April 2022

# QOS POLICIES FOR SERVICES ON THE SECURE INTERNET GATEWAY IN THE SDWAN DEPLOYMENTS

Niranjan M M

QOS POLICIES FOR SERVICES ON THE SECURE INTERNET GATEWAY
IN THE SDWAN DEPLOYMENTS

AUTHOR:

Niranjan M M

ABSTRACT

SD-WAN deployments in enterprise would consists of group of branch routers, which are software-defined branch (SD-Branch) routers. SD-Branch is a branch router that supports SD-WAN routing, security and other LAN access features that can be managed centrally. A thick branch is a high-end device that can incorporate all these features and provide required scale and performance for large enterprises. However, on a lean branch, not all security features can be switched on. These branch routers can integrate with third party Secure Internet Gateways (SIG) for securing the enterprise traffic. With the integration of SD-WAN and third-party SIG, all the traffic from the enterprise client's is forwarded to the SIG over the tunnel. The SD-WAN router at branch office is connected to the SIG over the WAN link and there would be bandwidth (i.e., throughput capacity) limitations for the traffic being routed over the tunnel (could be enforced by the service provider, tunnel limitations etc.,). Along with the services offered by the SIG to the SD-WAN customers, need a way to qualify customer requirements prior to onboarding, proactively monitor per-customer bandwidth consumption and provide SLA of the customer with low latency, dedicated throughput etc., Above requirements can be achieved by having QoS policies (min/max bandwidth allocation, rate limiting etc.,) for each service provided by the SIG. The techniques presented herein propose method to apply QoS policies for the services running on the SIG, so that SD-WAN deployments for the lean branch can have the same feature and functionality as that of thick branches without compromising on SLA.

DETAILED DESCRIPTION

SD-WAN deployments in enterprise would consists of group of branch routers, which are software-defined branch (SD-Branch) routers. SD-Branch is a branch router that supports SD-WAN routing, security and other LAN access features that can be managed centrally. A thick

branch is a high-end device that can incorporate all these features and provide required scale and performance for large enterprises. However, on a lean branch, not all security features can be switched on. These branch routers can integrate with third party Secure Internet Gateways (SIG) for securing the enterprise traffic.

The Direct Internet Access (DIA) traffic on SDWAN service VPN's may be tunnelled to SIG's for securing enterprise traffic. All LAN/Wi-Fi enabled enterprise client's traffic, based on routing or policy, will be forwarded to the SIG. In addition, SIG protects roaming/mobile users, BYOD use-cases as well. To achieve the same, SD-WAN branch/edge router would create IPSec/GRE tunnels to the SIG (a cloud-based security application stack) and locally breaking out internet bound packets from multiple service VPNs to this SIG through IPSec/GRE tunnel established, after which they carry-on towards their actual destinations. The return traffic is demultiplexed back to the source VPN's.

With the above integration of SD-WAN and third-party SIG, all the traffic from the enterprise client's is forwarded to the SIG over the tunnel. The SD-WAN router at branch office is connected to the SIG over the WAN link and there would be bandwidth (i.e., throughput capacity) limitations for the traffic being routed over the tunnel (could be enforced by the service provider, tunnel limitations etc.,). For example, let us say, tunnel throughput capacity is limited (measured up to 150Mbps). If branch office exceeds this limit, need to create new tunnel (and in-turn causes to have multiple tunnels from the same branch and for the same customer). With ever-increasing number of users, devices and services has an enormous impact on the over-all bandwidth (aka throughput capacity) limitations of the tunnel.  Hence, along with the services offered by the SIG to the SD-WAN customers, need a way to

- Qualify customer requirements prior to onboarding.
- Proactively monitor per-customer bandwidth consumption.
- Provide SLA of the customer with low latency, dedicated throughput etc.,

Above requirements can be achieved by having QoS policies (min/max bandwidth allocation, rate limiting etc.,) for each service provided by the SIG. Currently there are no ways to provide QoS policies to the services available on the SIG (to provide SLA to the customers with low latency, dedicated throughput etc.,).

The techniques presented herein propose method to apply QoS policies for the services running on the SIG, so that SD-WAN deployments for the lean branch can have the same feature

and functionality as that of thick branches (which are having high end device) without compromising on SLA. As part of this method, an administrator can allocate bandwidth to different services based on the requirements (ex: SLA of the customer). QoS Policing Engine (QPE) which runs on the SD-WAN Router periodically checks the bandwidth usage of every service. If any service is going beyond the allocated bandwidth, rate limiting would be applied. If bandwidth of any service is underutilised, then it would be re-distributed to the other services, which are requesting for more bandwidth. This method also considers the burst allowable limits for the sudden increase in the usages of service by the subscribers/clients.

This method considers the usage Continuous Learning (CL) component which keep track of effective bandwidth of the tunnel and services at periodic intervals. Based on the effective bandwidth usage, CL component will predict and feed QPE component with the over-commit factors. This will further increase effective bandwidth drastically. Also, this method includes the re-distribution of unused bandwidth across the services.

Bandwidth usage of the services is measured in the form of various parameters based on the type of traffic/application (real-time or non-real time), the nature of the traffic (average or burst) and the directions of the traffic (upstream or downstream). As we know services provided by the SIG mainly includes DNS, internet traffic (HTTP/HTTPS), DHCP server, Radius Server etc.,

Based on the services available, administrator can map services against the following traffic types (categories):

- Average non-real-time data rate – upstream
- Average real-time data rate – upstream
- Burst non-real-time data rate – upstream
- Burst real-time data rate – upstream
- Average non-real-time data rate – downstream
- Average real-time data rate – downstream
- Burst non-real-time data rate – downstream
- Burst real-time data rate - downstream

For example: Mapping is driven by administrator as per customer requirement.

a. DNS service can be mapped to traffic type as non-real time

b. Web service (http/https) can be mapped to traffic type as real-time

c. Radius service can be mapped to traffic type as real-time (to improve client association/join time)

Figure-1 depicts the over-all QoS policing for services between SD-WAN router and SIG.
- Multiple services would be available/provided by the SIG.
- Administrator configures bandwidth for each service (as per SLA of the customer).
- QoS Policing Engine (QPE) would be running on the SD-WAN routers and does the following functionality:
  - Communicate with each service to indicate the QoS values (i.e., bandwidth allocated) and/or rate limiting values of that service.
  - Gets the usage report back from all the services periodically.
  - Updates the CL component about the actual usages of each service.
  - Gets the feedback from CL component about the over-commit factors per each service.
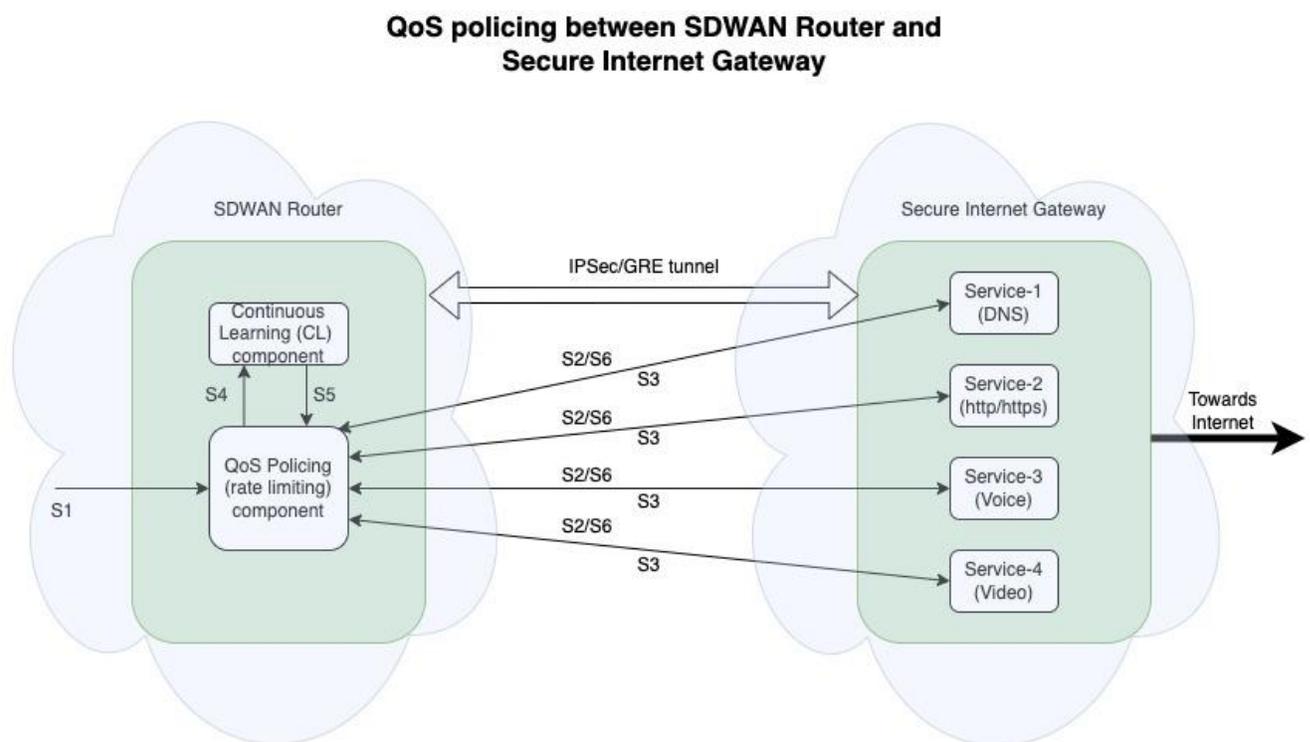


Figure-1

The techniques presented herein is explained in steps as below.

- [Step S1] As part of deployment and as per SLA, administrator configures bandwidth and/or rate limiting values for each service. QPE component uses these values initially to start with. It also considers initial over-commit factor as "1".

- [Step S2] As and when services are up and started using by SD-WAN Router, QPE component pushes the configured bandwidth and/or rate limiting values to each service.

- [Step S3] After every periodic interval, services report back to the QPE component about the actual usage values for each traffic type (aka category). They also indicate the new bandwidth requirements ("release" or "request-more"). QPE component accumulate all the released bandwidth to know the total free bandwidth.

- [Step S4] QPE component at this stage, updates the actual usages reported by each service along with timestamp to the CL component.

- [Step S5] Now, to allocate new bandwidth requirements, following set of steps repeated for each iteration for every service and for every category:

- [Step S5.1] Based on the given sample-sets, for every service, CL component predicts the probable effective bandwidth Usage Factor (bwUF).

- [Step S5.2] QPE uses configured bandwidth limit (say bwCfg) and hyperbolic it, to calculate (say bwHyp) the available bandwidth as below.

$$bwHyp = \frac{bwCfg}{bwUF}$$

This is how resultant bandwidth is over committed.

- [Step S5.3] QPE component also maintains accuracy factor for each service based on the predicted bandwidth and actual usage.

- Note: If the actual usage is around 10% plus or minus of what is predicted, consider as correct prediction. If not, consider as false positive.

(1) Accuracy factor per service shall be calculated as below:

$$Accuracy\ Factor\ (Af) = \frac{Number\ of\ Correct\ Predictions}{Total\ Predictions\ Made}$$

Note: initially Accuracy factor (Af) will be 0. Over a period, it keeps improving.

(2) We also derive 'Deterministic factor (Df)' from Accuracy factor (Af):

> Deterministic factor (Df) = 1 - Af

- [Step S5.4] In order to minimize the deviations, we compute the weighted average of bwHyp and bwCfg with Af and Df factor.

> bwAvail = (bwHyp * Af) + (bwCfg * Df)

- [Step S5.5] Finally QPE component will use bwAvail to actually allocate the new requirements for each service. It does it in two iterations:

  (Iteration-1) Allocate to all the services requested less than or equal to the configured bandwidth (bwCfg) and calculate the free bandwidth.

  (Iteration-2) Free bandwidth is allocated to the services requesting more in a proportionate manner.

- [Step S6] Allocated bandwidth (bwAlloc) values are pushed to all the services.


Example: Depicting the proposed method in allocating bandwidth and let us say 100 Mbps is the tunnel bandwidth.

a. Bandwidth Configured (bwCfg) for each Service is as below:

> bwCfg for DNS            = 20 Mbps
>
> bwCfg for DHCP         = 10 Mbps
>
> bwCfg for RADIUS       = 10 Mbps
>
> bwCfg for WEB (HTTP)   = 30 Mbps
>
> bwCfg for WEB (HTTPS)   = 30 Mbps

b. Effective Bandwidth Usage Factor (bwUF) is as below:

> bwUF for DNS            = 10 Mbps
>
> bwUF for DHCP          = 5 Mbps
>
> bwUF for RADIUS        = 5 Mbps
>
> bwUF for WEB (HTTP)    = 40 Mbps
>
> bwUF for WEB (HTTPS)    = 40 Mbps


c. Hyperbolic Bandwidth (bwHyp) to calculated available bandwidth is as below:

> bwHyp for DNS           = 20 Mbps / 10 Mbps = 2
>
> bwHyp for DHCP        = 10 Mbps / 5 Mbps = 2
>
> bwHyp for RADIUS      = 10 Mbps / 5 Mbps = 2

As WEB (HTTP an HTTPS) are more than the configured value, not used to calculate available bandwidth.

d. Af and Df calculation:

Let us say, Accuracy Factor (Af) = Number of Correct Predictions / Total predictions Made = 80 / 100 = 0.8.

Deterministic Factor (Df) = 1 - Af = 1 - 0.8 = 0.2

e. Now calculate the Bandwidth Available (bwAvail) on each service:

bwAvail for DNS     = (2 * 0.8) + (20 * 0.2) = 1.6 + 10 = 11.6 Mbps

bwAvail for DHCP    = (2 * 0.8) + (10 * 0.2) = 1.6 + 05 = 6.6 Mbps

bwAvail for RADIUS = (2 * 0.8) + (10 * 0.2) = 1.6 + 05 = 6.6 Mbps

Total 11.6 + 6.6 + 6.6 = 24.8 Mbps bandwidth available.

f. Allocate the Bandwidth (bwAlloc) to all services:

bwAlloc for DNS = 10 Mbps as requested by the service, even though configured value is 20 Mbps

bwAlloc for DHCP = 5 Mbps as requested by the service, even though configured value is 10 Mbps

bwAlloc for RADIUS = 5 Mbps as requested by the service, even though configured value is 10 Mbps

bwAlloc for WEB (HTTP) = 30 Mbps (Configured) + 10 Mbps (Requested by the service) = 40 Mbps

bwAlloc for WEB (HTTPS) = 30 Mbps (Configured) + 10 Mbps (Requested by the service) = 40 Mbps

This method allocates unused bandwidth for the services which are requesting for more if unused bandwidth is available.

In this example, 24.8 Mbps bandwidth is available as per Continuous Learning, it can be allocated to the services (HTTP and HTTPS) requesting for 10 Mbps more than configured value of 30 Mbps. If multiple services requesting for the bandwidth, then unused bandwidth will be allocated proportionately as per the requests.

The techniques presented herein provides QoS policies to the services on the SIG and hence provide SLA as per the customer requirements. This method provides efficient bandwidth allocation and rate limiting for the services. Moreover, this method helps in monitoring of bandwidth usage of each service and for each customer. Additionally, this method is applicable for any SD-WAN deployments with cloud services offered from other third party.