

ConEx
Internet-Draft
Intended status: Informational
Expires: August 28, 2013

B. Briscoe
BT
M. Sridharan
Microsoft
February 24, 2013

Network Performance Isolation in Data Centres using Congestion Policing
draft-briscoe-conex-data-centre-01

Abstract

This document describes how a multi-tenant (or multi-department) data centre operator can isolate tenants from network performance degradation due to each other's usage, but without losing the multiplexing benefits of a LAN-style network where anyone can use any amount of any resource. Zero per-tenant configuration and no implementation change is required on network equipment. Instead the solution is implemented with a simple change to the hypervisor (or container) beneath the tenant's virtual machines on every physical server connected to the network. These collectively enforce a very simple distributed contract - a single network allowance that each tenant can allocate among their virtual machines, even if distributed around the network. The solution uses layer-3 switches that support explicit congestion notification (ECN). It is best if the sending operating system supports congestion exposure (ConEx). Nonetheless, the operator can unilaterally deploy a complete solution while operating systems are being incrementally upgraded to support ConEx and ECN.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <http://datatracker.ietf.org/drafts/current/>.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on August 28, 2013.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	3
2. Features of the Solution	4
3. Outline Design	7
4. Performance Isolation: Intuition	9
4.1. Performance Isolation: The Problem	9
4.2. Why Congestion Policing Works	11
5. Design	13
5.1. Trustworthy Congestion Signals at Ingress	13
5.1.1. Tunnel Feedback vs. ConEx	14
5.1.2. ECN Recommended	14
5.1.3. Summary: Trustworthy Congestion Signals at Ingress	15
5.2. Switch/Router Support	16
5.3. Congestion Policing	17
5.4. Distributed Token Buckets	18
6. Incremental Deployment	19
6.1. Migration	19
6.2. Evolution	20
7. Related Approaches	20
8. Security Considerations	21
9. IANA Considerations (to be removed by RFC Editor)	21
10. Conclusions	21
11. Acknowledgments	21
12. Informative References	21
Appendix A. Summary of Changes between Drafts (to be removed by RFC Editor)	23

1. Introduction

A number of companies offer hosting of virtual machines on their data centre infrastructure--so-called infrastructure as a service (IaaS) or 'cloud computing'. A set amount of processing power, memory, storage and network are offered. Although processing power, memory and storage are relatively simple to allocate on the 'pay as you go' basis that has become common, the network is less easy to allocate, given it is a naturally distributed system.

This document describes how a data centre infrastructure provider can offer isolated network performance to each tenant by deploying congestion policing at every ingress to the data centre network, e.g. in all the hypervisors (or containers). The data packets pick up congestion information as they traverse the network, which is brought to the ingress using one of two approaches: feedback tunnels or ConEx (or a mix of the two). Then, these ingress congestion policers have sufficient information to limit the amount of congestion any tenant can cause anywhere in the whole meshed pool of data centre network resources. This isolates the network performance experienced by each tenant from the behaviour of all the others, without any tenant-related configuration on any of the switches.

How it works is very simple and quick to describe. Why this approach provides performance isolation may be more difficult to grasp. In particular, why it provides performance isolation across a network of links, even though there is no isolation mechanism in each link. Essentially, rather than limiting how much traffic can go where, traffic is allowed anywhere and the policer finds out whenever and wherever any traffic causes a small amount of congestion so that it can prevent heavier congestion.

This document explains how it works, while a companion document [conex-policing] builds up an intuition for why it works. Nonetheless to make this document self-contained, brief summaries of both the 'how' and the 'why' are given in sections 3 & 4. Then Section 5 gives details of the design and Section 6 explains the aspects of the design that enable incremental deployment. Finally Section 7 introduces other attempts to solve the network performance isolation problem and why they fall down in various ways.

The solution would also be just as applicable to isolate the network performance of different departments within the private data centre of an enterprise, which could be implemented without virtualisation. However, it will be described as a multi-tenant scenario, which is the more difficult case from a security point of view.

2. Features of the Solution

The following goals are met by the design, each of which is explained subsequently:

- o Performance isolation
- o No loss of LAN-like openness and multiplexing benefits
- o Zero tenant-related switch configuration
- o No change to existing switch implementations
- o Weighted performance differentiation
- o Ultra-Simple contract--per-tenant network-wide allowance
- o Sender constraint, but with transferrable allowance
- o Transport-agnostic
- o Extensible to wide-area and inter-data-centre interconnection

Performance Isolation with Openness of a LAN: The primary goal is to ensure that each tenant of a data centre receives a minimum assured performance from the whole network resource pool, but without losing the efficiency savings from multiplexed use of shared infrastructure (work-conserving). There is no need for partitioning or reservation of network resources.

Zero Tenant-Related Switch Configuration: Performance isolation is achieved with no per-tenant configuration of switches. All switch resources are potentially available to all tenants.

Separately, `_forwarding_` isolation may (or may not) be configured to ensure one tenant cannot receive traffic from another's virtual network. However, `_performance_` isolation is kept completely orthogonal, and adds nothing to the configuration complexity of the network.

No New Switch Implementation: Straightforward commodity switches (or routers) are sufficient. Bulk explicit congestion notification (ECN) is recommended, which is available in a growing range of layer-3 switches (a layer-3 switch does switching at layer-2, but it can use the Diffserv and ECN fields for traffic control if it can find an IP header).

Weighted Performance Differentiation: A tenant gets network performance in proportion to their allowance when constrained by others, with no constraint otherwise. Importantly, this assurance is not just instantaneous, but over time. And the assurance is not just localised to each link but network-wide. This will be explained later with reference to the numerical examples in [conex-policing].

Ultra-Simple Contract: The tenant needs to decide only two things: The peak bit-rate connecting each virtual machine to the network (as today) and an overall 'usage' allowance. This document focuses on the latter. A tenant just decides one number for this contracted allowance that can be shared between all the tenant's virtual machines (VMs). The 'usage' allowance is a measure of congestion-bit-rate, which will be explained later, but most tenants will just think of it as a number, where more is better.

Multi-machine: A tenant operating multiple VMs has no need to decide in advance which VMs will need more allowance and which less--an automated process can allocate the allowance across the VMs, shifting more to those that need it most, as they use it. Therefore, performance cannot be constrained by poor choice of allocations between VMs, removing a whole dimension from the problem that tenants face when choosing their traffic contract. The allocation process can be operated by the tenant, or provided by the data centre operator as part of an enhanced platform to complement the basic infrastructure (platform as a service or PaaS).

Sender Constraint with transferrable allowance: By default, constraints are always placed on data senders, determined by the sending party's traffic contract. Nonetheless, if the receiving party (or any other party) wishes to enhance performance it can arrange this with the sender at the expense of its own sending allowance.

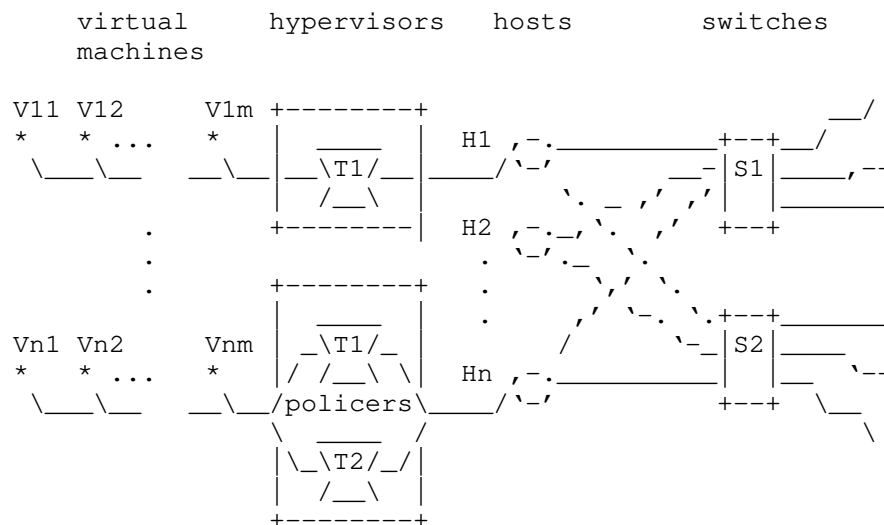
For instance, when a VM sends data to a storage facility the tenant that owns the VM consumes as much of their allowance as necessary to achieve the desired sending performance. But by default when that tenant later retrieves data from storage, the storage facility is the sender, so the storage facility consumes its allowance to determine performance in the reverse direction. Nonetheless, during the retrieval request, the storage facility can require that its sending 'costs' are covered by the receiving VM's allowance. The design of this feature is beyond the scope of this document, but the system provides all the hooks to build it at the application (or transport) layer.

Transport-Agnostic: In a well-provisioned network, enforcement of performance isolation rarely introduces constraints on network behaviour. However, it continually counts how much each tenant is limiting the performance of others, and it will intervene to enforce performance isolation against only those tenants who most persistently constrain others. By default, this intervention is oblivious to flows and to the protocols and algorithms being used above the IP layer. However, flow-aware or application-aware prioritisation can be built on top, either by the tenant or by the data centre operator as a complementary PaaS facility.

Interconnection: The solution is designed so that interconnected networks can ensure each is accountable for the performance degradation it contributes to in other networks. If necessary, one network has the information to intervene at its ingress to limit traffic from another network that is degrading performance. Alternatively, with the proposed protocols, networks can see sufficient information in traffic arriving at their borders to give their neighbours financial incentives to limit the traffic themselves.

The present document focuses on a single provider-scenario, but evolution to interconnection with other data centres over wide-area networks, and interconnection with access networks is briefly discussed in Section 6.2.

3. Outline Design



The two (or more) policers associated with tenant T1 act as one logical policer.

Figure 1: Edge Policing and the Hose Traffic Model

Edge policing: Traffic policing is located at the policy enforcement point where each sending host connects to the network, typically beneath the tenant's operating system in the hypervisor controlled by the infrastructure operator (Figure 1). In this respect, the approach has a similar arrangement to the Diffserv architecture with traffic policers forming a ring around the network [RFC2475].

(Multi-)Hose model: Each policer controls all traffic from the set of VMs associated with each tenant without regard to destination, similar to the Diffserv 'hose' model. If the tenant has VMs spread across multiple physical hosts, they are all constrained by one logical policer that feeds tokens to individual sub-policers within each hypervisor on each physical host (e.g. the two policers associated with tenant T1 in Figure 1). In other words, the network is treated as one resource pool.

Congestion policing: A congestion policer is very similar to a traditional bit-rate policer. A classifier associates each packet with the relevant tenant's meter to drain tokens from the associated token bucket, while at the same time the bucket fills with tokens at the tenant's contracted rate (Figure 2).

However, unlike a traditional policer, the tokens in a congestion policer represent congested bits (i.e. discarded or ECN-marked bits), not just any bits. So, the bits in ECN-marked packets in Figure 2 count as congested bits, while all other bits don't drain anything from the token bucket--unmarked packets are invisible to the meter. And a tenant's contracted fill rate (w_i for tenant T_i in Figure 2) is only the rate of congested bits, not all bits. Then if, on average, any tenant tries to cause more congestion than their allowance, the policer will focus discards on that tenant's traffic to prevent any further increase in congestion for everyone else.

The detail design section describes how congestion policers at the network ingress know the congestion that each packet will encounter in the network, as well as how the congestion policer limits both peak and average rates of congestion.

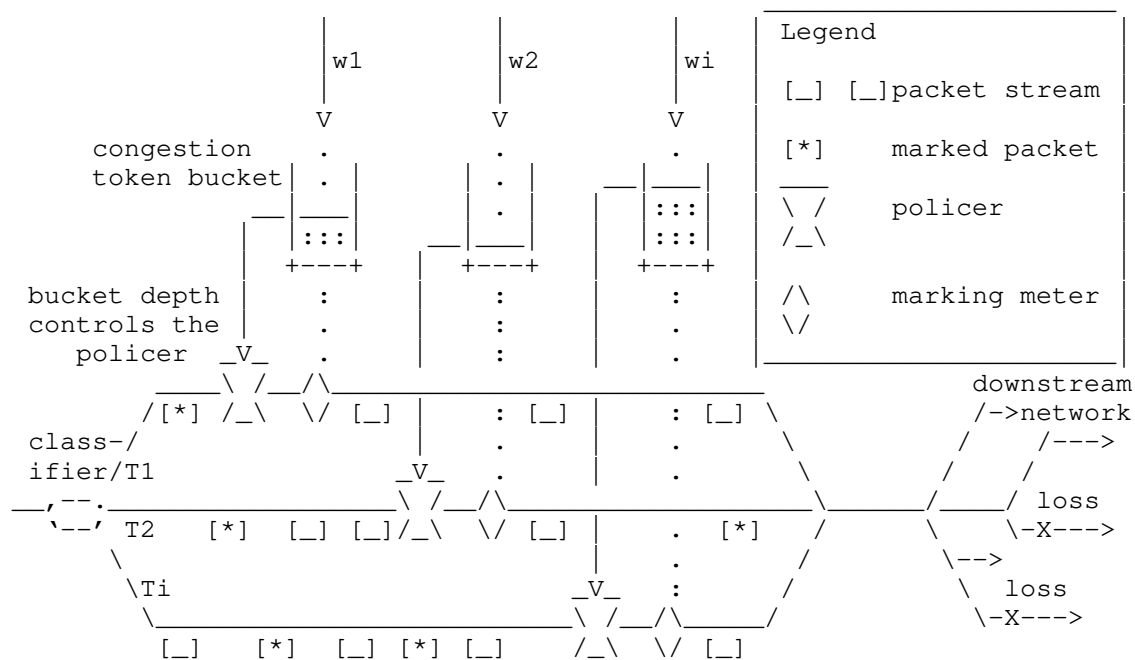


Figure 2: Bulk Congestion Policier Schematic

Optional Per-Flow policing: A congestion policer could be designed to focus policing on the particular data flow(s) contributing most to the excess congestion-bit-rate. However bulk per-tenant congestion policing is sufficient to protect other tenants, then each tenant can choose per-flow policing if it wants.

FIFO forwarding: If scheduling by traffic class is used in network buffers (for whatever reason), congestion policing can be used to isolate tenants from each other within each class. However, congestion policing will tend to keep queues short, therefore it is more likely that simple first-in first-out (FIFO) will be sufficient, with no need for any priority scheduling.

ECN marking recommended: All queues that might become congested should support bulk ECN marking. For any non-ECN-capable flows or packets, the solution enables ECN universally in the outer IP header of an edge-to-edge tunnel. It can use the edge-to-edge tunnel created by one of the network virtualisation overlay approaches, e.g. [nvgre, vxlan].

In the proposed approach, the network operator deploys capacity as usual--using previous experience to determine a reasonable contention ratio at every tier of the network. Then, the tenant contracts with the operator for the rate at which their congestion policer will allow them to contribute to congestion. [conex-policing] explains how the operator or tenant would determine an appropriate allowance.

4. Performance Isolation: Intuition

4.1. Performance Isolation: The Problem

Network performance isolation traditionally meant that each user could be sure of a minimum guaranteed bit-rate. Such assurances are useful if traffic from each tenant follows relatively predictable paths and is fairly constant. If traffic demand is more dynamic and unpredictable (both over time and across paths), minimum bit-rate assurances can still be given, but they have to be very small relative to the available capacity, because a large number of users might all want to simultaneously share any one link, even though they rarely all use it at the same time.

This either means the shared capacity has to be greatly overprovided so that the assured level is large enough, or the assured level has to be small. The former is unnecessarily expensive; the latter doesn't really give a sufficiently useful assurance.

Round robin or fair queuing are other forms of isolation that guarantee that each user will get $1/N$ of the capacity of each link, where N is the number of active users at each link. This is fine if the number of active users (N) sharing a link is fairly predictable. However, if large numbers of tenants do not typically share any one link but at any time they all could (as in a data centre), a $1/N$ assurance is fairly worthless. Again, given N is typically small but could be very large, either the shared capacity has to be expensively

overprovided, or the assured bit-rate has to be worthlessly small. The argument is no different for the weighted forms of these algorithms: WRR & WFQ).

Both these traditional forms of isolation try to give one tenant assured instantaneous bit-rate by constraining the instantaneous bit-rate of everyone else. This approach is flawed except in the special case when the load from every tenant on every link is continuous and fairly constant. The reality is usually very different: sources are on-off and the route taken varies, so that on any one link a source is more often off than on.

In these more realistic (non-constant) scenarios, the capacity available for any one tenant depends much more on how often everyone else uses a link, not just how much bit-rate everyone else would be entitled to if they did use it.

For instance, if 100 tenants are using a 1Gb/s link for 1% of the time, there is a good chance each will get the full 1Gb/s link capacity. But if just six of those tenants suddenly start using the link 50% of the time, whenever the other 94 tenants need the link, they will typically find 3 of these heavier tenants using it already. If a 1/N approach like round-robin were used, then the light tenants would suddenly get $1/4 * 1\text{Gb/s} = 250\text{Mb/s}$ on average. Round-robin cannot claim to isolate tenants from each other if they usually get 1Gb/s but sometimes they get 250Mb/s (and only 10Mb/s guaranteed in the worst case when all 100 tenants are active).

In contrast, congestion policing is the key to network performance isolation because it focuses policing only on those tenants that go fast over congested path(s) excessively and persistently over time. This keeps congestion below a design threshold everywhere so that everyone else can go fast. In this way, congestion policing takes account of highly variable loads (varying in time and varying across routes). And, if everyone's load happens to be constant, congestion policing converges on the same outcome as the traditional forms of isolation.

The other flaw in the traditional approaches to isolation, like WRR & WFQ, is that they actually prevent long-running flows from yielding to brief bursts from lighter tenants. A long-running flow can yield to brief flows and still complete nearly as soon as it would have otherwise (the brief flows complete sooner, freeing up the capacity for the longer flow sooner). However, WRR & WFQ prevent flows from even seeing the congestion signals that would allow them to co-ordinate between themselves, because they isolate each tenant completely into separate queues.

In summary, superficially, traditional approaches with separate queues sound good for isolation, but:

1. not when everyone's load is variable, so each tenant has no assurance about how many other queues there will be;
2. and not when each tenant can no longer even see the congestion signals from other tenants, so no-one's control algorithms can determine whether they would benefit most by pushing harder or yielding.

4.2. Why Congestion Policing Works

[conex-policing] explains why congestion policing works using numerical examples from a data centre and schematic traffic plots (in ASCII art). The bullets below provide a summary of that explanation, which builds from the simple case of long-running flows through a single link up to a full meshed network with on-off flows of different sizes and different behaviours:

- o Starting with the simple case of long-running flows focused on a single bottleneck link, tenants get weighted shares of the link, much like weighted round robin, but with no mechanism in any of the links. This is because losses (or ECN marks) are random, so if one tenant sends twice as much bit-rate it will suffer twice as many lost bits (or ECN-marked bits). So, at least for constant long-running flows, regulating congestion-bits gives the same outcome as regulating bits;
- o In the more realistic case where flows are not all long-running but a mix of short to very long, it is explained that bit-rate is not a sufficient metric for isolating performance; how often a tenant is sending (or not sending) is the significant factor for performance isolation, not whether bit-rate is shared equally whenever a source happens to be sending;
- o Although it might seem that data volume would be a good measure of how often a tenant is sending, we then show why it is not. For instance, a tenant can send a large volume of data but hardly affect the performance of others -- by being more responsive to congestion. Using congestion-volume (congestion-bit-rate over time) in a policer encourages large data senders to be more responsive (to yield), giving other tenants much higher performance while hardly affecting their own performance. Whereas, using straight volume as an allocation metric provides no distinction between high volume sources that yield and high volume sources that do not yield (the widespread behaviour today);

- o We then show that a policer based on the congestion-bit-rate metric works across a network of links treating it as a pool of capacity, whereas other approaches treat each link independently, which is why the proposed approach requires none of the configuration complexity on switches that is involved in other approaches.
- o We also show that a congestion policer can be arranged to limit bursts of congestion from sources that focus traffic onto a single link, even where one source may consist of a large aggregate of sources.
- o We show that a congestion policer rewards traffic that shifts to less congested paths (e.g. multipath TCP or virtual machine motion). This means congestion policing encourages and ultimately forces end-systems to balance their load over the whole pool of bandwidth. The network can attempt to balance the load, but bulk congestion policing is particularly designed to encourage end-systems to do the job, either at the transport layer with multipath TCP [RFC6356] or at the application layer by moving virtual machines or choosing peer virtual machines in a similar way to BitTorrent.
- o We show that congestion policing works on the pool of links, irrespective of whether individual links have significantly different capacities.
- o We show that a congestion policer allows a wide variety of responses to congestion (e.g. New Reno TCP, Cubic TCP, Compound TCP, Data Centre TCP and even unresponsive UDP traffic), while still encouraging and enforcing a sufficient response to congestion from all sources taken together.
- o Congestion policing can and will enforce a congestion response if a tenant persistently causes excessive congestion. This ensures that each tenant's minimum performance is isolated from the combined effects of everyone else. However, the purpose of congestion policing is not to intervene in everyone's rate control all the time. Rather it is encourage each tenant to avoid being policed -- to keep the aggregate of all their flows' responses to congestion within an overall envelope and balanced across the network.

[conex-policing] also includes a section that gives guidance on how to estimate appropriate fill rates and sizes for congestion token buckets.

5. Design

The design involves the following elements, each detailed in the following subsections:

1. Trustworthy Congestion Signals at Ingress
2. Switch/Router Support
3. Congestion Policing
4. Distributed Token Buckets

5.1. Trustworthy Congestion Signals at Ingress

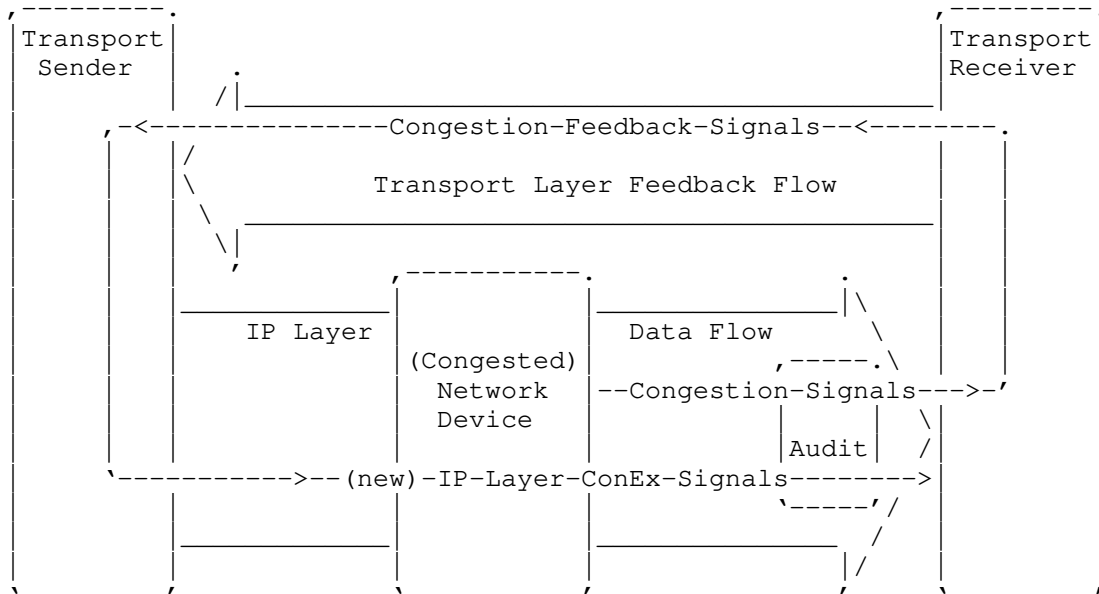


Figure 3: The ConEx Protocol in the Internet Architecture

The operator of the data centre infrastructure needs to trust this information, therefore it cannot just use the feedback in the end-to-end transport (e.g. TCP SACK or ECN echo congestion experienced flags) that might anyway be encrypted. Trusted congestion feedback may be implemented in either of the following two ways:

- a. Either as a shim in both sending and receiving hypervisors using an edge-to-edge (host-host) tunnel controlled by the infrastructure operator, with feedback messages reporting congestion back to the sending host's hypervisor (in addition to

the e2e feedback at the transport layer).

- b. Or in the sending operating system using the congestion exposure protocol (ConEx [ConEx-Abstract-Mech]) with a ConEx audit function at the egress edge to check ConEx signals against actual congestion signals (Figure 3);

5.1.1. Tunnel Feedback vs. ConEx

The feedback tunnel approach (a) is inefficient because it duplicates end-to-end feedback and it introduces at least a round trip's delay, whereas the ConEx approach (b) is more efficient and not delayed, because ConEx packets signal a conservative estimate of congestion in the upcoming round trip. Avoiding feedback delay is important for controlling congestion from aggregated short flows. However, ConEx signals will not necessarily be supported by the sending operating system.

Therefore, given ConEx IP packets are self-identifying, the best approach is to rely on ConEx signals when present and fill in with tunnelled feedback when not, on a packet-by-packet basis.

5.1.2. ECN Recommended

Both approaches are much easier if explicit congestion notification (ECN [RFC3168]) is enabled on network switches and if all packets are ECN-capable. For non-ECN-capable packets, ECN support can be turned on in the outer of an edge-to-edge tunnel. The reasons that ECN helps in each case are:

- a. Tunnel Feedback: To feed back congestion signals, the tunnel egress needs to be able to detect forward congestion signals in the first place. If the only symptom of congestion is dropped packets, the egress has to watch for gaps in the sequence space of the transport protocol, which cannot be guaranteed to be possible--the IP payload may be encrypted, or an unknown protocol, or parts of the flow may be sent over diverse paths. The tunnel ingress could add its own sequence numbers (as done by some pseudowire protocols), but it is easier to simply turn on ECN at the ingress so that the egress can detect ECN markings.
- b. ConEx: The audit function needs to be able to compare ConEx signals with actual congestion. So, as before, it needs to be able to detect congestion at the egress. Therefore the same arguments for ECN apply.

5.1.3. Summary: Trustworthy Congestion Signals at Ingress

The above cases can be arranged in a 2x2 matrix, to show when edge-to-edge tunnelling is needed and what function the tunnel would need to serve:

ConEx-capable?	ECN-capable: Y	ECN-capable: N
Y	No tunnel needed	ECN-enabled tunnel
N	Tunnel Feedback	ECN-enabled tunnel + Tunnel feedback

We can now summarise the steps necessary to ensure an ingress congestion policer obtains trustworthy congestion signals:

1. Sending operating system:

1. The sender SHOULD send ConEx-enabled and ECN-enabled packets whenever possible.
2. If the sender uses IPv6 it can signal ConEx in a destination option header [conex-destopt].
3. If the sender uses IPv4, it can signal ConEx markings by encoding them within the packet ID field as proposed in [ipv4-id-reuse].

2. Ingress edge:

1. If an arriving packet is either not ConEx-capable or not ECN-capable it SHOULD be tunnelled to the appropriate egress edge in an outer IP header.
2. A pre-existing edge-to-edge tunnel (e.g. [nvgre, vxlan]) can be used, irrespective of whether the packet is not ConEx-capable or not ECN-capable.
3. Incoming ConEx signals MUST be copied to the outer. For an incoming IPv4 packet, this implies copying the ID field. For an incoming IPv6 packet, this implies copying the Destination Option header.
4. In all cases, the tunnel ingress MUST use the normal mode of ECN tunnelling [RFC6040].

3. Directly after encapsulation (but not if the packet was not encapsulated):
 1. If and only if the ECN field of the outer header is not ECN-capable (Not-ECT i.e. 00) it MUST be made ECN-capable by remarking it to ECT(0), i.e. 01.
 2. If the outer ECN field carries any other value than 00, it should be left unchanged.
4. Directly before the edge egress (irrespective of whether the packet is encapsulated):
 1. If the outer IP header is ConEx-capable, it MUST be passed through a ConEx audit function
 2. If the packet is not ConEx-capable, it MUST be passed to a function that feeds back ECN marking statistics to the tunnel ingress. Such a function is also a requirement of [tunnel-cong-exp], which may be re-usable for this purpose {ToDo: to be confirmed}.
5. Egress Edge Decapsulator:
 1. Decapsulation must comply with [RFC6040]. This ensures that, a congestion experienced marking (CE or 11) on the outer will lead to the packet being dropped if the inner indicates that the endpoints will not understand ECN (i.e. the inner ECN field is Not-ECT or 00). Effectively the egress edge drops such packets on behalf of the congested upstream buffer that marked it because the packet appeared to be ECN-capable on the outside, but it is not ECN-capable on this inside. [RFC6040] was deliberately arranged like this so that it would drop such packets to give an equivalent congestion signal to the end-to-end transport.

5.2. Switch/Router Support

Network switches/routers do not need any modification. However, both congestion detection by the tunnel (approach a) and ConEx audit (approach b) are significantly easier if switches support ECN.

Once switches support ECN, Data centre TCP [DCTCP] could optionally be used (DCTCP requires ECN). It also requires modified sender and receiver TCP algorithms as well as a more aggressive configuration of the active queue management (AQM) in the L3 switches or routers.

5.3. Congestion Policing

Innovation in the design of congestion policers is expected and encouraged, but here we will describe one specific design to be concrete.

A bulk congestion policing function would most likely be implemented as a shim in the hypervisor. The hypervisor would create one instance of a bulk congestion policer per tenant on the physical machine, and it would ensure that all traffic sent by that tenant's VMs into the network would pass through the relevant congestion policer by associating every new virtual machine with the relevant policer.

A bulk congestion policing function has already been outlined in Section 3. To recap, it consists of a token bucket that is filled with congestion tokens at a constant rate. The bucket is drained by the size of every packet that carries a congestion marking. If the tunnel-feedback approach (a) were used, the bucket would be drained by congestion feedback from the tunnel egress, rather than markings on packets. If the ConEx approach (b) were used, the bucket would be drained by ConEx markings on the actual data packets being forwarded. A congestion policer will need to drain in response to either form of signal, because it is recommended that both approaches are used in combination.

Various more sophisticated congestion policer designs have been evaluated [CPolTrilogyExp]. In these experiments, it was found that it is better if the policer gradually increases discards as the bucket becomes empty. Also isolation between tenants is better if each tenant is policed based on the combination of two buckets, not one (Figure 4):

1. A deep bucket (that would take minutes or even hours to fill at the contracted fill-rate) that constrains the tenant's long-term average rate of congestion (w_i)
2. a very shallow bucket (e.g. only two or three MTU) that is filled considerably faster than the deep bucket ($c * w_i$), where $c = \sim 10$, which prevents a tenant storing up a large backlog of tokens then causing congestion in one large burst.

In this arrangement each marked packet drains tokens from both buckets, and the probability of policer discard is taken as the worse of the two buckets.

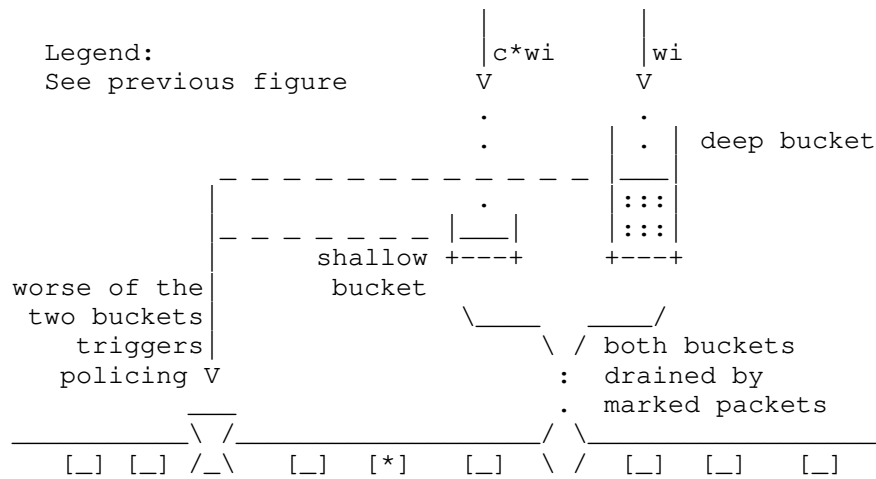


Figure 4: Dual Congestion Token Bucket (in place of each single bucket in the previous figure)

While the data centre network operator only needs to police congestion in bulk, tenants may wish to enforce their own limits on individual users or applications, as sub-limits of their overall allowance. Given all the information used for policing is readily available within the transport layer of their own operating system. Tenants can readily apply any such per-flow, per-user or per-application limitations. The tenant may operate their own fine-grained policing software, or such detailed control capabilities may be offered as part of the platform (platform as a service or PaaS).

5.4. Distributed Token Buckets

A customer may run virtual machines on multiple physical nodes, in which case at the time each VM is instantiated the data centre operator will deploy a congestion policer in the hypervisor on each node where the customer is running a VM. The DC operator can arrange for these congestion policers to collectively enforce the per-customer congestion allowance, as a distributed policer.

A function to distribute a customer's tokens to the policer associated with each of the customer's VMs would be needed. This could be similar to the distributed rate limiting of [DRL], which uses a gossip-like protocol to fill the sub-buckets. Alternatively, a logically centralised bucket of congestion tokens could be used. It could be replicated for reliability then there could be simple 1-1 communication between the central bucket and each local token bucket.

Importantly, congestion tokens can be freely reassigned between different VMs, because a congestion token is equivalent at any place or time in a network. In contrast, traditional bit-rate tokens cannot simply be reassigned from one VM to another without implications on the balance of network loading. This is because the parameters used for bit-rate policing depend on the topology and its capacity planning (open loop), whereas congestion policing complements the closed loop congestion avoidance system that adapts to the prevailing traffic and topology.

As well as distribution of tokens between the VMs of a tenant, it would similarly be feasible to allow transfer of tokens between tenants, also without breaking the performance isolation properties of the system. Secure token transfer mechanisms could be built above the underlying policing design described here, but that is beyond the current scope and therefore deferred to future work.

6. Incremental Deployment

6.1. Migration

A mechanism to bring trustworthy congestion signals to the ingress (Section 5.1) is critical to this performance isolation solution. Section 5.1.1 compares the two solutions: b) ConEx, which is efficient and it's timely enough to police short flows; and a) tunnel-feedback, which is neither. However, ConEx requires deployment in host operating systems first, while tunnel feedback can be deployed unilaterally by the data centre operator in all hypervisors (or containers), without requiring support in guest operating systems.

The section describes the steps necessary to support both approaches. This would provide an incremental deployment route with the best of both worlds: tunnel feedback could be deployed initially for unmodified guest OSs despite its weaknesses, and ConEx could gradually take over as it was deployed more widely in guest OSs. It is important not to deploy the tunnel feedback approach without checking for ConEx-capable packets, otherwise it will never be possible to migrate to ConEx. The advantages of being able to migrate to ConEx are:

- o no duplicate feedback channel between hypervisors (sending and forwarding a large proportion of tiny packets), which would cause considerable packet processing overhead
- o performance isolation includes the contribution to congestion from short (sub-round-trip-time) flows

6.2. Evolution

Initially, the approach would be confined to intra-data centre traffic. With the addition of ECN support on network equipment (at least bottleneck access routers) in the WAN between data centres, it could straightforwardly be extended to inter-data centre scenarios, including across interconnected backbone networks.

Once this approach becomes deployed within data centres and possibly across interconnects between data centres and enterprise LANs, the necessary support will be implemented in a wide range of equipment used in these scenarios. Similar equipment is also used in other networks (e.g. broadband access and backhaul), so that it would start to be possible for these other networks to deploy a similar approach.

7. Related Approaches

The Related Work section of [CongPol] provides a useful comparison of the approach proposed here against other attempts to solve similar problems.

When the hose model is used with Diffserv, capacity has to be considerably over-provisioned for all the unfortunate cases when multiple sources of traffic happen to coincide even though they are all in-contract at their respective ingress policers. Even so, every node within a Diffserv network also has to be configured to limit higher traffic classes to a maximum rate in case of really unusual traffic distributions that would starve lower priority classes. Therefore, for really important performance assurances, Diffserv is used in the 'pipe' model where the policer constrains traffic separately for each destination, and sufficient capacity is provided at each network node for the sum of all the peak contracted rates for paths crossing that node.

In contrast, the congestion policing approach is designed to give full performance assurances across a meshed network (the hose model), without having to divide a network up into pipes. If an unexpected distribution of traffic from all sources focuses on a congestion hotspot, it will increase the congestion-bit-rate seen by the policers of all sources contributing to the hot-spot. The congestion policers then focus on these sources, which in turn limits the severity of the hot-spot.

The critical improvement over Diffserv is that the ingress edges receive information about any congestion occurring in the middle, so they can limit how much congestion occurs, wherever it happens to occur. Previously Diffserv edge policers had to limit traffic generally in case it caused congestion, because they never knew

whether it would (open loop control).

Congestion policing mechanisms could be used to assure the performance of one data flow (the 'pipe' model), but this would involve unnecessary complexity, given the approach works well for the 'hose' model.

Therefore, congestion policing allows capacity to be provisioned for the average case, not for the near-worst case when many unlikely cases coincide. It assures performance for all traffic using just one traffic class, whereas Diffserv only assures performance for a small proportion of traffic by partitioning it off into higher priority classes and over-provisioning relative to the traffic contracts sold for for this class.

{ToDo: Refer to [conex-policing] for comparison with WRR & WFQ}

Seawall {ToDo} [Seawall]

8. Security Considerations

{ToDo}

9. IANA Considerations (to be removed by RFC Editor)

This document does not require actions by IANA.

10. Conclusions

{ToDo}

11. Acknowledgments

Thanks to Yu-Shun Wang for comments on some of the practicalities.

Bob Briscoe is part-funded by the European Community under its Seventh Framework Programme through the Trilogy 2 project (ICT-317756). The views expressed here are solely those of the author.

12. Informative References

[CPolTrilogyExp] Raiciu, C., Ed., "Progress on resource control", Trilogy EU 7th Framework Project ICT-216372 Deliverable 9, December 2009, <<http://trilogy-project.org/publications/deliverables.html>>.

[ConEx-Abstract-Mech] Mathis, M. and B. Briscoe, "Congestion

- Exposure (ConEx) Concepts and Abstract Mechanism", draft-ietf-conex-abstract-mech-03 (work in progress), October 2011.
- [CongPol] Jacquet, A., Briscoe, B., and T. Moncaster, "Policing Freedom to Use the Internet Resource Pool", Proc ACM Workshop on Re-Architecting the Internet (ReArch'08) , December 2008, <<http://bobbriscoe.net/projects/refb/#polfree>>.
- [DCTCP] Alizadeh, M., Greenberg, A., Maltz, D., Padhye, J., Patel, P., Prabhakar, B., Sengupta, S., and M. Sridharan, "Data Center TCP (DCTCP)", ACM SIGCOMM CCR 40(4)63--74, October 2010, <<http://portal.acm.org/citation.cfm?id=1851192>>.
- [DRL] Raghavan, B., Vishwanath, K., Ramabhadran, S., Yocum, K., and A. Snoeren, "Cloud control with distributed rate limiting", ACM SIGCOMM CCR 37(4)337--348, 2007, <<http://doi.acm.org/10.1145/1282427.1282419>>.
- [RFC2475] Blake, S., Black, D., Carlson, M., Davies, E., Wang, Z., and W. Weiss, "An Architecture for Differentiated Services", RFC 2475, December 1998.
- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", RFC 3168, September 2001.
- [RFC6040] Briscoe, B., "Tunnelling of Explicit Congestion Notification", RFC 6040, November 2010.
- [RFC6356] Raiciu, C., Handley, M., and D. Wischik, "Coupled Congestion Control for Multipath Transport Protocols", RFC 6356, October 2011.
- [Seawall] Shieh, A., Kandula, S., Greenberg, A., and C. Kim, "Seawall: Performance Isolation in Cloud Datacenter Networks", Proc 2nd USENIX Workshop on Hot Topics in Cloud Computing , June 2010, <<http://research.microsoft.com/en-us/projects/seawall/>>.

[conex-destopt]	Krishnan, S., Kuehlewind, M., and C. Ucendo, "IPv6 Destination Option for ConEx", draft-ietf-conex-destopt-03 (work in progress), September 2012.
[conex-policing]	Briscoe, B., "Network Performance Isolation using Congestion Policing", draft-briscoe-conex-policing-00 (work in progress), February 2013.
[ipv4-id-reuse]	Briscoe, B., "Reusing the IPv4 Identification Field in Atomic Packets", draft-briscoe-intarea-ipv4-id-reuse-02 (work in progress), October 2012.
[nvgre]	Sridhavan, M., Greenberg, A., Venkataramaiah, N., Wang, Y., Duda, K., Ganga, I., Lin, G., Pearson, M., Thaler, P., and C. Tumuluri, "NVGRE: Network Virtualization using Generic Routing Encapsulation", draft-sridharan-virtualization-nvgre-01 (work in progress), July 2012.
[tunnel-cong-exp]	Zhu, L., Zhang, H., and X. Gong, "Tunnel Congestion Exposure", draft-zhang-tsvwg-tunnel-congestion-exposure-00 (work in progress), October 2012.
[vxlan]	Mahalingam, M., Dutt, D., Duda, K., Agarwal, P., Kreeger, L., Sridhar, T., Bursell, M., and C. Wright, "VXLAN: A Framework for Overlaying Virtualized Layer 2 Networks over Layer 3 Networks", draft-mahalingam-dutt-dcops-vxlan-03 (work in progress), February 2013.

Appendix A. Summary of Changes between Drafts (to be removed by RFC Editor)

Detailed changes are available from
<http://tools.ietf.org/html/draft-briscoe-conex-data-centre>

From briscoe-conex-data-centre-00 to briscoe-conex-data-centre-01:

- * Took out text Section 4 "Performance Isolation Intuition" and Section 6. "Parameter Setting" into a separate draft [conex-policing] and instead included only a summary in these sections, referring out for details.

- * Considerably updated Section 5 "Design"
- * Clarifications and updates throughout, including addition of diagrams

From briscoe-conex-initial-deploy-02 to
briscoe-conex-data-centre-00:

- * Split off data-centre scenario as a separate document, by popular request.

Authors' Addresses

Bob Briscoe
BT
B54/77, Adastral Park
Martlesham Heath
Ipswich IP5 3RE
UK

Phone: +44 1473 645196
EMail: bob.briscoe@bt.com
URI: <http://bobbriscoe.net/>

Murari Sridharan
Microsoft
1 Microsoft Way
Redmond, WA 98052

Phone:
Fax:
EMail: muraris@microsoft.com
URI:

ConEx
Internet-Draft
Intended status: Informational
Expires: August 22, 2013

B. Briscoe
BT
February 18, 2013

Network Performance Isolation using Congestion Policing
draft-briscoe-conex-policing-00

Abstract

This document describes why policing using congestion information can isolate users from network performance degradation due to each other's usage, but without losing the multiplexing benefits of a LAN-style network where anyone can use any amount of any resource. Extensive numerical examples and diagrams are given. The document is agnostic to how the congestion information reaches the policer. The congestion exposure (ConEX) protocol is recommended, but other tunnel feedback mechanisms have been proposed.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <http://datatracker.ietf.org/drafts/current/>.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on August 22, 2013.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must

include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1. Introduction	3
2. Example Bulk Congestion Policer	4
3. Network Performance Isolation: Intuition	8
3.1. The Problem	8
3.2. Approach	10
3.3. Terminology	11
3.4. Simple Boundary Model of Congestion Control	11
3.5. Long-Running Flows	12
3.6. On-Off Flows	14
3.6.1. Numerical Examples Without Policing	15
3.6.2. Congestion Policing of On-Off Flows	19
3.7. Weighted Congestion Controls	20
3.8. A Network of Links	22
3.8.1. Numerical Example: Isolation from Focused Load	23
3.8.2. Isolation in the Short-Term	24
3.8.3. Encouraging Load Balancing	24
3.9. Tenants of Different Sizes	25
3.10. Links of Different Sizes	25
3.11. Diverse Congestion Control Algorithms	26
4. Parameter Setting	28
5. Security Considerations	29
6. IANA Considerations	30
7. Conclusions	30
8. Acknowledgments	30
9. Informative References	30
Appendix A. Summary of Changes between Drafts	32

1. Introduction

This document is informative, not normative. It describes why congestion policing isolates the network performance of different parties who are using a network, and why it is more useful than other forms of performance isolation such as peak and committed information rate policing, weighted round-robin or weighted fair-queueing.

The main point is that congestion policing isolates performance even when the load applied by everyone is highly variable, where variation may be over time, or by spreading traffic across different paths (the hose model). Unlike other isolation techniques, congestion policing is an edge mechanism that works over a pool of links across whole network, not just at one link. If the load from everyone happens to be constant, and there is a single bottleneck, congestion policing gives a nearly identical outcome to weighted round-robin or weighted fair-queueing at that bottleneck. But otherwise, congestion policing is a far more general solution.

The document complements others that describe various network arrangements that could use congestion policing, such as in a broadband access network [I-D.briscoe-conex-initial-deploy], in a mobile access network [I-D.ietf-conex-mobile] or in a data centre network [I-D.briscoe-conex-data-centre]. Nonetheless, all the numerical examples use the data centre scenario.

The key to the solution is the use of congestion-bit-rate rather than bit-rate as the policing metric. How this works is very simple and quick to describe Section 2.

However, it is much more difficult to understand why this approach provides performance isolation. In particular, why it provides performance isolation across a network of links, even though there is apparently no isolation mechanism in each link. Section 3 builds up an intuition for why the approach works, and why other approaches fall down in different ways. The explanation builds as follows:

- o Starting with the simple case of long-running flows focused any one bottleneck link in the network, tenants get weighted shares of the link, much like weighted round robin, but with no mechanism in any of the links;
- o In the more realistic case where flows are not all long-running but a mix of short to very long, it is explained that bit-rate is not a sufficient metric for isolating performance; how often a tenant is not sending is the significant factor for performance isolation, not whether bit-rate is shared equally whenever a source happens to be sending;

- o Although it might seem that data volume would be a good measure of how often a tenant does not send, we then show that a tenant can send a large volume of data but hardly affect the performance of others -- by being more responsive to congestion. Using congestion-volume (congestion-bit-rate over time) in a policer encourages large data senders to do this, giving other tenants much higher performance while hardly affecting their own performance. Whereas, using straight volume as an allocation metric provides no distinction between high volume sources that yield and high volume sources that do not yield (the widespread behaviour today);
- o We then show that a policer based on the congestion-bit-rate metric works across a network of links treating it as a pool of capacity, whereas other approaches treat each link independently, which is why the proposed approach requires none of the configuration complexity on switches that is involved in other approaches.
- o We also show that a congestion policer can be arranged to limit bursts of congestion from sources that focus traffic onto a single link, even where one source may consist of a large aggregate of sources.
- o We show that a congestion policer rewards traffic that shifts to less congested paths (e.g. multipath TCP or virtual machine motion).
- o We show that congestion policing works on the pool of links, irrespective of whether individual links have significantly different capacities.
- o We show that a congestion policer allows a wide variety of responses to congestion (e.g. New Reno TCP [RFC5681], Cubic TCP, Compound TCP, Data Centre TCP [DCTCP] and even unresponsive UDP traffic), while still encouraging and enforcing a sufficient response to congestion from all sources taken together, so that the performance of each application is sufficiently isolated from others.

2. Example Bulk Congestion Policer

In order to explain why a congestion policer works, we first describe a particular design of congestion policer, called a bulk congestion policer. The aim is to give a concrete example to motivate the protocol changes necessary to implement such a policer. This is merely an informative document, so there is no intention to standardise this or any specific congestion policing algorithm. The

aim is for different implementers to be free to develop their own improvements to this design.

A number of other documents describe deployment arrangements that include a congestion policer, for instance in a broadband access network [I-D.briscoe-conex-initial-deploy], in a mobile access network [I-D.ietf-conex-mobile] and in a data centre network [I-D.briscoe-conex-data-centre]. The congestion policer described in the present document could be deployed in any of these arrangements. However, to be concrete, the example link capacities, topology and terminology used in the present document assume a multi-tenant data centre scenario [I-D.briscoe-conex-data-centre]. It is believed that similar examples could be drawn up for the other scenarios, just using different numbers for link capacities, etc and replacing some terminology, e.g. substituting 'end-customer' for 'tenant' throughout.

Figure 1 illustrates a bulk congestion policer. It is very similar to a regular token bucket for policing bit-rate except the tokens do not represent bits, rather they represent 'congestion-bits'. Using Congestion Exposure (ConEx) [I-D.ietf-conex-abstract-mech] is probably the easiest way to understand what congestion-bits are. However, there are other valid ways to signal congestion information to a congestion policer without ConEx (e.g. tunnelled feedback described in [I-D.briscoe-conex-data-centre]).

Using ConEx as an example, a ConEx-based congestion token bucket ignores any packets unless they are ConEx-marked. The ConEx protocol includes audit capabilities that force the sender to reveal information it receives about losses anywhere on the downstream path (e.g. at the points marked X in Figure 1). The sender signals this information with a ConEx-marking in the IP header of the packets it sends into the network, using an arrangement called re-inserted feedback or re-feedback. So, every time the sender detects a loss of a 1500B packet, it has to apply a ConEx-marking to a subsequent 1500B packet (or it has to ConEx-mark at least as many bytes made up of smaller packets). In Figure 1 ConEx-marked packets are shown tagged with a '*'.

A bulk congestion policer is deployed at every point where traffic enters an operator's network (similar to the Diffserv architecture). At any one point of entry, the operator assigns a congestion token bucket to each tenant that is sending traffic into the network. The operator assigns a contracted token-fill-rate w_1 to tenant 1, w_2 to tenant 2, and in general w_i to tenant with index i , where $i = 1, 2, \dots$. This contracted token-fill-rate, w_i , can be thought of like a weight, where a higher weight allows a tenant to contribute more traffic through a downstream congested link or to contribute traffic into

more congested links, as we shall see.

Tenants wanting to send more traffic can be assigned a higher weight (probably paying more for it). Typically, tenants might sign-up to a traffic contract in two parts:

1. a peak rate, which represents the capacity of their access to the network
2. a congestion token-fill rate, which limits how much traffic can be sent that limits what others are sending.

In Figure 1, as each packet arrives, the classifier assigns it to the congestion token bucket of the relevant tenant. The token bucket ignores all non-ConEx-marked packets - they do not drain any tokens from the bucket. Otherwise, if a packet is ConEx-marked and ,say, 1500B in size it will drain 1500B of tokens from the congestion token bucket. If the tenant is contributing to more congestion than its contracted congestion-fill-rate, the bucket level will decrease, and if it persists, the bucket level will become low enough to cause the policer to start discarding a proportion of all that tenant's traffic (whether ConEx-marked or not).

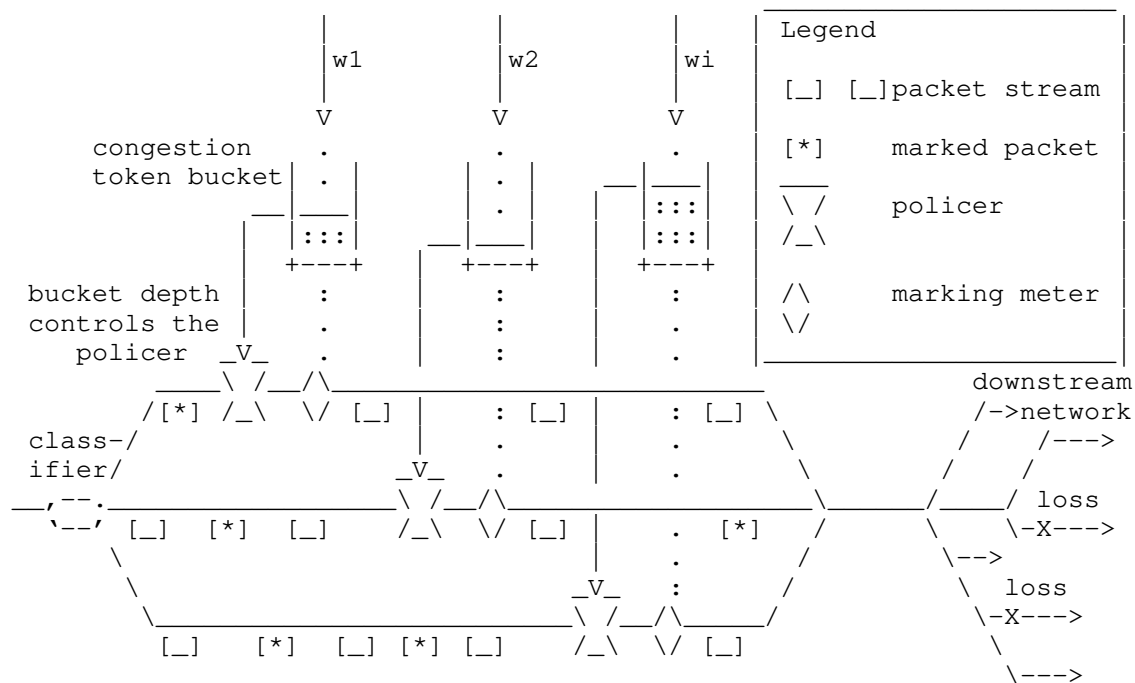


Figure 1: Bulk Congestion Policer Schematic

Note that this is called a bulk congestion policer because it polices all the flows of a tenant even if they spread out in a hose pattern across numerous downstream networks and even if only a few of these flows are primarily to blame for contributing heavily to downstream congestion, while others find plenty of capacity downstream. This approach is motivated by simplicity rather than precision (if precision were required other policer designs could be used). Simplicity might be more important than precision if capacity were generally well-provisioned and the policer was intended to intervene only rarely to limit unusually extreme behaviour. Also, a network might take such a simple but somewhat heavy-handed approach in order to motivate tenants to implement more complex flow-specific controls themselves [CongPol].

Various more sophisticated congestion policer designs have been evaluated [CPolTrilogyExp]. In these experiments, it was found that it is better if the policer gradually increases discards as the bucket becomes empty. Also isolation between tenants is better if each tenant is policed based on the combination of two buckets, not one (Figure 2):

1. A deep bucket (that would take minutes or even hours to fill at the contracted fill-rate) that constrains the tenant's long-term average rate of congestion (w_i)
2. a very shallow bucket (e.g. only two or three MTU) that is filled considerably faster than the deep bucket ($c * w_i$), where $c = \sim 10$, which prevents a tenant storing up a large backlog of tokens then causing congestion in one large burst.

In this arrangement each marked packet drains tokens from both buckets, and the probability of policer discard is taken as the worse of the two buckets.

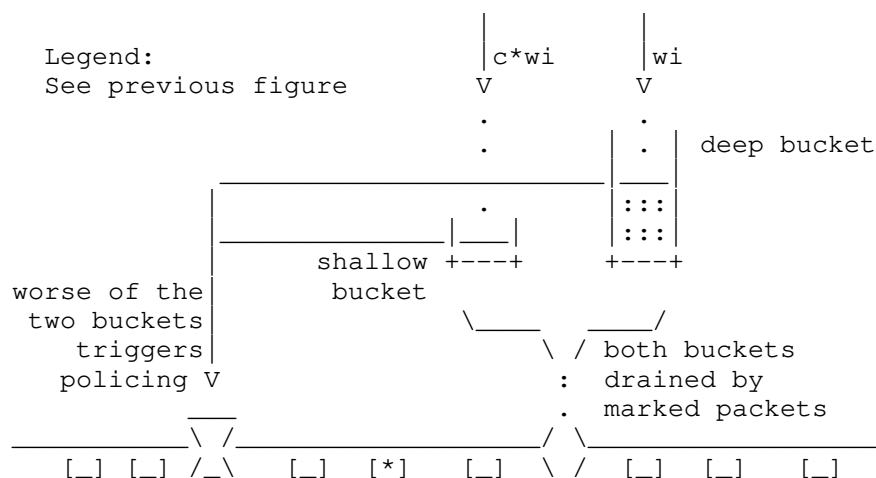


Figure 2: Dual Congestion Token Bucket (in place of each single bucket in the previous figure)

This design of congestion policer will be sufficient to understand the rest of the document. However, it should be remembered that this is only an example, not an ideal design. Also it should be remembered that other mechanisms such as tunnel feedback can be used instead of ConEx.

3. Network Performance Isolation: Intuition

3.1. The Problem

Network performance isolation traditionally meant that each user could be sure of a minimum guaranteed bit-rate. Such assurances are useful if traffic from each tenant follows relatively predictable paths and is fairly constant. If traffic demand is more dynamic and unpredictable (both over time and across paths), minimum bit-rate assurances can still be given, but they have to be very small relative to the available capacity, because a large number of users might all want to simultaneously share any one link, even though they rarely all use it at the same time.

This either means the shared capacity has to be greatly overprovided so that the assured level is large enough, or the assured level has to be small. The former is unnecessarily expensive; the latter doesn't really give a sufficiently useful assurance.

Round robin or fair queuing are other forms of isolation that guarantee that each user will get $1/N$ of the capacity of each link,

where N is the number of active users at each link. This is fine if the number of active users (N) sharing a link is fairly predictable. However, if large numbers of tenants do not typically share any one link but at any time they all could (as in a data centre), a $1/N$ assurance is fairly worthless. Again, given N is typically small but could be very large, either the shared capacity has to be expensively overprovided, or the assured bit-rate has to be worthlessly small. The argument is no different for the weighted forms of these algorithms: WRR & WFQ).

Both these traditional forms of isolation try to give one tenant assured instantaneous bit-rate by constraining the instantaneous bit-rate of everyone else. This approach is flawed except in the special case when the load from every tenant on every link is continuous and fairly constant. The reality is usually very different: sources are on-off and the route taken varies, so that on any one link a source is more often off than on.

In these more realistic (non-constant) scenarios, the capacity available for any one tenant depends much more on how often everyone else uses a link, not just how much bit-rate everyone else would be entitled to if they did use it.

For instance, if 100 tenants are using a 1Gb/s link for 1% of the time, there is a good chance each will get the full 1Gb/s link capacity. But if just six of those tenants suddenly start using the link 50% of the time, whenever the other 94 tenants need the link, they will typically find 3 of these heavier tenants using it already. If a $1/N$ approach like round-robin were used, then the light tenants would suddenly get $1/4 * 1\text{Gb/s} = 250\text{Mb/s}$ on average. Round-robin cannot claim to isolate tenants from each other if they usually get 1Gb/s but sometimes they get 250Mb/s (and only 10Mb/s guaranteed in the worst case when all 100 tenants are active).

In contrast, congestion policing is the key to network performance isolation because it focuses policing only on those tenants that go fast over congested path(s) excessively and persistently over time. This keeps congestion below a design threshold everywhere so that everyone else can go fast. In this way, congestion policing takes account of highly variable loads (varying in time and varying across routes). And, if everyone's load happens to be constant, congestion policing converges on the same outcome as the traditional forms of isolation.

The other flaw in the traditional approaches to isolation like WRR & WFQ is that they actually prevent long-running flows from yielding to brief bursts from lighter tenants. A long-running flow can yield to brief flows and still complete nearly as soon as it would have

otherwise (the brief flows complete sooner, freeing up the capacity for the longer flow sooner--see Section 3.7). However, WRR & WFQ prevent flows from even seeing the congestion signals that would allow them to co-ordinate between themselves, because they isolate each tenant completely into separate queues.

In summary, superficially, traditional approaches with separate queues sound good for isolation, but:

1. not when everyone's load is variable, so each tenant has no assurance about how many other queues there will be;
2. and not when each tenant can no longer even see the congestion signals from other tenants, so no-one's control algorithms can determine whether they would benefit most by pushing harder or yielding.

3.2. Approach

It has not been easy to find a way to give the intuition on why congestion policing isolates performance, particularly across a networks of links not just on a single link. The approach used in this section, is to describe the system as if everyone is using the congestion response they would be forced to use if congestion policing had to intervene. We therefore call this the boundary model of congestion control. It is a very simple congestion response, so it is much easier to understand than if we introduced all the square root terms and other complexity of New Reno TCP's response. And it means we don't have to try to describe a mix of responses.

The purpose of congestion policing is not to intervene in everyone's rate control all the time. Rather it is encourage each tenant to avoid being policed -- to keep the aggregate of all their flows' rates within an overall envelope that is sufficiently responsive to congestion. Nonetheless, congestion policing can and will enforce a congestion response if a particular tenant sends traffic that is completely unresponsive to congestion. This upper bound enforced by the congestion policer is what ensures that each tenant's minimum performance is isolated from the combined effect of everyone else.

We cannot emphasise enough that the intention is not to make individual flows conform to this boundary response to congestion. Indeed the intention is to allow a diverse evolving mix of congestion responses, but constrained in total within a simple overall envelope.

After describing and further justifying using the simple boundary model of congestion control, we start by considering long-running flows sharing one link. Then we will consider on-off traffic, before

widening the scope from one link to a network of links and to links of different sizes. Then we will depart from the initial simplified model of congestion control and consider diverse congestion control algorithms, including completely unresponsive end-systems.

Formal analysis to back-up the intuition provided by this section will be made available in a more extensive companion technical report [conex-dc_tr].

3.3. Terminology

The term congestion probability will be used as a generisation of either loss probability or ECN marking probability.

It is necessary to carefully distinguish:

- o congestion bit-rate (or just congestion rate), which is an absolute measure of the rate of congested bits
- o congestion probability, which is a relative measure of the proportion of congested bits to all bits

Sometimes people loosely talk of loss rate when they mean loss probability (measured in %). In this document, we will keep to strict terminology, so loss rate would be measured in lost packets per second, or lost bits per second.

3.4. Simple Boundary Model of Congestion Control

The boundary model of congestion control ensures a flow's bit-rate is inversely proportional to the congestion level that it detects. For instance, if congestion probability doubles, the flow's bit-rate halves. This is called a scalable congestion control because it maintains the same rate of congestion signals (marked or dropped packets) no matter how fast it goes. Examples from the research community are Relentless TCP [Relentless] and Scalable TCP [STCP].

In production systems, TCP algorithms based on New Reno [RFC5681] have been widely replaced by alternatives closer to this scalable ideal (e.g. Cubic TCP in Linux and Compound TCP in Windows), because at high rates New Reno generated congestion signals too infrequently to track available capacity fast enough [RFC3649]. More recent TCP updates (e.g. Data Centre TCP [DCTCP] in Windows 8 and available in Linux) are becoming closer still to the scalable ideal.

For instance, consider a scenario where a flow with scalable congestion control is alone in a 1Gb/s link, then another similar flow from another tenant joins it. Both will push up the congestion

probability, which will push down their rates until they together fit into the link. Because the flow's rate has to halve to accommodate the new flow, congestion probability will double (lets say from 0.002% to 0.004%), by our initial assumption of a two similar scalable congestion controls. When it is alone on the link, the congestion-bit-rate of the flow is 20kb/s ($=1\text{Gb/s} * 0.002\%$), and when it shares the link it is still 20kb/s ($= 500\text{Mb/s} * 0.004\%$).

In summary, a congestion control can be considered scalable if the bit-rate of packets carrying congestion signals (the congestion-bit-rate) always stays the same no matter how much capacity it finds available. This ensures there will always be enough signals in a round trip time to keep the dynamics under control.

Reminder: Making individual flows conform to this boundary or scalable response to congestion is a non-goal. Although we start this explanation with this specific simple end-system congestion response, this is just to aid intuition.

3.5. Long-Running Flows

Table 1 shows various scenarios where each of five tenants has contracted for 40kb/s of congestion-bit-rate in order to share a 1Gb/s link. In order to help intuition, we start with the (unlikely) scenario where all their flows are long-running. A number of long-running flows will try to use all the link capacity, so for simplicity we take utilisation as a round figure of 100%.

In the case we have just described (scenario A) neither tenant's policer is intervening at all, because both their congestion allowances are 40kb/s and each sends only one flow that contributes 20kb/s of congestion -- half the allowance.

Tenant	contracted congestion- bit-rate kb/s	scenario A # : Mb/s	scenario B # : Mb/s	scenario C # : Mb/s	scenario D # : Mb/s
(a)	40	1 : 500	5 : 250	5 : 200	5 : 250
(b)	40	1 : 500	3 : 250	3 : 200	2 : 250
(c)	40	- : ---	3 : 250	3 : 200	2 : 250
(d)	40	- : ---	2 : 250	2 : 200	1 : 125
(e)	40	- : ---	- : ---	2 : 200	1 : 125
	Congestion probability	0.004%	0.016%	0.02%	0.016%

Table 1: Bit-rates that a congestion policer allocates to five tenants sharing a 1Gb/s link with various numbers (#) of long-running flows all using 'scalable congestion control'

Scenario B shows a case where four of the tenants all send 2 or more long-running flows. Recall, from the assumption of scalable controls, that each flow always contributes 20kb/s no matter how fast it goes. Therefore the policers of tenants (a-c) limit them to two flows-worth of congestion ($2 \times 20\text{kb/s} = 40\text{kb/s}$). Tenant (d) is only asking for 2 flows, so it gets them without being policed, and all four get the same quarter share of the link.

Scenario C is similar, except the fifth tenant (e) joins in, so they all get equal $1/5$ shares of the link.

In Scenario D, only tenant (a) sends more than two flows, so (a)'s policer limits it to two flows-worth of congestion, and everyone else gets the number of flows-worth that they ask for. This means that tenants (d&e) get less than everyone else, which is fine because they asked for less than they would have been allowed. (Similarly, in Scenarios A & B, some of the tenants are inactive, so they get zero, which is also less than they could have had if they had wanted.)

With lots of long-running flows, as in scenarios B & C, congestion policing seems to emulate per-tenant round robin scheduling, equalising the bit-rate of each tenant, no matter how many flows they run. By configuring different contracted allowances for each tenant, it can easily be seen that congestion policing could emulate weighted round robin (WRR), with the relative sizes of the allowances acting as the weights.

Scenario D departs from round-robin. This is deliberate, the idea being that tenants are free to take less than their share in the short term, which allows them to take more at other times, as we will see in Section 3.7. In Scenario D, policing focuses only on the tenant (a) that is continually exceeding its contract. This policer focuses discard solely on tenant a's traffic so that it cannot cause any more congestion at the shared link (shown as 0.016% in the last row).

To summarise so far, ingress congestion policers control congestion-bit-rate in order to indirectly assure a minimum bit-rate per tenant. With lots of long-running flows, the outcome is somewhat similar to WRR, but without the need for any mechanism in each queue.

However, aiming to behave like WRR is only useful when all flows are infinitely long. When load is variable because transfers are of finite size, congestion policing properly isolates one tenant's

performance from the behaviour of others -- unlike WRR would, as will now be explained.

3.6. On-Off Flows

Figure 3 compares two example scenarios where tenant 'b' regularly sends small files in the top chart and the same size files but more often in the bottom chart (a higher 'on-off ratio'). This is the typical behaviour of a Web server when more clients request more files at peak time. Meanwhile, in this example, tenant c's behaviour doesn't change between the two scenarios -- it sends a couple of large files, each starting at the same time in both cases.

The capacity of the link that 'b' and 'c' share is shown as the full height of the plot. The files sent by 'b' are shown as little rectangles. 'b' can go at the full bit-rate of the link when 'c' is not sending, which is represented by the tall thin rectangles part-way through the trace labelled 'b'. We assume for simplicity that 'b' and 'c' divide up the bit-rate equally. So, when both 'b' and 'c' are sending, the 'b' rectangles are half the height (bit-rate) and twice the duration relative to when 'b' sends alone. The area of a file to be transferred stays the same, whether tall and thin or short and fat, because the area represents the size of the file (bit-rate x duration = file size). The files from 'c' look like inverted castellations, because 'c' uses half the link rate while each file from 'b' completes, then 'c' can fill the link until 'b' starts the next file. The cross-hatched areas represent idle times when no-one is sending.

For this simple scenario we ignore start-up dynamics and just focus on the rate and duration of flows that are long enough to stabilise, which is why they can be represented as simple rectangles. We will introduce the effect of flow startups later.

In the bottom case, where 'b' sends more often, the gaps between b's transfers are smaller, so 'c' has less opportunity to use the whole line rate. This squeezes out the time it takes for 'c' to complete its file transfers (recall a file will always have the same area which represents its size). Although 'c' finishes later, it still starts the next flow at the same time. In turn, this means 'c' is sending during a greater proportion of b's transfers, which extends b's average completion time too.

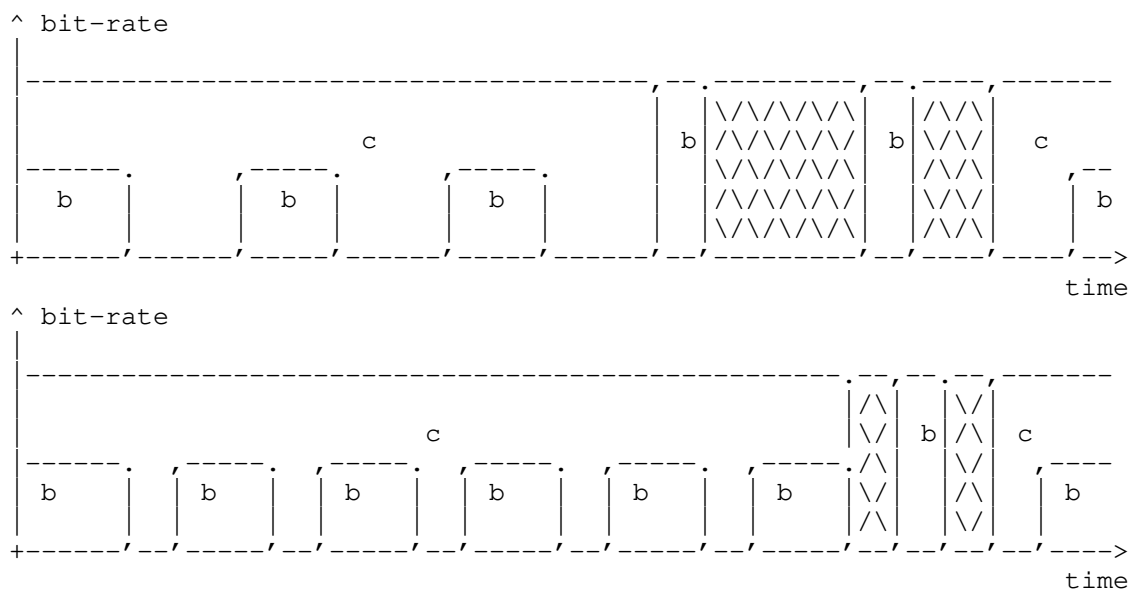


Figure 3: In the lower case, the on-off ratio of 'b' has increased, which extends all the completion times of 'c' and 'b'

Round-robin would do little if anything to isolate 'c' from the effect of 'b' sending files more often. Round-robin is designed to force 'b' and 'c' to share the capacity equally when they are both active. But in both scenarios they already share capacity equally when they are both active. The difference is in how often they are active. Round-robin and other traditional fair queuing techniques don't have any memory to sense that 'b' has been active more of the time.

In contrast, a congestion policer can tell when one tenant is sending files more frequently, by measuring the rate at which the tenant is contributing to congestion. Our aim is to show that policers will be able to isolate performance properly by using the right metric (congestion bit-rate), rather than using the wrong metric (bit-rate), which doesn't sense whether the load over time is large or small.

3.6.1. Numerical Examples Without Policing

The usefulness of the congestion bit-rate metric will now be illustrated with the numerical examples in Table 2. The scenarios illustrate what the congestion bit-rate would be without any policing or scheduling action in the network. Then this metric can be monitored and limited by a policer, to prevent one tenant from harming the performance of others.

The 2nd & 3rd columns (file-size and inter-arrival time) fully represent the behaviour of each tenant in each scenario. All the other columns merely characterise the outcome in various ways. The inter-arrival time (T) is the average time between starting one file and the next. For instance, tenant 'b' sends a 16Mb file every 200ms on average. The formula in the heading of some columns shows how the column was derived from other columns.

Scenario E is contrived so that the three tenants all offer the same load to the network, even though they send files of very different size (S). The files sent by tenant 'a' are 100 times smaller than those of tenant 'b', but 'a' sends them 100 times more often. In turn, b's files are 100 times smaller than c's, but 'b' in turn sends them 100 times more often. Graphically, the scenario would look similar to Figure 3, except with three sizes of file, not just two. Scenarios E-G are designed to roughly represent various distributions of file sizes found in data centres, but still to be simple enough to facilitate intuition, even though each tenant would not normally send just one size file.

The average completion time (t) and the maximum were calculated from a fairly simple analytical model (documented in a companion technical report [conex-dc_tr]). Using one data point as an example, it can be seen that a 1600Mb (200MB) file from tenant 'c' completes in 1905ms (about 1.9s). The files that are 100 times smaller complete 100 times more quickly on average. In fact, in this scenario with equal loads, each tenant perceives that their files are being transferred at the same rate of 840Mb/s on average (file-size divided by completion time, as shown in the apparent bit-rate column). Thus on average all three tenants perceive they are getting 84% of the 1Gb/s link on average (due to the benefit of multiplexing and utilisation being low at 240Mb/s / 1Gb/s = 24% in this case).

The completion times of the smaller files vary significantly, depending on whether a larger file transfer is proceeding at the same time. We have already seen this effect in Figure 3, where, when tenant b's files share with 'c', they take twice as long to complete as when they don't. This is why the maximum completion time is greater than the average for the small files, whereas there is imperceptible variance for the largest files.

The final column shows how congestion bit-rate will be a useful metric to enforce performance isolation (as we said, the figures illustrate the situation before any enforcement mechanism is added). In the case of equal loads (scenario E), average congestion bit-rates are all equal. In scenarios F and G average congestion bit-rates are higher, because all tenants are placing much more load on the network over time, even though each still sends at equal rates to others when

they are active together. Figure 3 illustrated a similar effect in the difference between the top and bottom scenarios.

The maximum instantaneous congestion bit-rate is nearly always 20kb/s. That is because, by definition, all the tenants are using scalable congestion controls with a constant congestion rate of 20kb/s, and each tenant's flows rarely overlap because the per-tenant load in these scenarios is fairly low. As we saw in Section 3.4, the congestion rate of a particular scalable congestion control is always the same, no matter how many other flows it competes with.

Once it is understood that the congestion bit-rate of one scalable flow is always 'w' and doesn't change whenever a flow is active, it becomes clear what the congestion bit-rate will be when averaged over time; it will simply be 'w' multiplied by the proportion of time that the tenant's file transfers are active. That is, $w \cdot t / T$. For instance, in scenario E, on average tenant b's flows start 200ms apart, but they complete in 19ms. So they are active for $19 / 200 = 10\%$ of the time (rounded). A tenant that causes a congestion bit-rate of 20kb/s for 10% of the time will have an average congestion-bit-rate of 2kb/s, as shown.

To summarise so far, no matter how many more files other tenants transfer at the same time, each scalable flow still contributes to congestion at the same rate, but it contributes for more of the time, because it squeezes out into the gap before the next flow from that tenant starts.

Tenant	File size	Ave. inter-arrival	Ave. load	Completion time	Apparent bit-rate	Congestion bit-rate
	S	T	S/T	ave : max t	ave : min S/t	ave : max w*t/T
	Mb	ms	Mb/s	ms	Mb/s	kb/s
Scenario E						
a	0.16	2	80	0.19 : 0.48	840 : 333	2 : 20
b	16	200	80	19 : 35	840 : 460	2 : 20
c	1600	20000	80	1905 : 1905	840 : 840	2 : 20
			240			
Scenario F						
a	0.16	0.67	240	0.31 : 0.48	516 : 333	9 : 20
b	16	50	320	29 : 42	557 : 380	11 : 20
c	1600	10000	160	3636 : 3636	440 : 440	7 : 20
			720			
Scenario G						
a	0.16	0.67	240	0.33 : 0.64	481 : 250	10 : 20
b	16	40	400	32 : 46	505 : 345	16 : 40
c	1600	10000	160	4543 : 4543	352 : 352	9 : 20
			800			

Single link of capacity 1Gb/s. Each tenant uses a scalable congestion control which contributes a congestion-bit-rate for each flow of $w = 20\text{kb/s}$.

Table 2: How the effect on others of various file-transfer behaviours can be measured by the resulting congestion-bit-rate

In scenario F, clients have increased the rate they request files from tenants a, b and c respectively by 3x, 4x and 2x relative to scenario E. The tenants send the same size files but 3x, 4x and 2x more often. For instance tenant 'b' is sending 16Mb files four times as often as before, and they now take longer as well -- nearly 29ms rather than 19ms -- because the other tenants are active more often too, so completion gets squeezed to later. Consequently, tenant 'b' is now sending 57% of the time, so its congestion-bit-rate is $20\text{kb/s} * 57\% = 11\text{kb/s}$. This is nearly 6x higher than in scenario E, reflecting both b's own increase by 4x and that this increase coincides with everyone else increasing their load.

In scenario G, tenant 'b' increases even more, to 5x the load it offered in scenario E. This results in average utilisation of 800Mb/s / 1Gb/s = 80%, compared to 72% in scenario F and only 24% in scenario E. 'b' sends the same files but 5x more often, so its load rises 5x.

Completion times rise for everyone due to the overall rise in load, but the congestion rates of 'a' and 'c' don't rise anything like as much as that of 'b', because they still leave large gaps between files. For instance, tenant 'c' completes each large file transfer in 4.5s (compared to 1.9s in scenario E), but it still only sends files every 10s. So 'c' only sends 45% of the time, which is reflected in its congestion bit-rate of 20kb/s * 45% = 9kb/s.

In contrast, on average tenant 'b' can only complete each medium-sized file transfer in 32ms (compared to 19ms in scenario E), but on average it starts sending another file after 40ms. So 'b' sends 79% of the time, which is reflected in its congestion bit-rate of 20kb/s * 79% = 16kb/s (rounded).

However, during the 45% of the time that 'c' sends a large file, b's completion time is higher than average (as shown in Figure 3). In fact, as shown in the maximum completion time column, 'b' completes in 46ms, but it starts sending a new file after 40ms, which is before the previous one has completed. Therefore, during each of c's large files, 'b' sends 46/40 = 116% of the time on average.

This actually means 'b' is overlapping two files for 16% of the time on average and sending one file for the remaining 84%. Whenever two file transfers overlap, 'b' will be causing 2 x 20kb/s = 40kb/s of congestion, which explains why tenant b in scenario G is the only case with a maximum congestion rate of 40kb/s rather than 20kb/s as in every other case. Over the duration of c's large files, 'b' would therefore cause congestion at an average rate of 20kb/s * 84% + 40kb/s * 16% = 23kb/s (or more simply 20kb/s * 116% = 23kb/s). Of course, when 'c' is not sending a large file, 'b' will contribute less to congestion, which is why its average congestion rate is 16kb/s overall, as discussed earlier.

3.6.2. Congestion Policing of On-Off Flows

Still referring to the numerical examples in Table 2, we will now discuss the effect of limiting each tenant with a congestion policer.

The network operator might have deployed congestion policers to cap each tenant's average congestion rate to 16kb/s. None of the tenants are exceeding this limit in any of the scenarios, but tenant 'b' is just shy of it in scenario G. Therefore all the tenants would be free to behave in all sorts of ways like those of scenarios E-G, but they

would be prevented from degrading the performance of the other tenants beyond the point reached by tenant 'b' in scenario G. If tenant 'b' added more load, the policer would prevent the extra load entering the network by focusing drop solely on tenant 'b', preventing the other tenants from experiencing any more congestion due to tenant 'b'. Then tenants 'a' and 'c' would be assured the (average) apparent bit-rates shown, whatever the behaviour of 'b'.

If 'a' added more load, 'c' would not suffer. Instead 'b' would go over limit and its rate would be trimmed during congestion peaks, sacrificing some of its lead to 'a'. Similarly, if 'c' added more load, 'b' would be made to sacrifice some of its performance, so that 'a' would not suffer. Further, if more tenants arrived to share the same link, the policer would force 'b' to sacrifice performance in favour of the additional tenants.

There is nothing special about a policer limit of 16kb/s. The example when discussing infinite flows used a limit of 40kb/s per tenant. And some tenants can be given higher limits than others (e.g. at an additional charge). If the operator gives out congestion limits that together add up to a higher amount but it doesn't increase the link capacity, it merely allows the tenants to apply more load (e.g. more files of the same size in the same time), but each with lower bit-rate.

3.7. Weighted Congestion Controls

At high speed, congestion controls such as Cubic TCP, Data Centre TCP, Compound TCP etc all contribute to congestion at widely differing rates, which is called their 'aggressiveness' or 'weight'. So far, we have made the simplifying assumption of a scalable congestion control algorithm that contributes to congestion at a constant rate of $w = 20\text{kb/s}$. We now assume tenant 'c' uses a similar congestion control to before, but with different parameters in the algorithm so that its weight is still constant, but much lower, at $w = 2.2\text{kb/s}$.

Tenant 'b' still uses $w = 20\text{kb/s}$ for its smaller files, so when the two compete for the 1Gb/s link, they will share it in proportion to their weights, 20:2.2 (or 90%:10%). That is, when competing, 'b' and 'c' will respectively get $(20/22.2)*1\text{Gb/s} = 900\text{Mb/s}$ and $(2.2/22.2)*1\text{Gb/s} = 100\text{Mb/s}$ of the 1Gb/s link. Figure 4 shows the situation before (upper) and after (lower) this change.

When the two compete, 'b' transfers each file 9/5 faster than before (900Mb/s rather than 500Mb/s), so it completes them in 5/9 of the time. 'b' still contributes congestion at the same rate of 20kb/s, but for 5/9 less time than before. Therefore, relative to before,

'b' uses up its allowance 5/9 as quickly.

Even though 'b' completes each small file much faster, 'c' should complete its file transfer hardly any later than before. Although tenant 'b' goes faster, as each 'b' file finishes, it gets out of the way sooner. Irrespective of its lower weight, 'c' can still use the full link when 'b' is not present, because weights only determine shares between active flows. Because 'c' has the whole link sooner it can catch up to the same point it would have reached in its download if 'b' had been slower but longer. Tenant 'c' will probably lose some completion time because it has to accelerate and decelerate more. But, whenever it is sending a file, 'c' gains $(20\text{kb/s} - 2\text{kb/s}) = 18\text{kb}$ of allowance every second, which it can use for other transfers.

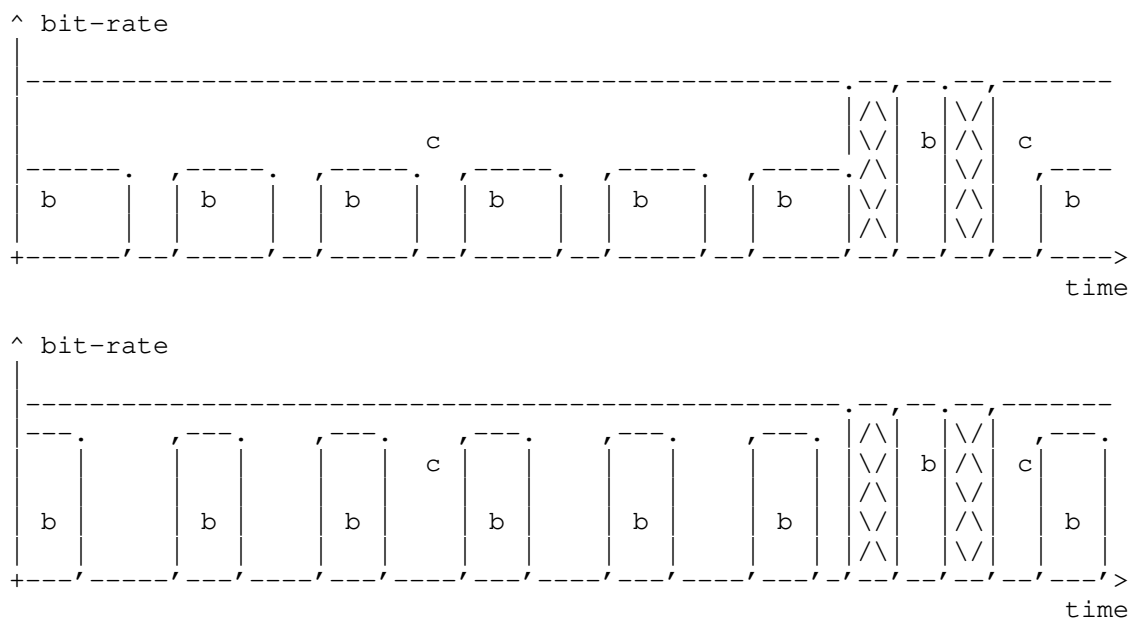


Figure 4: Weighted congestion controls with equal weights (upper) and unequal (lower)

It seems too good to be true that both tenants gain so much and lose so little by 'c' reducing its aggressiveness. The gains are unlikely to be as perfect as this simple model predicts, but we believe they will be nearly as substantial.

It might seem that everyone can keep gaining by everyone agreeing to reduce their weights, ad infinitum. However, the lower the weight, the less signals the congestion control gets, so it starts to lose

its control during dynamics. Nonetheless, congestion policing should encourage congestion control designs to keep reducing their weights, but they will have to stop when they reach the minimum necessary congestion in order to maintain sufficient control signals.

3.8. A Network of Links

So far we have only considered a single link. Congestion policing at the network edge is designed to work across a network of links, treating them all as a pool of resources, as we shall now explain. We will use the dual-homed data centre topology shown in Figure 5 (stretching the bounds of ASCII art) as a simple example of a pool of resources.

In this case where there are 48 servers (H_1, H_2, \dots, H_n where $n=48$) on the left, with on average 8 virtual machines (VMs) running on each (e.g. server n is running V_{n1}, V_{n2}, \dots to V_{nm} where $m = 8$). Each server is connected by two 1Gb/s links, one to each top-of-rack switch S_1 & S_2 . To the right of the switches, there are 6 links of 10Gb/s each, connecting onwards to customer networks or to the rest of the data centre. There is a total of $48 * 2 * 1\text{Gb/s} = 96\text{Gb/s}$ capacity between the 48 servers and the 2 switches, but there is only $6 * 10\text{Gb/s} = 60\text{Gb/s}$ to the right of the switches. Nonetheless, data centres are often designed with some level of contention like this, because at the ToR switches a proportion of the traffic from certain hosts turns round locally towards other hosts in the same rack.

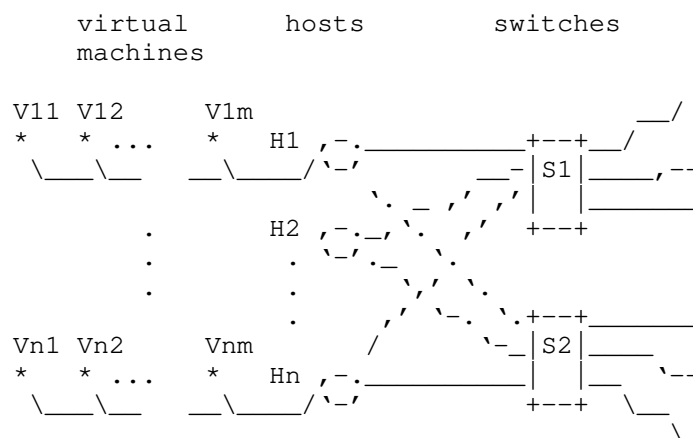


Figure 5: Dual-Homed Topology -- a Simple Resource Pool

The congestion policer proposed in this document is based on the 'hose' model, where a tenant's congestion allowance can be used for sending data over any path, including many paths at once. Therefore,

any one of the virtual machines on the left can use its allowance to contribute to congestion on any or all of the 6 links on the right (or any other link in the diagram actually, including those from the server to the switches and those turning back to other hosts).

Nonetheless, if congestion policers are to enforce performance isolation, they should stop one tenant squeezing the capacity available to another tenant who needs to use a particular bottleneck link or links (e.g. to reach a particular receiver). Policing should work whether the offending tenant is acting deliberately or merely carelessly.

3.8.1. Numerical Example: Isolation from Focused Load

Consider the following scenario: two tenants have equal congestion allowances of 40kb/s, and let's imagine they each run servers that send about 1 flow per second. All their flows use the same scalable congestion control and all use the same aggressiveness or weight of 20kb/s.

Tenant A has clients all over the network, so it spreads its traffic over numerous links, and its flows usually complete within 1 second. Therefore, on average, A has one flow running at any one instant, so it consumes 20kb/s from its 40kb/s congestion allowance.

Then tenant B decides (carelessly or deliberately) to concentrate all its flows into one of the links that A sometimes uses for one of its flows. All the flows through this new bottleneck (A's, B's and those from other tenants) will still each consume 20kb/s of allowance while they are active (Section 3.6), but they will all take longer to complete. Let's say they take 15 times longer to complete on average (i.e. 15s not 1s). And let's assume 5% of tenant A's flows pass through this link. Then each tenant will consume congestion allowance at:

- o Tenant A: $(95\% * 1 * 20\text{kb/s}) + (5\% * 15 * 20\text{kb/s}) = 34\text{kb/s}$

- o Tenant B: $(100\% * 15 * 20\text{kb/s}) = 300\text{kb/s}$

Tenant B's congestion token bucket will therefore drain 7.5 times faster than it is filling ($300 / 40 = 7.5$) so it will rapidly empty and start dropping some of tenant B's traffic before it enters the network. In fact, once tenant B's bucket is empty, it will be limited to no more than 40kb/s of congestion at the bottleneck link, rather than 300kb/s otherwise. In effect, the policer shifts packet drops in excess of the 40kb/s allowance out of the shared link and into itself, to focus excess discards only on tenant B.

Once tenant B's bucket empties, tenant A's completion time for the flows through this bottleneck would not return to 1s, but it would be limited to about 2s (not 15s), and it would consume its congestion allowance at about 21kb/s, well within its limit of 40kb/s.

3.8.2. Isolation in the Short-Term

It seems problematic that tenant B's congestion policer takes time to empty before it kicks in, during which time other tenants suffer reduced performance. This is deliberate--to allow some give and take over time. Nonetheless, it is possible to limit this problem as well, using the dual token bucket already described in Figure 5. The second shallow bucket sets an immediate limit on how much tenant B can harm the performance of others. The fill-rate of the shallow bucket is c times that of the deeper bucket (let's say $c = 8$), so tenant B is immediately limited to $8 * 40\text{kb/s} = 320\text{kb/s}$ of congestion. In the scenario described above, tenant B only reaches 300kb/s but at least this backstop prevents congestion from tenant B ever exceeding 320kb/s.

3.8.3. Encouraging Load Balancing

A policer may not even have to directly intervene for tenants to be protected; it might encourage load balancing to remove the problem first. Load balancing might either be provided by the network (usually just random), or some of the tenants might themselves actively shift traffic off the increasingly congested bottleneck and onto other paths. Some of them might be using the multipath TCP protocol (MPTCP -- see experimental [RFC6356]) that would achieve this automatically, or ultimately if the congestion signals persisted, automatic VM placement algorithms might shift their virtual machine to a different endpoint to circumvent the congestion hot-spot completely. Even if some tenants were not using MPTCP or could not shift easily, others shifting away would achieve the same outcome. Essentially, the deterrent effect of congestion policers encourages everyone to even out congestion across the network, shifting load away from hot spots. Then performance isolation becomes an emergent property of everyone's behaviour, due to the deterrent effect of policers, rather than always through explicit policer intervention.

In such an environment, a policer is needed that shares out the whole pool of network resources, not one that just controls shares of each link individually.

In contrast, enforcement mechanisms based on scheduling algorithms like WRR or WFQ have to be deployed at each link, and each one works in ignorance of how much of other links the tenant is using. These

algorithms were designed for networks with a single known bottleneck per customer (e.g. many access networks). However, in data centres there are many potential bottlenecks and each tenant generally only uses a share of a small number of them. A mechanism like WRR would not isolate anyone's performance if it gave every tenant the right to use the same share of all the links in the network, without regard to how many they were using.

3.9. Tenants of Different Sizes

The scenario in Section 3.8 assumed tenants A & B had equal congestion allowances. If one tenant bought a huge allowance (e.g. 500kb/s) while another tenant bought a small allowance (e.g. 10kb/s), it would seem that the former could focus all its traffic on one link to swamp the latter.

At this point, it needs to be pointed out that congestion allowances do not preclude the need for decent capacity planning. If we further assume that tenant A has bought the same access capacity as tenant B, 50 times more congestion allowance implies tenant A expects to be provided with 50 times the contention ratio of tenant B, i.e. any single links they share will be big enough for B to still have enough if it gets 1 share for every 50 used by A.

However, there is still a problem in that A can store up congestion allowance and cause much more than 500kb/s of congestion in one burst. The dual token bucket arrangement (Figure 5) limits tenant A to $c * 500\text{kb/s}$ of congestion, but if $c=8$ as before, $8 * 500\text{kb/s} = 4\text{Mb/s}$ of congestion in a burst would seriously swamp tenant B if all focused on one link where tenant B was quietly consuming its 10kb/s allowance.

The answer to this problem is that the burst limit c can be smaller for larger tenants, tending to 1 for the largest tenants (i.e. a single shallow bucket would suffice for huge tenants). The reasoning is that the c factor allows a tenant some give-and-take over time, but the aggregate of traffic from a larger tenant consists of enough flows that it has sufficient give-and-take within its own aggregate without needing give-and-take from others.

3.10. Links of Different Sizes

Congestion policing treats a Mb/s of capacity in one link as identical to a Mb/s of capacity in another link, even if the size of each link is different. For instance, consider the case where one of the three links to the right of each switch in Figure 5 were upgraded to 40Gb/s while the other two remained at 10Gb/s (perhaps to accommodate the extra traffic from a couple of the dual homed 1Gb/s

servers being upgraded to dual-homed 10Gb/s).

A congestion control algorithm running at the same rate will cause the same level of congestion probability, whatever size link it is sharing. For example:

- o If 50 equal flows share a 10Gb/s link ($10\text{Gb/s} / 50 = 200\text{Mb/s}$ each) they will cause 0.01% congestion probability;
- o If 200 of the same flows share a 40Gb/s link ($40\text{Gb/s} / 200 = 200\text{Mb/s}$ each) they will still cause 0.01% congestion probability;

This is because the congestion probability is determined by the congestion control algorithms, not by the link.

Therefore, if an average of 300 flows were spread across the above links (1x 40Gb/s and 2 x 10Gb/s), the numbers on each link would tend towards respectively 200:50:50, so that each flow would get 200Mb/s and each link would have 0.01% congestion on it. Sometimes, there might be more flows on one link, resulting in less than 200Mb/s per flow and congestion higher than 0.01%. However, whenever the congestion level was greater on one link than another, congestion policing would encourage flows to balance out the congestion level across the links (as long as some flows could use congestion balancing mechanisms like MPTCP).

In summary, all the outcomes of congestion policing described so far (emulating WRR, etc) apply across a pool of diverse link sizes just as much as they apply to single links.

3.11. Diverse Congestion Control Algorithms

Throughout this explanation we have assumed a scalable congestion control algorithm, which we justified (Section 3.4) as the 'boundary' case if congestion policing had to intervene, which is all that is relevant when considering whether the policer can enforce performance isolation.

The defining difference between the scalable congestion we have assumed and the congestion controls in widespread production operating systems (New Reno, Compound, Cubic, Data Centre TCP etc) is the way congestion probability decreases as flow-rate increases (for a long-running flow). With a scalable congestion control, if flow-rate doubles, congestion probability halves. Whereas, with most production congestion controls, if flow-rate doubles, congestion probability reduces to less than half. For instance, New Reno TCP reduces congestion to a quarter. The responses of Cubic and Compound are closer to the ideal scalable control than to New Reno, but they

do not depart too far from New Reno to ensure they can co-exist happily with it.

Congestion policing still works, whether or not the congestion controls in daily use by tenants fit the scalable model. A bulk congestion policer constrains the sum of all the congestion controls being used by a tenant so that they collectively remain below an aggregate envelope that is itself shaped like the sum of many scalable algorithms. Bulk congestion policers will constrain the overall congestion effect (the sum) of any mix of algorithms within it, including flows that are completely unresponsive to congestion.

This will now be explained using Figure 6 (if the ASCII art is incomprehensible, see a similar plot at Figure 3 of [CongPol]).

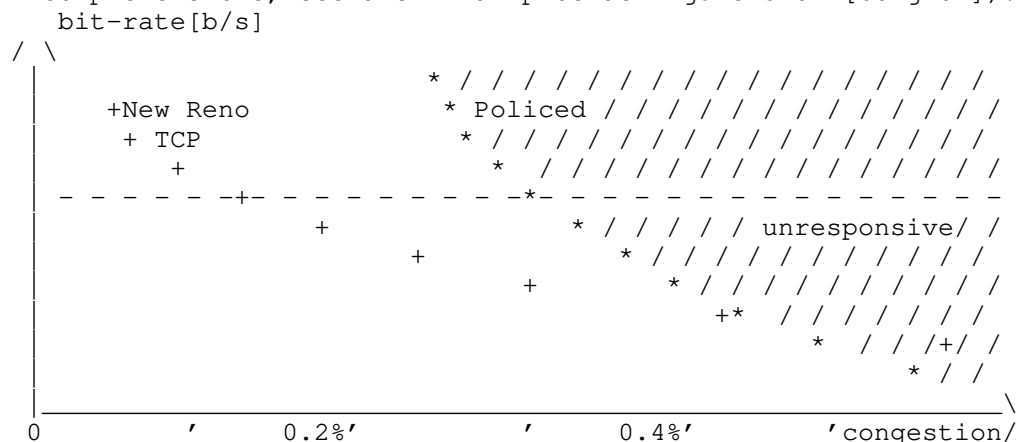


Figure 6: Schematic Plot of Bulk Policing of Traffic with Different Congestion Responses

The congestion policer prevents the aggregate of a tenant's traffic from operating in the hatched region shown in Figure 6 (there is deliberately no vertical scale, because the plot is schematic only). The tenant can have either high congestion or high rate, but not both at the same time. It can be seen that the single New Reno TCP flow responds to congestion in a similar way, but with a flatter slope. It can also be seen that, at 0.25% prevailing congestion, this policer would allow roughly two simultaneous New Reno TCP flows without intervening. At about 0.45% congestion it would allow only one, and above 0.45% it would start to police even one flow. But at much lower congestion levels typical in a high-capacity data centre (e.g. 0.01%) the policer would allow numerous New Reno TCP flows (extrapolating the curve upwards).

The horizontal plot represents an unresponsive (UDP) flow. The

congestion policer would allow one of these flows to run unimpeded up to roughly 0.3% congestion, when the policer would effectively take control and impose its own congestion response on this flow, dropping more UDP packets for higher levels of prevailing congestion. If a tenant ran two of these UDP flows (or one at twice the rate), the policer would leave them alone unless prevailing congestion was above about 0.22% when it would intervene (by extrapolation of the visible plots).

The New Reno TCP curve follows the characteristic k/\sqrt{p} rule, while the 'Policed' curve follows a simple w/p rule, where k is the New Reno constant and w is the constant contracted fill-rate of a particular policer. A scalable TCP (not shown) would follow a v/p rule where v is the weight of the individual flow. If we imagine a scenario where $w = 8v$, the policer would allow 8 of these scalable TCP flows (averaged over time), and this would be true at any congestion level, because the curvature of 8 scalable TCP's would exactly match the 'Policed' curve.

So far, the discussion has focused on the congestion avoidance phase of each congestion response. As each TCP flow starts, it typically increases exponentially (TCP slow-start) until it overshoots available capacity causing a spike of loss due to the feedback delay over the following round-trip. The ConEx audit mechanism requires the source to send credit markings in advance to cover the risk of these spikes of loss. Credit marked ConEx packets drain tokens from the congestion policer similarly to regular ConEx marked packets, which reduces the available congestion allowance for other flows. If the tenant starts new flows fast enough, these credit markings alone will empty the token bucket so that the policer starts to discard some packets from the start-up phase before they can enter the network. (The tunnel feedback method for signalling congestion back to the policer lacks such a credit mechanism, which is an important advantage of the ConEx protocol, given Internet traffic generally consists of a high proportion of short flows.)

A tenant could run a mix of multiple different types of TCP and UDP flows, as long as, when all stacked up, the sum of them all still fitted under the 'Policed' curve at the prevailing level of congestion.

4. Parameter Setting

The primary parameter to set for each tenant is the contracted fill-rate of their congestion token bucket. For a server sending predominantly long-running flows (e.g. a video server or an indexing robot filling search databases), this fill-rate can be determined from the required number of simultaneous TCP flows, and the typical

congestion-bit-rate of one TCP flow.

For a scalable TCP flow, the congestion bit-rate is constant. But currently (2013) TCP Cubic is predominant, and the congestion bit-rate of one Cubic flow is not constant, but depends somewhat on bit-rate and RTT. For example, a Cubic flow with 100Mb/s throughput and 50ms RTT, has a congestion-bit-rate of ~2kbs, whereas at 200Mb/s and 5ms RTT, its congestion-bit-rate is ~6kb/s. Therefore, a reasonable rule of thumb for Cubic's congestion-bit-rate at rates and RTTs of about this order is 5kb/s.

Therefore, if a tenant wanted to serve no more on average than 12 simultaneous TCP Cubic flows, their contracted congestion-rate should be roughly $12 * 5\text{kb/s} = 60\text{kb/s}$.

The contracted fill rate of a congestion policer should not need to change as flow-rates and link capacities scale, assuming the tenant still utilises the higher capacity to the same extent.

A congestion policer based on the dual token bucket in Figure 2 also needs to specify the c factor that determined the maximum congestion-rate allowed, rather than the average. A theoretical approach to setting this factor is being worked on but, in the mean time, a figure of about 10 seems reasonable for a tenant sending only a few simultaneous flows (see Section 3.8.2), while Section 3.9 explains that the c factor should tend to reduce down to 1 for very large tenants.

The dual token bucket makes isolation fairly insensitive to the depth of the deeper bucket, because the shallower bucket constrains the congestion burst rate, which makes the congestion burst size less important. It is believed that the deeper bucket depth could be minutes or even hours of token fill. The depth of the shallower bucket should be just a few MTU--more experimentation is needed to determine how few.

Indeed, all the above recommendations on parameter setting are best considered as initial values to be used in experiments to determine whether other values would be better.

No recommendations can yet be given for parameter settings for a tenant running a large proportion of short flows. This requires further experimentation.

5. Security Considerations

This document is about performance isolation, therefore the whole document concerns security. Where ConEx is used, the specification

of the ConEx Abstract Mechanism [I-D.ietf-conex-abstract-mech] should be consulted for security matters, particularly the design of audit to assure the integrity of ConEx signals.

6. IANA Considerations

This document does not require actions by IANA.

7. Conclusions

{ToDo}

8. Acknowledgments

9. Informative References

- [CPolTrilogyExp] Raiciu, C., Ed., "Progress on resource control", Trilogy EU 7th Framework Project ICT-216372 Deliverable 9, December 2009, <<http://trilogy-project.org/publications/deliverables.html>>.
- [CongPol] Jacquet, A., Briscoe, B., and T. Moncaster, "Policing Freedom to Use the Internet Resource Pool", Proc ACM Workshop on Re-Architecting the Internet (ReArch'08) , December 2008, <<http://bobbbriscoe.net/projects/refb/#polfree>>.
- [DCTCP] Alizadeh, M., Greenberg, A., Maltz, D., Padhye, J., Patel, P., Prabhakar, B., Sengupta, S., and M. Sridharan, "Data Center TCP (DCTCP)", ACM SIGCOMM CCR 40(4)63--74, October 2010, <<http://portal.acm.org/citation.cfm?id=1851192>>.
- [I-D.briscoe-conex-data-centre] Briscoe, B. and M. Sridharan, "Network Performance Isolation in Data Centres using Congestion Policing", draft-briscoe-conex-data-centre-01 (work in progress), February 2013.

- [I-D.briscoe-conex-initial-deploy] Briscoe, B., "Initial Congestion Exposure (ConEx) Deployment Examples", draft-briscoe-conex-initial-deploy-03 (work in progress), July 2012.
- [I-D.ietf-conex-abstract-mech] Mathis, M. and B. Briscoe, "Congestion Exposure (ConEx) Concepts and Abstract Mechanism", draft-ietf-conex-abstract-mech-06 (work in progress), October 2012.
- [I-D.ietf-conex-mobile] Kutscher, D., Mir, F., Winter, R., Krishnan, S., Zhang, Y., and C. Bernardos, "Mobile Communication Congestion Exposure Scenario", draft-ietf-conex-mobile-01 (work in progress), January 2013.
- [RFC3649] Floyd, S., "HighSpeed TCP for Large Congestion Windows", RFC 3649, December 2003.
- [RFC5681] Allman, M., Paxson, V., and E. Blanton, "TCP Congestion Control", RFC 5681, September 2009.
- [RFC6356] Raiciu, C., Handley, M., and D. Wischik, "Coupled Congestion Control for Multipath Transport Protocols", RFC 6356, October 2011.
- [Relentless] Mathis, M., "Relentless Congestion Control", Proc. Int'l Wkshp on Protocols for Future, Large-scale & Diverse Network Transports (PFLDNeT'09) , May 2009, <<http://www.hpcc.jp/pfldnet2009/Program.html>>.
- [STCP] Kelly, T., "Scalable TCP: Improving Performance in Highspeed Wide Area Networks", ACM SIGCOMM Computer Communication Review 32(2),

April 2003, <<http://www-lce.eng.cam.ac.uk/~ctk21/scalable/>>.

[conex-dc_tr]

Briscoe, "Network Performance Isolation in Data Centres by Congestion Exposure to Edge Policers", BT Technical Report TR-DES8-2011-004, November 2011.

Work in progress

Appendix A. Summary of Changes between Drafts

This draft was separated out from the -00 version of [I-D.briscoe-conex-data-centre], taken mostly from what was section 4. Changes from that draft are listed below:

From section 4 of draft-briscoe-conex-data-centre-00:

- * New Introduction
- * Added section 2. "Example Bulk Congestion Policers"
- * Rearranged and clarified the previous introductory text that preceded what is now Section 3.4 into three subsections:
 - 3.1 "The Problem"
 - 3.2 "Approach"
 - 3.3 "Terminology"
- * Corrected numerical examples:
 - + Section 3.4 "Long-RunningFlows": 0.04% to 0.004% and 400kb/s to 40kb/s
 - + Section 3.6.1 "Numerical Examples Without Policing": 10kb/s to 20kb/s
 - + section 3.7 "Weighted Congestion Controls": 22.2 to 20
- * Clarified where necessary, especially the descriptions of the scenarios in 3.7 "Weighted Congestion Controls" and 3.8 "A Network of Links".

- * Section 3.8 "A Network of Links": added a numerical example including congestion burst limiting that was not covered at all previously
- * Added Section 3.9 "Tenants of Different Sizes"
- * Filled in absent text in:
 - + Section 3.11 "Diverse Congestion Controls"
 - + Section 4. "Parameters Settings"
 - + Section 5. "Security Considerations"
- * Minor corrections throughout and updated references.

Author's Address

Bob Briscoe
BT
B54/77, Adastral Park
Martlesham Heath
Ipswich IP5 3RE
UK

Phone: +44 1473 645196
EMail: bob.briscoe@bt.com
URI: <http://bobbriscoe.net/>

