP service in the same timeframe of the start of GÉANT, 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.
- 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 to start the service without enabling shaping in the core and it may render 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 allowed capacity, to avoid initial packet loss due to policing when entering the Diffserv domain, and to ensure a fair share of the 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, as an initial estimate 5% or less, will be made available for Premium IP service to minimise the probability of instantaneous burstiness at aggregation nodes and to avoid 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 packet in the core and just one or 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 IP source and destination addresses and packets will be policed according to the agreed capacity. Packets exceeding the allowed capacity 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 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 agreed capacity.
- Marking will not require a single value on all Diffserv domains, although uniformity it is strongly suggested.
- The link between 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 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.

## 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 on the agreed service capacity, 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 capacity.

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

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

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 capacity and hence leading to packet loss.

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 correctly agree on  the Premium IP service capacity.

#### 3.2.2 Fair sharing between Elastic and Anelastic flows

The fairness of the sharing of allowed capacity between flows with elastic and anelastic protocols, like a mixture of TCP and UDP flows, is not considered a responsibility of the network, like the shaping of the initial flows. If shaping is correctly implemented in the source nodes the issue is of minimal relevance.

#### 3.2.3 Policing

Microflow policing will be done as close as possible to the source of the flow, in the first Diffserv domain. The minimum set of required parameters are the IP source and IP destination addresses and the allowed capacity.
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.

Policing will be done the first time using a token bucket of minimal depth of two MTU. In case already at initial policer multiple flows are expected, a depth of three or more MTU is suggested.
Packet exceeding the allowed capacity will be discarded.
Packets with a partial match will not be discarded.
Packet exhibiting a DSCP or IP Precedence value equal to the Premium IP service value for that domain and not matching the policing rule for IP addresses 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.

It is suggested that policing is performed only on ingress traffic and never on egress traffic.
When entering successive Diffserv domain the policing function 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 capacity per each pair of 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 just needs to maintain a global table of agreed aggregated capacity between each pair of NRENs which use the service. The table does not need to be symmetric. It is suggested that the rule enforces a capacity limit which is greater than the contracted sum of all the values between each pair of NRENs, as computed from the table.

An additional extension to this simplification is that policing is not performed on traffic coming from a trusted, DiffServ domain. For example for premium traffic flowing 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 set 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 capacities on the two links will be twice the agreed value if the rules set are identical, and the appropriate values of the capacity 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 flows coming from different interfaces on the same egress bucket, even if the links are unloaded. In addition. 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 and never at egress and the 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 too low and experimental tests supports the conclusion [QTP-D6.2]. It is suggested that the depth is 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.

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, 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 IP source and destination addresses. 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 only when the tag is present, which is discouraged on a LAN due to security concern.
Admission control and classification must be enable 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, a identical value on all domain is strongly suggested.
Packets undergoing classification for the first time, which exhibit a correct tag, but wrong IP addresses will have the tag reset to zero and will be treated as best effort.
All other packets will have the IP precedence bits left untouched.

### 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 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 be 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 exceeds 5% of the speed link. The computation should take into account the link speed between domains and capacity 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 address allows to compute in each node of the domain the maximum amount of Premium IP traffic that may flow through it.. It is suggested that each domain build a matrix to compute and account the capacity subscribed between each pair of its border links.

For GÉANT, for example, it will be the matrix of NREN to NREN 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.

### 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 signalling protocol for both admission and policy propagation, but it must comply with the listed tasks (admit, mark, police, queue and propagate according to 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 address,
    - marking of premium packets with agreed DSCP or IP Precedence value
    - policing according to a token bucket of depth of 2 MTU and agreed capacity
- 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 the path
- propagate packets on links with a different domain according to EF PHB

The domain MAY
- propagate the rules for aggregated traffic using signalling techniques like QoS policy propagation over BGP.
- shape the ingress traffic
- shape in selected or all transport nodes inside the domain
- enable shaping at egress from the domain
- police at egress from the domain
- police at each ingress to the domain according to a series of polices defined for each  domain ingress-egress pair (aggregate policing)

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

The domain MAY
- propagate the rules for aggregated traffic using signalling techniques like QoS policy propagation over BGP.
- also require a valid IP source and destination address pair at ingress to the domain
- shape the ingress traffic
- shape in selected or all transport nodes inside the domain
- enable shaping at egress
- 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 using the EF PHB along all hops of the path
- propagate packets on links with a different domain according to 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 capacity 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

The domain MAY
- propagate the rules for aggregated traffic using signalling techniques like QoS policy propagation over BGP.
- also require a valid IP source and destination address pair at ingress to the domain
- shape the ingress traffic
- shape in selected or all transport nodes inside the domain
- enable shaping at egress
- 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 the path
- propagate packets on links with a different domain according to 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 capacity 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
- propagate the rules for aggregated traffic using signalling techniques like QoS policy propagation over BGP.
- also require a valid IP source and destination address 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 at egress
- 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 the path

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 capacity 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
- propagate the rules for aggregated traffic using signalling techniques like QoS policy propagation over BGP.
- also require a valid IP source and destination address 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 at egress
- police at egress from the domain

## 4 PRACTICAL RULES

This section contains some practical suggestion to simplify the implementation or to signal alternatives:
- The basic unit of measure is bits per seconds an all the relevant numbers should adopt this metric.
- 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 a VLAN with appropriate QoS guarantees can provide the connection between the EF host and the Diffserv border router.
- Being the IP name service fundamental to any application, it suggested to evaluate the importance of setting up Premium IP name server to be configured as default name server for IP Premium hosts and with its traffic marked as Premium.

## 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 at Gigabit link speed an be capable of performing at least all of the basic set of actions
- The investigation of the application of the model to much lower speed links and LANs is not yet completed. Indications from current measurement do not foresee problems, provided that the hardware is powerful enough for the connected links speed.

The only limitation so far encountered is a scaling problem due to the hardware bounds on the number of rules to be applied to ingress traffic. The limitation is currently thought to be around few hundredths or thousands of rules per border router.

The limitation implies a maximum number on incoming flows per border routers, both near sources, implying rules based on IP addresses, and between domains, where flow aggregates are policed.

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

On the general side, the service, to be considered useful, should be implemented in a large percentage of connected countries in a reasonable time. Experimental evidence on the usefulness of the service and safety of the implementation should be available at the start of the GÉANT.

## 6 SECURITY CONSIDERATIONS

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

[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-01.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. (Status: INFORMATIONAL)

[Keshav97]       Keshav, S., "An Engineering Approach to Computer Networking", Addison
                 Wesley, January 1997.

[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-1889]       RFC 1889, RTP: A Transport Protocol for Real-Time Applications

[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-2211]       Specification of the Controlled-Load Network Element Service. J. Wroclawski.
                 September 1997.

[RFC-2212]       Specification of Guaranteed Quality of Service. S. Shenker, C. Partridge, R.
                 Guerin. September 1997.

[RFC-2330]       Paxson, V. , Almes, G., Mahdavi, J. and M. Mathis, "Framework for IP
                 Performance Metrics", RFC 2330, May 1998

[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.
                 (Format: TXT=50576 bytes) (Obsoletes RFC1455, RFC1349) (Status:
                 PROPOSED STANDARD)

[RFC-2475]       RFC2475 An Architecture for Differentiated Service. S. Blake, D. Black, M.
                 Carlson, E. Davies, Z. Wang, W. Weiss. December 1998. (Status:
                 INFORMATIONAL)

[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

[Shreedhar95]        Shreedhar, M. and G. Varghese, "Efficient FairQueueing Using Deficit Round Robin", SIGCOMM 1995, pages 231-242

[Sreenivasamurthy]   Sreenivasamurthy, D., "Implementation and Evaluation of support for Differentiated Services mapping to ABR service in an Edge/Core Network", Thesis, University of Kansas.

[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

| | |
|---|---|
| ADSL | Asymmetric Digital Subscriber Link |
| ATM | Asynchronous Transfer Mode |
| CU | Currently Unused |
| BGP | Border Gateway Protocol |
| DSCP | Differentiated Services Code Point |
| DoS | Denial of Service |
| ECN | Explicit Congestion Notification |
| 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 | Instantaneous 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 |