

HOWTO du routage avancé et du contrôle de trafic sous Linux

Bert Hubert

bert.hubert@netherlabs.nl

Thomas Graf (Auteur d'une section)

tgraf%suug.ch

Gregory Maxwell (Auteur d'une section)

Remco van Mook (Auteur d'une section)

remco@virtu.nl

Martijn van Oosterhout (Auteur d'une section)

kleptog@cupid.suninternet.com

Paul B. Schroeder (Auteur d'une section)

paulsch@us.ibm.com

Jasper Spaans (Auteur d'une section)

jasper@spaans.ds9a.nl

Pedro Larroy (Auteur d'une section)

piotr%member.fsf.org

Laurent Foucher

foucher(at)gch.iut-tlse3.fr

Philippe Latu

philippe.latu(at)inetdoc.net

Guillaume Allègre

Guillaume.Allegre(at)imag.fr

Alexey Kuznetsov

kuznet@ms2.inr.ac.ru

Andi Kleen

ak@muc.de

Esteve Camps

marvin@grn.es

Pedro Larroy Tovar

piotr%member.fsf.org

Junk Alins

juanjo@mat.upc.es

Joe Van Andel

Michael Babcock

mbabcock@fibrespeed.net

Christopher Barton

cpbarton%uiuc.edu

Peter Bieringer

pb:bieringer.de

Ard van Breemen

ard%kwaak.net

Ron Brinker

service%emcis.com

?ukasz Bromirski

l.bromirski@mr0vka.eu.org

Lennert Buytenhek

buytenh@gnu.org

Esteve Camps

esteve@hades.udg.es

Ricardo Javier Cardenes

ricardo%conysis.com

Stef Coene

stef.coene@docum.org

Don Cohen

don-lartc%isis.cs3-inc.com

Jonathan Corbet

lwn%lwn.net

Gerry Creager

gerry%cs.tamu.edu

Marco Davids

marco@sara.nl

Jonathan Day

jd9812@my-deja.com

Martin Devera

devik@cdi.cz
Hannes Ebner
he%fli4l.de
Derek Fawcus
dfawcus%cisco.com
David Fries
dfries%mail.win.org
Stephan Gehring
Stephan.Gehring@bechtle.de
Jacek Glinkowski
jglinkow%hns.com
Andrea Glorioso
sama%perchetopi.org
Thomas Graaf
tgraf%suug.ch
Sandy Harris
sandy%storm.ca
Nadeem Hasan
nhasan@usa.net
Erik Hensema
erik%hensema.xs4all.nl
Vik Heyndrickx
vik.heyndrickx@edchq.com
Spauldo Da Hippie
spauldo%usa.net
Koos van den Hout
koos@kzdoos.xs4all.nl
Ayotunde Itayemi
aitayemi:metrong.com
Dave Johnson
dj@www.uk.linux.org
Jose Luis Domingo Lopez
jdomingo@24x7linux.com
Robert Lowe
robert.h.lowe@lawrence.edu
William Stearns
wstearns@pobox.com
Chris Wilson
chris@netservers.co.uk
Lazar Yanackiev
Lyanackiev%gmx.net
Pedro Larroy
piotr%member.fsf.org

<http://www.inetdoc.net>

Une approche pratique d'iproute2, de la mise en forme du trafic et un peu de netfilter.

Table des matières

1. Dédicace	1
2. Introduction	2
2.1. Conditions de distribution et Mise en garde	2
2.2. Connaissances préalables	2
2.3. Ce que Linux peut faire pour vous	2
2.4. Notes diverses	3
2.5. Accès, CVS et propositions de mises à jour	3
2.6. Liste de diffusion	3
2.7. Plan du document	3
3. Introduction à iproute2	4
3.1. Pourquoi iproute2 ?	4
3.2. Un tour d'horizon d'iproute2	4
3.3. Prérequis	4
3.4. Explorer votre configuration courante	4
3.4.1. ip nous montre nos liens	4
3.4.2. ip nous montre nos adresses IP	5
3.4.3. ip nous montre nos routes	5
3.5. ARP	5
4. Règles - bases de données des politiques de routage	7
4.1. Politique de routage simple par l'adresse source	7
4.2. Routage avec plusieurs accès Internet/fournisseurs d'accès	7
4.2.1. Accès séparé	8
4.2.2. Balance de charge	8
5. GRE et autres tunnels	10
5.1. Quelques remarques générales à propos des tunnels :	10
5.2. IP dans un tunnel IP	10
5.3. Le tunnel GRE	10
5.3.1. Le tunnel IPv4	10
5.3.2. Le tunnel IPv6	11
5.4. Tunnels dans l'espace utilisateur	11
6. Tunnel IPv6 avec Cisco et/ou une dorsale IPv6 (6bone)	12
6.1. Tunnel IPv6	12
7. IPSEC: IP sécurisé à travers Internet	14
7.1. Introduction sur la gestion manuelle des clés	14
7.2. Gestion automatique des clés	16
7.2.1. Théorie	17
7.2.2. Exemple	17
7.2.2.1. Problèmes et défauts connus	18
7.2.3. Gestion automatique des clés en utilisant les certificats X.509	19
7.2.3.1. Construire un certificat X.509 pour votre hôte	19
7.2.3.2. Configuration et lancement	19
7.2.3.3. Comment configurer des tunnels sécurisés	20
7.3. tunnels IPSEC	20
7.4. Autre logiciel IPSEC	21
7.5. Interopérabilité d'IPSEC avec d'autres systèmes	21
7.5.1. Windows	21
7.5.2. Check Point VPN-1 NG	21
8. Routage multidistribution (<i>multicast</i>)	22
9. Gestionnaires de mise en file d'attente pour l'administration de la bande passante	23
9.1. Explication sur les files d'attente et la gestion de la mise en file d'attente	23
9.2. Gestionnaires de mise en file d'attente simples, sans classes	23
9.2.1. <i>pfifo_fast</i>	23
9.2.1.1. Paramètres & usage	23
9.2.2. Filtre à seau de jetons (<i>Token Bucket Filter</i>)	25
9.2.2.1. Paramètres & usage	25
9.2.2.2. Configuration simple	26
9.2.3. Mise en file d'attente stochastiquement équitable (<i>Stochastic Fairness Queueing</i>)	26
9.2.3.1. Paramètres & usage	27
9.2.3.2. Configuration simple	27
9.3. Conseils pour le choix de la file d'attente	27
9.4. terminologie	27
9.5. Gestionnaires de file d'attente basés sur les classes	29
9.5.1. Flux à l'intérieur des gestionnaires basés sur des classes & à l'intérieur des classes	29
9.5.2. La famille des gestionnaires de mise en file d'attente : racines, descripteurs, descendances et parents	29
9.5.2.1. Comment les filtres sont utilisés pour classifier le trafic	30
9.5.2.2. Comment les paquets sont retirés de la file d'attente et envoyés vers le matériel	30
9.5.3. Le gestionnaire de mise en file d'attente PRIO	30

9.5.3.1. Paramètres PRIO & usage	31
9.5.3.2. Configuration simple	31
9.5.4. Le célèbre gestionnaire de mise en file d'attente CBQ	32
9.5.4.1. Mise en forme CBQ en détail	32
9.5.4.2. Le comportement <i>CBQ classful</i>	33
9.5.4.3. Paramètres CBQ qui déterminent le partage & le prêt du lien	34
9.5.4.4. Configuration simple	34
9.5.4.5. D'autres paramètres CBQ : split & defmap	35
9.5.5. Seau de jetons à contrôle hiérarchique (<i>Hierarchical Token Bucket</i>)	36
9.5.5.1. Configuration simple	36
9.6. Classifier des paquets avec des filtres	36
9.6.1. Quelques exemples simples de filtrage	37
9.6.2. Toutes les commandes de filtres dont vous aurez normalement besoin	37
9.7. Le périphérique de file d'attente intermédiaire (The Intermediate queueing device (IMQ))	38
9.7.1. Configuration simple	38
10. Équilibrage de charge sur plusieurs interfaces	40
10.1. Avertissement	40
11. Netfilter et iproute - marquage de paquets	41
12. Filtres avancés pour la (re-)classification des paquets	42
12.1. Le classificateur u32	42
12.1.1. Le sélecteur U32	42
12.1.2. Sélecteurs généraux	43
12.1.3. Les sélecteurs spécifiques	44
12.2. Le classificateur route	44
12.3. Les filtres de réglementation (<i>Policing filters</i>)	44
12.3.1. Techniques de réglementation	45
12.3.1.1. Avec l'estimateur du noyau	45
12.3.1.2. Avec le <i>Token Bucket Filter</i>	45
12.3.2. Actions de dépassement de limite (<i>Overlimit actions</i>)	45
12.3.3. Exemples	45
12.4. Filtres hachés pour un filtrage massif très rapide	45
12.5. Filtrer le trafic IPv6	46
12.5.1. Comment se fait-il que ces filtres tc IPv6 ne fonctionnent pas ?	46
12.5.2. Marquer les paquets IPv6 en utilisant ip6tables	47
12.5.3. Utiliser le sélecteur u32 pour repérer le paquet IPv6	47
13. Paramètres réseau du noyau	48
13.1. Filtrage de Chemin Inverse (<i>Reverse Path Filtering</i>)	48
13.2. Configurations obscures	48
13.2.1. ipv4 générique	48
13.2.2. Configuration des périphériques	51
13.2.3. Politique de voisinage	52
13.2.4. Configuration du routage	53
14. Gestionnaires de mise en file d'attente avancés & moins communs	54
14.1. bfifo/pfifo	54
14.1.1. Paramètres & usage	54
14.2. Algorithme Clark-Shenker-Zhang (CSZ)	54
14.3. DSMARK	54
14.3.1. Introduction	54
14.3.2. A quoi DSMARK est-il relié ?	54
14.3.3. Guide des services différenciés	55
14.3.4. Travailler avec DSMARK	55
14.3.5. Comment SCH_DSMARK travaille ?	55
14.3.6. Le filtre TC INDEX	56
14.4. Gestionnaire de mise en file d'attente d'entrée (<i>Ingress qdisc</i>)	57
14.4.1. Paramètres & usage	57
14.5. <i>Random Early Detection</i> (RED)	57
14.6. Generic Random Early Detection	58
14.7. Emulation VC/ATM	58
14.8. Weighted Round Robin (WRR)	58
15. Recettes de cuisine	60
15.1. Faire tourner plusieurs sites avec différentes SLA (autorisations)	60
15.2. Protéger votre machine des inondations SYN	60
15.3. Limiter le débit ICMP pour empêcher les dénis de service	61
15.4. Donner la priorité au trafic interactif	61
15.5. Cache web transparent utilisant netfilter, iproute2, ipchains et squid	62
15.5.1. Schéma du trafic après l'implémentation	64
15.6. Circonvenir aux problèmes de la découverte du MTU de chemin en configurant un MTU par routes	64
15.6.1. Solution	64
15.7. Circonvenir aux problèmes de la découverte du MTU de chemin en imposant le MSS (pour les utilisateurs de l'ADSL, du câble, de PPPoE & PptP)	65

15.8. Le Conditionneur de Trafic Ultime : Faible temps de latence, Téléchargement vers l'amont et l'aval rapide	65
15.8.1. Pourquoi cela ne marche t-il pas bien par défaut ?	66
15.8.2. Le script (CBQ)	67
15.8.3. Le script (HTB)	68
15.9. Limitation du débit pour un hôte ou un masque de sous-réseau	69
15.10. Exemple d'une solution de traduction d'adresse avec de la QoS	69
15.10.1. Commençons l'optimisation de cette rare bande passante	70
15.10.2. Classification des paquets	70
15.10.3. Améliorer notre configuration	71
15.10.4. Rendre tout ceci actif au démarrage	72
16. Construire des ponts et des pseudo ponts avec du Proxy ARP	73
16.1. Etat des ponts et iptables	73
16.2. Pont et mise en forme	73
16.3. Pseudo-pont avec du Proxy-ARP	73
16.3.1. ARP & Proxy-ARP	73
16.3.2. Implémentez-le	73
17. Routage Dynamique - OSPF et BGP	75
17.1. Configurer OSPF avec Zebra	75
17.1.1. Prérequis	76
17.1.2. Configurer Zebra	76
17.1.3. Exécuter Zebra	77
17.2. Configurer BGP4 avec Zebra	78
17.2.1. schéma réseau (Exemple)	78
17.2.2. Configuration (Exemple)	78
17.2.3. Vérification de la configuration	79
18. Autres possibilités	80
19. Lectures supplémentaires	82
20. Remerciements	83

Ce document est dédié à beaucoup de gens ; dans ma tentative de tous me les rappeler, je peux en citer quelques-uns :

- Rusty Russell
- Alexey N. Kuznetsov
- La fine équipe de Google
- L'équipe de Casema Internet

Bienvenue, cher lecteur.

Ce document a pour but de vous éclairer sur la manière de faire du routage avancé avec Linux 2.2/2.4. Méconnus par les utilisateurs, les outils standard de ces noyaux permettent de faire des choses spectaculaires. Les commandes comme **route** et **ifconfig** sont des interfaces vraiment pauvres par rapport à la grande puissance potentielle d'iproute2.

J'espère que ce HOWTO deviendra aussi lisible que ceux de Rusty Russell, très réputé (parmi d'autres choses) pour son netfilter.

Vous pouvez nous contacter en nous écrivant à [l'équipe HOWTO](mailto:HOWTO@ds9a.nl)¹. Cependant, postez, s'il vous plaît, vos questions sur la liste de diffusion (voir la section correspondante) pour celles qui ne sont pas directement liées à ce HOWTO.

Avant de vous perdre dans ce HOWTO, si la seule chose que vous souhaitez faire est de la simple mise en forme de trafic, allez directement au chapitre *Autres possibilités*, et lisez ce qui concerne CBQ.init.

2.1. Conditions de distribution et Mise en garde

Ce document est distribué dans l'espoir qu'il sera utile et utilisé, mais SANS AUCUNE GARANTIE ; sans même une garantie implicite de qualité légale et marchande ni aptitude à un quelconque usage.

En un mot, si votre dorsale STM-64 est tombée ou distribuée de la pornographie à vos estimés clients, cela n'est pas de notre faute. Désolé.

Copyright (c) 2001 par Bert Hubert, Gregory Maxwell et Martijn van Oosterhout, Remco van Mook, Paul B. Schroeder et autres. Ce document ne peut être distribué qu'en respectant les termes et les conditions exposés dans la Open Publication License, v1.0 ou supérieure (la dernière version est actuellement disponible sur <http://www.opencontent.org/openpub/>).

Copiez et distribuez (vendez ou donnez) librement ce document, dans n'importe quel format. Les demandes de corrections et/ou de commentaires sont à adresser à la personne qui maintient ce document.

Il est aussi demandé que, si vous publiez cet HOWTO sur un support papier, vous en envoyiez des exemplaires aux auteurs pour une « relecture critique » :-)

2.2. Connaissances préalables

Comme le titre l'implique, ceci est un HOWTO « avancé ». Bien qu'il ne soit pas besoin d'être un expert réseau, certains pré-requis sont nécessaires.

Voici d'autres références qui pourront vous aider à en apprendre plus :

[Rusty Russell's networking-concepts-HOWTO](#)²

Très bonne introduction, expliquant ce qu'est un réseau, et comment on le connecte à d'autres réseaux.

Linux Networking-HOWTO (ex Net-3 HOWTO)

Excellent document, bien que très bavard. Il vous apprendra beaucoup de choses qui sont déjà configurées si vous êtes capable de vous connecter à Internet. Il peut éventuellement être situé à `/usr/doc/HOWTO/NET-HOWTO.txt`, mais peut également être trouvé [en ligne](#)³

2.3. Ce que Linux peut faire pour vous

Une petite liste des choses qui sont possibles :

- Limiter la bande passante pour certains ordinateurs
- Limiter la bande passante VERS certains ordinateurs
- Vous aider à partager équitablement votre bande passante
- Protéger votre réseau des attaques de type Déni de Service
- Protéger Internet de vos clients
- Multiplexer plusieurs serveurs en un seul, pour l'équilibrage de charge ou une disponibilité améliorée
- Restreindre l'accès à vos ordinateurs
- Limiter l'accès de vos utilisateurs vers d'autres hôtes
- Faire du routage basé sur l'ID utilisateur (eh oui !), l'adresse MAC, l'adresse IP source, le port, le type de service, l'heure ou le contenu.

Peu de personnes utilisent couramment ces fonctionnalités avancées. Il y a plusieurs raisons à cela. Bien que la documentation soit fournie, la prise en main est difficile. Les commandes de contrôle du trafic ne sont pratiquement pas documentées.

¹ <mailto:HOWTO@ds9a.nl>

² <http://netfilter.org/documentation/HOWTO/networking-concepts-HOWTO.html>

³ <http://www.linuxports.com/howto/networking>

2.4. Notes diverses

Il y a plusieurs choses qui doivent être notées au sujet de ce document. Bien que j'en ai écrit la majeure partie, je ne veux vraiment pas qu'il reste tel quel. Je crois beaucoup à l'Open Source, je vous encourage donc à envoyer des remarques, des mises à jour, des corrections, etc. N'hésitez pas à m'avertir des coquilles ou d'erreurs pures et simples. Si mon anglais vous paraît parfois peu naturel, ayez en tête, s'il vous plaît, que l'anglais n'est pas ma langue natale. N'hésitez pas à m'envoyer vos suggestions [NdT : en anglais !]

Si vous pensez que vous êtes plus qualifié que moi pour maintenir une section ou si vous pensez que vous pouvez écrire et maintenir de nouvelles sections, vous êtes le bienvenu. La version SGML de ce HOWTO est disponible via CVS. J'envisage que d'autres personnes puissent travailler dessus.

Pour vous aider, vous trouverez beaucoup de mentions FIXME (NdT : A CORRIGER). Les corrections sont toujours les bienvenues. Si vous trouvez une mention FIXME, vous saurez que vous êtes en territoire inconnu. Cela ne veut pas dire qu'il n'y a pas d'erreurs ailleurs, faites donc très attention. Si vous avez validé quelque chose, faites-nous le savoir, ce qui nous permettra de retirer la mention FIXME.

Je prendrai quelques libertés tout au long de cet HOWTO. Par exemple, je pars de l'hypothèse d'une connexion Internet à 10 Mbits, bien que je sache très bien que cela ne soit pas vraiment courant.

2.5. Accès, CVS et propositions de mises à jour

L'adresse canonique de cet HOWTO est [Ici](#)⁴.

Nous avons maintenant un CVS en accès anonyme disponible depuis le monde entier. Cela est intéressant pour plusieurs raisons. Vous pouvez facilement télécharger les nouvelles versions de ce HOWTO et soumettre des mises à jour.

En outre, cela permet aux auteurs de travailler sur la source de façon indépendante, ce qui est une bonne chose aussi.

```
$ export CVSROOT=:pserver:anon@outpost.ds9a.nl:/var/cvsroot
$ cvs login
CVS password: [enter 'cvs' (sans les caractères ')]
$ cvs co 2.4routing
cvs server: Updating 2.4routing
U 2.4routing/lartc.db
```

Si vous avez fait des changements et que vous vouliez contribuer au HOWTO, exécutez `cvs -z3 diff -uBb`, et envoyez-nous le résultat par courrier électronique de façon à pouvoir facilement intégrer les modifications. Merci ! Au fait, soyez sûr que vous avez édité le fichier `.db`, les autres documents étant générés à partir de celui-ci.

Un fichier Makefile est fourni pour vous aider à créer des fichiers PostScript, dvi, pdf, html et texte. Vous pouvez avoir à installer les `docbook`, `docbook-utils`, `ghostscript` et `tetex` pour obtenir tous les formats de sortie.

Faites attention de ne pas éditer le fichier `2.4routing.sgml` ! Il contient une ancienne version du HOWTO. Le bon fichier est `lartc.db`.

2.6. Liste de diffusion

Les auteurs reçoivent de plus en plus de courriers électroniques à propos de cet HOWTO. Vu l'intérêt de la communauté, il a été décidé la mise en place d'une liste de diffusion où les personnes pourront discuter du routage avancé et du contrôle de trafic. Vous pouvez vous abonner à la liste [ici](#)⁵.

Il devra être noté que les auteurs sont très hésitants à répondre à des questions qui n'ont pas été posées sur la liste. Nous aimerions que la liste devienne une sorte de base de connaissance. Si vous avez une question, recherchez, s'il vous plaît, d'abord dans l'archive, et ensuite postez-là dans la liste de diffusion.

2.7. Plan du document

Nous allons essayer de faire des manipulations intéressantes dès le début, ce qui veut dire que tout ne sera pas expliqué en détail tout de suite. Veuillez passer sur ces détails, et accepter de considérer qu'ils deviendront clairs par la suite.

Le routage et le filtrage sont deux choses distinctes. Le filtrage est très bien documenté dans le HOWTO de Rusty, disponible [ici](#) :

- [The netfilter/iptables HOWTO's](#)⁶

Nous nous focaliserons principalement sur ce qu'il est possible de faire en combinant netfilter et iproute2.

⁴ <http://www.ds9a.nl/lartc>

⁵ <http://mailman.ds9a.nl/mailman/listinfo/lartc>

⁶ <http://netfilter.org/documentation/>

3.1. Pourquoi iproute2 ?

La plupart des distributions Linux et des UNIX utilisent couramment les vénérables commandes **arp**, **ifconfig** et **route**. Bien que ces outils fonctionnent, ils montrent quelques comportements inattendus avec les noyaux Linux des séries 2.2 et plus. Par exemple, les tunnels GRE font partie intégrante du routage de nos jours, mais ils nécessitent des outils complètement différents.

Avec iproute2, les tunnels font partie intégrante des outils.

Les noyaux Linux des séries 2.2 et plus ont un sous-système réseau complètement réécrit. Ce nouveau codage de la partie réseau apporte à Linux des performances et des fonctionnalités qui n'ont pratiquement pas d'équivalent parmi les autres systèmes d'exploitation. En fait, le nouveau logiciel de filtrage, routage et de classification possède plus de fonctionnalités que les logiciels fournis sur beaucoup de routeurs dédiés, de pare-feu et de produits de mise en forme (*shaping*) du trafic.

Dans les systèmes d'exploitation existants, au fur et à mesure que de nouveaux concepts réseau apparaissaient, les développeurs sont parvenus à les greffer sur les structures existantes. Ce travail constant d'empilage de couches a conduit à des codes réseau aux comportements étranges, un peu comme les langues humaines. Dans le passé, Linux émulait le mode de fonctionnement de SunOS, ce qui n'était pas l'idéal.

La nouvelle structure d'iproute2 a permis de formuler clairement des fonctionnalités impossibles à implémenter dans le sous-système réseau précédent.

3.2. Un tour d'horizon d'iproute2

Linux possède un système sophistiqué d'allocation de bande passante appelé Contrôle de trafic (*Traffic Control*). Ce système supporte différentes méthodes pour classer, ranger par ordre de priorité, partager et limiter le trafic entrant et sortant.

Nous commencerons par un petit tour d'horizon des possibilités d'iproute2.

3.3. Prérequis

Vous devez être sûr que vous avez installé les outils utilisateur (NdT : userland tools, par opposition à la partie « noyau » d'iproute2). Le paquet concerné s'appelle iproute sur RedHat et Debian. Autrement, il peut être trouvé à <ftp://ftp.inr.ac.ru/ip-routing/iproute2-2.2.4-now-ss?????.tar.gz>.

Vous pouvez aussi essayer iproute2-current.tar.gz¹ pour la dernière version.

Certains éléments d'iproute vous imposent l'activation de certaines options du noyau. Il devra également être noté que toutes les versions de RedHat jusqu'à la version 6.2 incluse n'ont pas les fonctionnalités du contrôle de trafic activées dans le noyau fourni par défaut.

RedHat 7.2 contient tous les éléments par défaut.

Soyez également sûr que vous avez le support netlink, même si vous devez choisir de compiler votre propre noyau ; iproute2 en a besoin.

3.4. Explorer votre configuration courante

Cela peut vous paraître surprenant, mais iproute2 est déjà configuré ! Les commandes courantes **ifconfig** et **route** utilisent déjà les appels système avancés d'iproute2, mais essentiellement avec les options par défaut (c'est-à-dire ennuyeuses).

L'outil **ip** est central, et nous allons lui demander de nous montrer les interfaces.

3.4.1. ip nous montre nos liens

```
[ahu@home ahu]$ ip link list
1: lo: <LOOPBACK,UP> mtu 3924 qdisc noqueue
   link/loopback 00:00:00:00:00:00 brd 00:00:00:00:00:00
2: dummy: <BROADCAST,NOARP> mtu 1500 qdisc noop
   link/ether 00:00:00:00:00:00 brd ff:ff:ff:ff:ff:ff
3: eth0: <BROADCAST,MULTICAST,PROMISC,UP> mtu 1400 qdisc pfifo_fast qlen 100
   link/ether 48:54:e8:2a:47:16 brd ff:ff:ff:ff:ff:ff
4: eth1: <BROADCAST,MULTICAST,PROMISC,UP> mtu 1500 qdisc pfifo_fast qlen 100
   link/ether 00:e0:4c:39:24:78 brd ff:ff:ff:ff:ff:ff
3764: ppp0: <POINTOPOINT,MULTICAST,NOARP,UP> mtu 1492 qdisc pfifo_fast qlen 10
   link/ppp
```

La sortie peut varier, mais voici ce qui est affiché pour mon routeur NAT (NdT : traduction d'adresse) chez moi. J'expliquerai seulement une partie de la sortie, dans la mesure où tout n'est pas directement pertinent.

La première interface que nous voyons est l'interface loopback. Bien que votre ordinateur puisse fonctionner sans, je vous le déconseille. La taille de MTU (unité maximum de transmission) est de 3924 octets, et loopback n'est pas supposé être mis en file d'attente, ce qui prend tout son sens dans la mesure où cette interface est le fruit de l'imagination de votre noyau.

Je vais passer sur l'interface dummy pour l'instant, et il se peut qu'elle ne soit pas présente sur votre ordinateur. Il y a ensuite mes deux interfaces physiques, l'une du côté de mon modem câble, l'autre servant mon segment ethernet à la maison. De plus, nous voyons une interface ppp0.

¹ <ftp://ftp.inr.ac.ru/ip-routing/iproute2-current.tar.gz>

Notons l'absence d'adresses IP. Iproute déconnecte les concepts de « liens » et « d'adresses IP ». Avec l'*IP aliasing*, le concept de l'adresse IP canonique est devenu, de toute façon, sans signification.

ip nous montre bien, cependant, l'adresse MAC, l'identifiant matériel de nos interfaces ethernet.

3.4.2. ip nous montre nos adresses IP

```
[ahu@home ahu]$ ip address show
1: lo: <LOOPBACK,UP> mtu 3924 qdisc noqueue
   link/loopback 00:00:00:00:00:00 brd 00:00:00:00:00:00
   inet 127.0.0.1/8 brd 127.255.255.255 scope host lo
2: dummy: <BROADCAST,NOARP> mtu 1500 qdisc noop
   link/ether 00:00:00:00:00:00 brd ff:ff:ff:ff:ff:ff
3: eth0: <BROADCAST,MULTICAST,PROMISC,UP> mtu 1400 qdisc pfifo_fast qlen 100
   link/ether 48:54:e8:2a:47:16 brd ff:ff:ff:ff:ff:ff
   inet 10.0.0.1/8 brd 10.255.255.255 scope global eth0
4: eth1: <BROADCAST,MULTICAST,PROMISC,UP> mtu 1500 qdisc pfifo_fast qlen 100
   link/ether 00:e0:4c:39:24:78 brd ff:ff:ff:ff:ff:ff
3764: ppp0: <POINTOPOINT,MULTICAST,NOARP,UP> mtu 1492 qdisc pfifo_fast qlen 10
   link/ppp
   inet 212.64.94.251 peer 212.64.94.1/32 scope global ppp0
```

Cela contient plus d'informations : **ip** montre toutes nos adresses, et à quelles cartes elles appartiennent. `inet` signifie Internet (IPv4). Il y a beaucoup d'autres familles d'adresses, mais elles ne nous concernent pas pour le moment.

Examinons l'interface `eth0` de plus près. Il est dit qu'elle est reliée à l'adresse internet `10.0.0.1/8`. Qu'est-ce que cela signifie ? Le `/8` désigne le nombre de bits réservés à l'adresse réseau. Il y a 32 bits, donc il reste 24 bits pour désigner une partie de notre réseau. Les 8 premiers bits de `10.0.0.1` correspondent à `10.0.0.0`, notre adresse réseau, et notre masque de sous-réseau est `255.0.0.0`.

Les autres bits repèrent des machines directement connectées à cette interface. Donc, `10.250.3.13` est directement disponible sur `eth0`, comme l'est `10.0.0.1` dans notre exemple.

Avec `ppp0`, le même concept existe, bien que les nombres soient différents. Son adresse est `212.64.94.251`, sans masque de sous-réseau. Cela signifie que vous avez une liaison point à point et que toutes les adresses, à l'exception de `212.64.94.251`, sont distantes. Il y a cependant plus d'informations. En effet, on nous dit que de l'autre côté du lien, il n'y a encore qu'une seule adresse, `212.64.94.1`. Le `/32` nous précise qu'il n'y a pas de « bits réseau ».

Il est absolument vital que vous compreniez ces concepts. Référez-vous à la documentation mentionnée au début de ce HOWTO si vous avez des doutes.

Vous pouvez aussi noter `qdisc`, qui désigne la gestion de la mise en file d'attente (*Queueing Discipline*). Cela deviendra vital plus tard.

3.4.3. ip nous montre nos routes

Nous savons maintenant comment trouver les adresses `10.x.y.z`, et nous sommes capables d'atteindre `212.64.94.1`. Cela n'est cependant pas suffisant, et nous avons besoin d'instructions pour atteindre le monde. L'Internet est disponible via notre connexion PPP, et il se trouve que `212.64.94.1` est prêt à propager nos paquets à travers le monde, et à nous renvoyer le résultat.

```
[ahu@home ahu]$ ip route show
212.64.94.1 dev ppp0 proto kernel scope link src 212.64.94.251
10.0.0.0/8 dev eth0 proto kernel scope link src 10.0.0.1
127.0.0.0/8 dev lo scope link
default via 212.64.94.1 dev ppp0
```

Cela se comprend de soi-même. Les 4 premières lignes donnent explicitement ce qui était sous-entendu par **ip address show**, la dernière ligne nous indiquant que le reste du monde peut être trouvé via `212.64.94.1`, notre passerelle par défaut. Nous pouvons voir que c'est une passerelle à cause du mot « `via` », qui nous indique que nous avons besoin d'envoyer les paquets vers `212.64.94.1`, et que c'est elle qui se chargera de tout.

En référence, voici ce que l'ancien utilitaire **route** nous propose :

```
[ahu@home ahu]$ route -n
Kernel IP routing table
Destination Gateway Genmask Flags Metric Ref Use
Iface
212.64.94.1 0.0.0.0 255.255.255.255 UH 0 0 0 ppp0
10.0.0.0 0.0.0.0 255.0.0.0 U 0 0 0 eth0
127.0.0.0 0.0.0.0 255.0.0.0 U 0 0 0 lo
0.0.0.0 212.64.94.1 0.0.0.0 UG 0 0 0 ppp0
```

3.5. ARP

ARP est le Protocole de Résolution d'Adresse (*Address Resolution Protocol*). Il est décrit dans le [RFC 826](http://www.faqs.org/rfcs/rfc826.html)². ARP est utilisé par une machine d'un réseau local pour retrouver l'adresse matérielle (la localisation) d'une autre machine sur le même réseau. Les machines sur Internet sont généralement connues par leur nom auquel correspond une adresse IP. C'est ainsi qu'une machine sur le réseau `foo.com` est capable de

² <http://www.faqs.org/rfcs/rfc826.html>

communiquer avec une autre machine qui est sur le réseau `bar.net`. Une adresse IP, cependant, ne peut pas vous indiquer la localisation physique de la machine. C'est ici que le protocole ARP entre en jeu.

Prenons un exemple très simple. Supposons que j'aie un réseau composé de plusieurs machines, dont la machine `foo` d'adresse IP `10.0.0.1` et la machine `bar` qui a l'adresse IP `10.0.0.2`. Maintenant, `foo` veut envoyer un **ping** vers `bar` pour voir s'il est actif, mais `foo` n'a aucune indication sur la localisation de `bar`. Donc, si `foo` décide d'envoyer un **ping** vers `bar`, il a besoin d'envoyer une requête ARP. Cette requête ARP est une façon pour `foo` de crier sur le réseau « `Bar (10.0.0.2) ! Où es-tu ?` ». Par conséquent, toutes les machines sur le réseau entendront `foo` crier, mais seul `bar (10.0.0.2)` répondra. `Bar` enverra une réponse ARP directement à `foo` ; ce qui revient à dire : « `Foo (10.0.0.1) ! je suis ici, à l'adresse 00:60:94:E:08:12` ». Après cette simple transaction utilisée pour localiser son ami sur le réseau, `foo` est capable de communiquer avec `bar` jusqu'à ce qu'il (le cache ARP de `foo`) oublie où `bar` est situé (typiquement au bout de 15 minutes sur Unix).

Maintenant, voyons comment cela fonctionne. Vous pouvez consulter votre cache (table) ARP (*neighbor*) comme ceci :

```
[root@espa041 /home/src/iputils]# ip neigh show
9.3.76.42 dev eth0 lladdr 00:60:08:3f:e9:f9 nud reachable
9.3.76.1 dev eth0 lladdr 00:06:29:21:73:c8 nud reachable
```

Comme vous pouvez le voir, ma machine `espa041 (9.3.76.41)` sait où trouver `espa042 (9.3.76.42)` et `espagate (9.3.76.1)`. Maintenant, ajoutons une autre machine dans le cache ARP.

```
[root@espa041 /home/paulsch/.gnome-desktop]# ping -c 1 espa043
PING espa043.austin.ibm.com (9.3.76.43) from 9.3.76.41 : 56(84) bytes of data.
64 bytes from 9.3.76.43: icmp_seq=0 ttl=255 time=0.9 ms
```

```
1 packets transmitted, 1 packets received, 0% packet loss
round-trip min/avg/max = 0.9/0.9/0.9 ms
```

```
[root@espa041 /home/src/iputils]# ip neigh show
9.3.76.43 dev eth0 lladdr 00:06:29:21:80:20 nud reachable
9.3.76.42 dev eth0 lladdr 00:60:08:3f:e9:f9 nud reachable
9.3.76.1 dev eth0 lladdr 00:06:29:21:73:c8 nud reachable
```

Par conséquent, lorsque `espa041` a essayé de contacter `espa043`, l'adresse matérielle de `espa043` (sa localisation) a alors été ajoutée dans le cache ARP. Donc, tant que la durée de vie de l'entrée correspondant à `espa043` dans le cache ARP n'est pas dépassée, `espa041` sait localiser `espa043` et n'a plus besoin d'envoyer de requête ARP.

Maintenant, effaçons `espa043` de notre cache ARP.

```
[root@espa041 /home/src/iputils]# ip neigh delete 9.3.76.43 dev eth0
[root@espa041 /home/src/iputils]# ip neigh show
9.3.76.43 dev eth0 nud failed
9.3.76.42 dev eth0 lladdr 00:60:08:3f:e9:f9 nud reachable
9.3.76.1 dev eth0 lladdr 00:06:29:21:73:c8 nud stale
```

Maintenant, `espa041` a à nouveau oublié la localisation d'`espa043` et aura besoin d'envoyer une autre requête ARP la prochaine fois qu'il voudra communiquer avec lui. Vous pouvez aussi voir ci-dessus que l'état d'`espagate (9.3.76.1)` est passé en *stale*. Cela signifie que la localisation connue est encore valide, mais qu'elle devra être confirmée à la première transaction avec cette machine.

Si vous avez un routeur important, il se peut que vous vouliez satisfaire les besoins de différentes personnes, qui peuvent être traitées différemment. Les bases de données des politiques de routage vous aident à faire cela, en gérant plusieurs ensembles de tables de routage.

Si vous voulez utiliser cette fonctionnalité, assurez-vous que le noyau est compilé avec les options IP : `Advanced router` et `IP : policy routing`.

Quand le noyau doit prendre une décision de routage, il recherche quelle table consulter. Par défaut, il y a trois tables. L'ancien outil **route** modifie les tables principale (*main*) et locale (*local*), comme le fait l'outil **ip** (par défaut).

Les règles par défaut :

```
[ahu@home ahu]$ ip rule list
0: from all lookup local
32766: from all lookup main
32767: from all lookup default
```

Ceci liste la priorité de toutes les règles. Nous voyons que toutes les règles sont appliquées à tous les paquets (*from all*). Nous avons vu la table *main* précédemment, sa sortie s'effectuant avec `ip route ls`, mais les tables *local* et *default* sont nouvelles.

Si nous voulons faire des choses fantaisistes, nous pouvons créer des règles qui pointent vers des tables différentes et qui nous permettent de redéfinir les règles de routage du système.

Pour savoir exactement ce que fait le noyau en présence d'un assortiment de règles plus complet, référez-vous à la documentation `ip-cref` d'Alexey [NdT : dans le paquet `iproute2` de votre distribution].

4.1. Politique de routage simple par l'adresse source

Prenons encore une fois un exemple réel. J'ai 2 modems câble, connectés à un routeur Linux NAT (*masquerading*). Les personnes habitant avec moi me paient pour avoir accès à Internet. Supposons qu'un de mes co-locataires consulte seulement hotmail et veuille payer moins. C'est d'accord pour moi, mais il utilisera le modem le plus lent.

Le modem câble « rapide » est connu sous `212.64.94.251` et est en liaison PPP avec `212.64.94.1`. Le modem câble « lent » est connu sous diverses adresses IP : `212.64.78.148` dans notre exemple avec un lien vers `195.96.98.253`.

La table locale :

```
[ahu@home ahu]$ ip route list table local
broadcast 127.255.255.255 dev lo proto kernel scope link src 127.0.0.1
local 10.0.0.1 dev eth0 proto kernel scope host src 10.0.0.1
broadcast 10.0.0.0 dev eth0 proto kernel scope link src 10.0.0.1
local 212.64.94.251 dev ppp0 proto kernel scope host src 212.64.94.251
broadcast 10.255.255.255 dev eth0 proto kernel scope link src 10.0.0.1
broadcast 127.0.0.0 dev lo proto kernel scope link src 127.0.0.1
local 212.64.78.148 dev ppp2 proto kernel scope host src 212.64.78.148
local 127.0.0.1 dev lo proto kernel scope host src 127.0.0.1
local 127.0.0.0/8 dev lo proto kernel scope host src 127.0.0.1
```

Il y a beaucoup de choses évidentes, mais aussi des choses qui ont besoin d'être précisées quelque peu, ce que nous allons faire. La table de routage par défaut est vide.

Regardons la table principale (*main*) :

```
[ahu@home ahu]$ ip route list table main
195.96.98.253 dev ppp2 proto kernel scope link src 212.64.78.148
212.64.94.1 dev ppp0 proto kernel scope link src 212.64.94.251
10.0.0.0/8 dev eth0 proto kernel scope link src 10.0.0.1
127.0.0.0/8 dev lo scope link
default via 212.64.94.1 dev ppp0
```

Maintenant, nous générons une nouvelle règle que nous appellerons *John*, pour notre hypothétique co-locataire. Bien que nous puissions travailler avec des nombres IP purs, il est plus facile d'ajouter notre table dans le fichier `/etc/iproute2/rt_tables`.

```
# echo 200 John >> /etc/iproute2/rt_tables
# ip rule add from 10.0.0.10 table John
# ip rule ls
0: from all lookup local
32765: from 10.0.0.10 lookup John
32766: from all lookup main
32767: from all lookup default
```

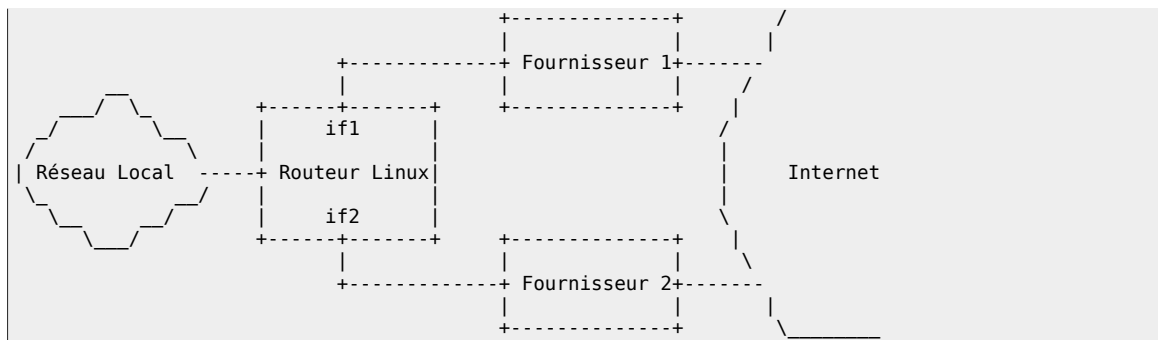
Maintenant, tout ce qu'il reste à faire est de générer la table *John*, et de vider le cache des routes :

```
# ip route add default via 195.96.98.253 dev ppp2 table John
# ip route flush cache
```

Et voilà qui est fait. Il ne reste plus, comme exercice laissé au lecteur, qu'à implémenter cela dans **ip-up**.

4.2. Routage avec plusieurs accès Internet/fournisseurs d'accès

Une configuration classique est la suivante, où deux fournisseurs d'accès permettent la connexion d'un réseau local (ou même d'une simple machine) à Internet.



Il y a généralement deux questions à se poser pour cette configuration.

4.2.1. Accès séparé

La première est de savoir comment router les réponses aux paquets entrants par un fournisseur particulier, disons le Fournisseur 1, vers ce même fournisseur.

Commençons par définir quelques symboles. **\$IF1** sera le nom de la première interface (if1 sur la figure au-dessus) et **\$IF2** le nom de la deuxième interface. **\$IP1** sera alors l'adresse IP associée à **\$IF1** et **\$IP2** sera l'adresse IP associée à **\$IF2**. **\$P1** sera l'adresse IP de la passerelle du fournisseur d'accès 1 et **\$P2** sera l'adresse IP de la passerelle du fournisseur d'accès 2. Enfin, **\$P1_NET** sera l'adresse réseau à l'intérieur duquel se situe **\$P1** et **\$P2_NET** sera l'adresse réseau à l'intérieur duquel se situe **\$P2**.

Deux tables de routage supplémentaires sont créées, par exemple **T1** et **T2**. Celles-ci sont ajoutées dans /etc/iproute2/rt_tables. La configuration du routage dans ces tables s'effectue de la façon suivante :

```
ip route add $P1_NET dev $IF1 src $IP1 table T1
ip route add default via $P1 table T1
ip route add $P2_NET dev $IF2 src $IP2 table T2
ip route add default via $P2 table T2
```

Rien de vraiment spectaculaire. Une route est simplement positionnée vers la passerelle et une route par défaut via cette passerelle est mise en place, comme nous le ferions dans le cas d'un seul fournisseur d'accès. Ici, les routes sont placées dans des tables séparées, une par fournisseur d'accès. Il est à noter que la route vers le réseau suffit, dans la mesure où elle indique comment trouver n'importe quel hôte dans ce réseau, ce qui inclut la passerelle.

La table de routage principale est maintenant configurée. C'est une bonne idée de router les éléments à destination d'un voisin direct à travers l'interface connectée à ce voisin. Notez les arguments "src" qui assurent que la bonne adresse IP source sera choisie.

```
ip route add $P1_NET dev $IF1 src $IP1
ip route add $P2_NET dev $IF2 src $IP2
```

Indiquez maintenant votre préférence pour votre route par défaut :

```
ip route add default via $P1
```

Vous configurez ensuite les règles de routage. Celles-ci définissent la table qui sera vraiment choisie pour le routage. Il faut s'assurer que le routage s'effectue à travers une interface donnée si vous avez l'adresse source correspondante :

```
ip rule add from $IP1 table T1
ip rule add from $IP2 table T2
```

Cet ensemble de commandes vous assure que toutes les réponses au trafic entrant sur une interface particulière seront envoyées par cette interface.



Avertissement

Notes d'un lecteur : si **\$P0_NET** est le réseau local et **\$IF0** est son interface, alors les entrées suivantes sont désirables :

```
ip route add $P0_NET dev $IF0 table T1
ip route add $P2_NET dev $IF2 table T1
ip route add 127.0.0.0/8 dev lo table T1
ip route add $P0_NET dev $IF0 table T2
ip route add $P1_NET dev $IF1 table T2
ip route add 127.0.0.0/8 dev lo table T2
```

Nous avons maintenant une configuration très basique. Elle marchera pour tous les processus exécutés sur le routeur lui-même, ainsi que pour le réseau local si celui-ci est masqué. Si ce n'est pas le cas, soit vous avez une plage d'adresses IP pour chaque fournisseur d'accès, soit vous masquez vers l'un des deux fournisseurs d'accès. Dans les deux cas, vous ajouterez des règles indiquant, en fonction de l'adresse IP de la machine du réseau local, vers quel fournisseur vous allez router.

4.2.2. Balance de charge

La seconde question concerne la balance de charge du trafic sortant vers les deux fournisseurs d'accès. Ceci n'est pas vraiment très dur si vous avez déjà configuré l'accès séparé comme décrit ci-dessus.

Au lieu de choisir l'un des deux fournisseurs d'accès comme route par défaut, celle-ci peut être une route multi-chemin. Par défaut, le noyau répartira les routes vers les deux fournisseurs d'accès. Ceci est réalisé de la façon suivante (construit également sur l'exemple de la section de l'accès séparé) :

```
ip route add default scope global nexthop via $P1 dev $IF1 weight 1 \  
nexthop via $P2 dev $IF2 weight 1
```

Ceci réalisera la balance des routes vers les deux fournisseurs. Les paramètres **weight** peuvent permettre de favoriser un fournisseur par rapport à un autre.

Il est à noter que la balance de charge ne sera pas parfaite dans la mesure où elle est basée sur les routes et que celles-ci sont mises dans des caches. Ceci signifie que les routes vers les sites les plus souvent utilisés passeront toujours par le même fournisseur d'accès.

De plus, si vous voulez vraiment mettre en oeuvre ceci, vous devriez également aller consulter les mises à jour de Julien Anastasov à <http://www.ssi.bg/~ja/#routes>¹ Elles rendront le travail plus facile.

¹ <http://www.ssi.bg/~ja/#routes>

Il y a trois sortes de tunnels sous Linux : l'IP dans un tunnel IP, le tunnel GRE et les tunnels qui existent en dehors du noyau (comme PPTP, par exemple).

5.1. Quelques remarques générales à propos des tunnels :

Les tunnels peuvent faire des choses très inhabituelles et vraiment sympas. Ils peuvent aussi absolument tout détraquer si vous ne les avez pas configurés correctement. Ne définissez pas votre route par défaut sur un tunnel, à moins que vous ne sachiez *EXACTEMENT* ce que vous faites.

De plus, le passage par un tunnel augmente le poids des en-têtes (*overhead*), puisqu'un en-tête IP supplémentaire est nécessaire. Typiquement, ce surcoût est de 20 octets par paquet. Donc, si la taille maximum de votre paquet sur votre réseau (MTU) est de 1500 octets, un paquet qui est envoyé à travers un tunnel sera limité à une taille de 1480 octets. Ce n'est pas nécessairement un problème, mais soyez sûr d'avoir bien étudié la fragmentation et le réassemblage des paquets IP quand vous prévoyez de relier des réseaux de grande taille par des tunnels. Et bien sûr, la manière la plus rapide de creuser un tunnel est de creuser des deux côtés.

5.2. IP dans un tunnel IP

Ce type de tunnel est disponible dans Linux depuis un long moment. Il nécessite deux modules, **ipip.o** et **new_tunnel.o**.

Disons que vous avez trois réseaux : 2 réseaux internes A et B, et un réseau intermédiaire C (ou disons Internet). Les caractéristiques du réseau A sont :

```
réseau 10.0.1.0
masque de sous-réseau 255.255.255.0
routeur 10.0.1.1
```

Le routeur a l'adresse 172.16.17.18 sur le réseau C.

et le réseau B :

```
réseau 10.0.2.0
masque de sous-réseau 255.255.255.0
routeur 10.0.2.1
```

Le routeur a l'adresse 172.19.20.21 sur le réseau C.

En ce qui concerne le réseau C, nous supposons qu'il transmettra n'importe quel paquet de A vers B et vice-versa. Il est également possible d'utiliser l'Internet pour cela.

Voici ce qu'il faut faire :

Premièrement, assurez-vous que les modules soient installés :

```
insmod ipip.o
insmod new_tunnel.o
```

Ensuite, sur le routeur du réseau A, faites la chose suivante :

```
ifconfig tunl0 10.0.1.1 pointopoint 172.19.20.21
route add -net 10.0.2.0 netmask 255.255.255.0 dev tunl0
```

et sur le routeur du réseau B :

```
ifconfig tunl0 10.0.2.1 pointopoint 172.16.17.18
route add -net 10.0.1.0 netmask 255.255.255.0 dev tunl0
```

Et quand vous aurez terminé avec votre tunnel :

```
ifconfig tunl0 down
```

Vite fait, bien fait. Vous ne pouvez pas transmettre les paquets de diffusion (*broadcast*), ni le trafic IPv6 à travers un tunnel IP-IP. Vous ne pouvez connecter que deux réseaux IPv4 qui, normalement, ne seraient pas capables de se « parler », c'est tout. Dans la mesure où la compatibilité a été conservée, ce code tourne depuis un bon moment, et il reste compatible depuis les noyaux 1.3. Le tunnel Linux IP dans IP ne fonctionne pas avec d'autres systèmes d'exploitation ou routeurs, pour autant que je sache. C'est simple, ça marche. Utilisez-le si vous le pouvez, autrement utilisez GRE.

5.3. Le tunnel GRE

GRE est un protocole de tunnel qui a été à l'origine développé par Cisco, et qui peut réaliser plus de choses que le tunnel IP dans IP. Par exemple, vous pouvez aussi transporter du trafic multi-diffusion (*multicast*) et de l'IPv6 à travers un tunnel GRE.

Dans Linux, vous aurez besoin du module `ip_gre.o`.

5.3.1. Le tunnel IPv4

Dans un premier temps, intéressons-nous au tunnel IPv4 :

Disons que vous avez trois réseaux : 2 réseaux internes A et B, et un réseau intermédiaire C (ou disons internet).

Les caractéristiques du réseau A sont :

```
réseau 10.0.1.0
masque de sous-réseau 255.255.255.0
routeur 10.0.1.1
```

Le routeur a l'adresse 172.16.17.18 sur le réseau C. Appelons ce réseau neta.

Et pour le réseau B :

```
réseau 10.0.2.0
masque de sous-réseau 255.255.255.0
routeur 10.0.2.1
```

Le routeur a l'adresse 172.19.20.21 sur le réseau C. Appelons ce réseau netb.

En ce qui concerne le réseau C, nous supposons qu'il transmettra n'importe quels paquets de A vers B et vice-versa. Comment et pourquoi, on s'en moque.

Sur le routeur du réseau A, nous faisons la chose suivante :

```
ip tunnel add netb mode gre remote 172.19.20.21 local 172.16.17.18 ttl 255
ip link set netb up
ip addr add 10.0.1.1 dev netb
ip route add 10.0.2.0/24 dev netb
```

Discutons un peu de cela. Sur la ligne 1, nous avons ajouté un périphérique tunnel, que nous avons appelé netb (ce qui est évident, dans la mesure où c'est là que nous voulons aller). De plus, nous lui avons dit d'utiliser le protocole GRE (mode gre), que l'adresse distante est 172.19.20.21 (le routeur de l'autre côté), que nos paquets « tunnelés » devront être générés à partir de 172.16.17.18 (ce qui autorise votre serveur à avoir plusieurs adresses IP sur le réseau C et ainsi vous permet de choisir laquelle sera utilisée pour votre tunnel) et que le champ TTL de vos paquets sera fixé à 255 (ttl 255).

La deuxième ligne active le périphérique.

Sur la troisième ligne, nous avons donné à cette nouvelle interface l'adresse 10.0.1.1. C'est bon pour de petits réseaux, mais quand vous commencez une exploitation minière (*BEAUCOUP* de tunnels !), vous pouvez utiliser une autre gamme d'adresses IP pour vos interfaces « tunnel » (dans cet exemple, vous pourriez utiliser 10.0.3.0).

Sur la quatrième ligne, nous positionnons une route pour le réseau B. Notez la notation différente pour le masque de sous-réseau. Si vous n'êtes pas familiarisé avec cette notation, voici comment ça marche : vous écrivez le masque de sous-réseau sous sa forme binaire, et vous comptez tous les 1. Si vous ne savez pas comment faire cela, rappelez-vous juste que 255.0.0.0 est /8, 255.255.0.0 est /16 et 255.255.255.0 est /24. Et 255.255.254.0 est /23, au cas où ça vous intéresserait.

Mais arrêtons ici, et continuons avec le routeur du réseau B.

```
ip tunnel add neta mode gre remote 172.16.17.18 local 172.19.20.21 ttl 255
ip link set neta up
ip addr add 10.0.2.1 dev neta
ip route add 10.0.1.0/24 dev neta
```

Et quand vous voudrez retirer le tunnel sur le routeur A :

```
ip link set netb down
ip tunnel del netb
```

Bien sûr, vous pouvez remplacer netb par neta pour le routeur B.

5.3.2. Le tunnel IPv6

Voir la section 6 pour une courte description des adresses IPv6.

À propos des tunnels.

Supposons que vous ayez le réseau IPv6 suivant, et que vous vouliez le connecter à une dorsale IPv6 (6bone) ou à un ami.

```
Réseau 3ffe:406:5:1:5:a:2:1/96
```

Votre adresse IPv4 est 172.16.17.18, et le routeur 6bone a une adresse IPv4 172.22.23.24.

```
ip tunnel add sixbone mode sit remote 172.22.23.24 local 172.16.17.18 ttl 255
ip link set sixbone up
ip addr add 3ffe:406:5:1:5:a:2:1/96 dev sixbone
ip route add 3ffe::/15 dev sixbone
```

Voyons cela de plus près. Sur la première ligne, nous avons créé un périphérique tunnel appelé sixbone. Nous lui avons affecté le mode sit (qui est le tunnel IPv6 sur IPv4) et lui avons dit où l'on va (remote) et d'où l'on vient (local). TTL est configuré à son maximum : 255. Ensuite, nous avons rendu le périphérique actif (up). Puis, nous avons ajouté notre propre adresse réseau et configuré une route pour 3ffe::/15 à travers le tunnel.

Les tunnels GRE constituent actuellement le type de tunnel préféré. C'est un standard qui est largement adopté, même à l'extérieur de la communauté Linux, ce qui constitue une bonne raison de l'utiliser.

5.4. Tunnels dans l'espace utilisateur

Il y a des dizaines d'implémentations de tunnels à l'extérieur du noyau. Les plus connues sont bien sûr PPP et PPTP, mais il y en a bien plus (certaines propriétaires, certaines sécurisés, d'autres qui n'utilisent pas IP), qui dépassent le cadre de ce HOWTO.

Par Marco Davids <marco@sara.nl>

NOTE au mainteneur :

En ce qui me concerne, ce tunnel IPv6-IPv4 n'est pas, par définition, un tunnel GRE. Vous pouvez réaliser un tunnel IPv6 sur IPv4 au moyen des périphériques tunnels GRE (tunnels GRE *N'IMPORTE QUOI* vers IPv4), mais le périphérique utilisé ici (sit) ne permet que des tunnels IPv6 sur IPv4, ce qui est quelque chose de différent.

6.1. Tunnel IPv6

Voici une autre application des possibilités de tunnels de Linux. Celle-ci est populaire parmi les premiers adeptes d'IPv6 ou les pionniers si vous préférez. L'exemple pratique décrit ci-dessous n'est certainement pas la seule manière de réaliser un tunnel IPv6. Cependant, c'est la méthode qui est souvent utilisée pour réaliser un tunnel entre Linux et un routeur Cisco IPv6 et l'expérience m'a appris que c'est ce type d'équipement que beaucoup de personnes ont. Dix contre un que ceci s'appliquera aussi pour vous ;-).

De petites choses à propos des adresses IPv6 :

Les adresses IPv6 sont, en comparaison avec les adresses IPv4, vraiment grandes : 128 bits contre 32 bits. Et ceci nous fournit la chose dont nous avons besoin : beaucoup, beaucoup d'adresses IP : 340,282,266,920,938,463,463,374,607,431,768,211,465 pour être précis. A part ceci, IPv6 (ou IPng génération suivante (*Next Generation*)) est supposé fournir des tables de routage plus petites sur les routeurs des dorsales Internet, une configuration plus simple des équipements, une meilleure sécurité au niveau IP et un meilleur support pour la Qualité de Service (QoS).

Un exemple : 2002:836b:9820:0000:0000:0000:836b:9886

Ecrire les adresses IPv6 peut être un peu lourd. Il existe donc des règles qui rendent la vie plus facile :

- Ne pas utiliser les zéros de tête, comme dans IPv4 ;
- Utiliser des double-points de séparation tous les 16 bits ou 2 octets ;
- Quand vous avez beaucoup de zéros consécutifs, vous pouvez écrire ::. Vous ne pouvez, cependant, faire cela qu'une seule fois par adresse et seulement pour une longueur de 16 bits.

L'adresse 2002:836b:9820:0000:0000:0000:836b:9886 peut être écrite 2002:836b:9820::836b:9886, ce qui est plus amical.

Un autre exemple : l'adresse 3ffe:0000:0000:0000:0000:0000:34A1:F32C peut être écrite 3ffe::20:34A1:F32C, ce qui est beaucoup plus court.

IPv6 a pour but d'être le successeur de l'actuel IPv4. Dans la mesure où cette technologie est relativement récente, il n'y a pas encore de réseau natif IPv6 à l'échelle mondiale. Pour permettre un développement rapide, la dorsale IPv6 (6bone) a été introduite.

Les réseaux natifs IPv6 sont interconnectés grâce à l'encapsulation du protocole IPv6 dans des paquets IPv4, qui sont envoyés à travers l'infrastructure IPv4 existante, d'un site IPv6 à un autre.

C'est dans cette situation que l'on monte un tunnel.

Pour être capable d'utiliser IPv6, vous devrez avoir un noyau qui le supporte. Il y a beaucoup de bons documents qui expliquent la manière de réaliser cela. Mais, tout se résume à quelques étapes :

- Récupérer une distribution Linux récente, avec une glibc convenable.
- Récupérer alors les sources à jour du noyau.

Si tout cela est fait, vous pouvez alors poursuivre en compilant un noyau supportant l'IPv6 :

- Aller dans `/usr/src/linux` et tapez :
- **make menuconfig**
- Choisir Networking Options
- Sélectionner The IPv6 protocol, IPv6: enable EUI-64 token format, IPv6: disable provider based addresses

ASTUCE :Ne compiler pas ces options en tant que module. Ceci ne marchera souvent pas bien.

En d'autres termes, compilez IPv6 directement dans votre noyau. Vous pouvez alors sauvegarder votre configuration comme d'habitude et entreprendre la compilation de votre noyau.

ASTUCE: Avant de faire cela, modifier votre Makefile comme suit : `EXTRAVERSION = -x ; --> ; EXTRAVERSION = -x-IPv6`

Il y a beaucoup de bonnes documentations sur la compilation et l'installation d'un noyau. Cependant, ce document ne traite pas de ce sujet. Si vous rencontrez des problèmes à ce niveau, allez et recherchez dans la documentation des renseignements sur la compilation du noyau Linux correspondant à vos propres spécifications.

Le fichier `/usr/src/linux/README` peut constituer un bon départ. Après avoir réalisé tout ceci, et redémarré avec votre nouveau noyau flambant neuf, vous pouvez lancer la commande `/sbin/ifconfig -a` et noter un nouveau périphérique `sit0`. SIT signifie Transition Simple d'Internet (*Simple Internet Transition*). Vous pouvez vous auto complimenter : vous avez maintenant franchi une étape importante vers IP, la prochaine génération ;-)

Passons maintenant à l'étape suivante. Vous voulez connecter votre hôte ou peut-être même tout votre réseau LAN à d'autres réseaux IPv6. Cela pourrait être la dorsale IPv6 « 6bone » qui a été spécialement mise en place dans ce but particulier.

Supposons que vous avez le réseau IPv6 suivant : `3ffe:604:6:8::/64` et que vous vouliez le connecter à une dorsale IPv6 ou à un ami. Notez, s'il vous plaît, que la notation sous-réseau `/64` est similaire à celle des adresses IPv4.

Votre adresse IPv4 est `145.100.24.181` et le routeur 6bone a l'adresse IPv4 `145.100.1.5`.

```
# ip tunnel add sixbone mode sit remote 145.100.1.5 [local 145.100.24.181 ttl 225]
# ip link set sixbone up
# ip addr add 3FFE:604:6:7::2/96 dev sixbone
# ip route add 3ffe::0/15 dev sixbone
```

Discutons de ceci. Dans la première ligne, nous avons créé un périphérique appelé `sixbone`. Nous lui avons donné l'attribut `sit` (`mode sit`) (qui est le tunnel IPv6 dans IPv4) et nous lui avons dit où aller (`remote`) et d'où nous venions (`local`). TTL est configuré à son maximum, 255.

Ensuite, nous avons rendu le périphérique actif (`up`). Après cela, nous avons ajouté notre propre adresse réseau et configuré une route pour `3ffe::/15` (qui est actuellement la totalité du 6bone) à travers le tunnel. Si la machine sur laquelle vous mettez en place tout ceci est votre passerelle IPv6, ajoutez alors les lignes suivantes :

```
# echo 1 >/proc/sys/net/ipv6/conf/all/forwarding
# /usr/local/sbin/radvd
```

En dernière instruction, `radvd` est un démon d'annonce comme `zebra` qui permet de supporter les fonctionnalités d'autoconfiguration d'IPv6. Recherchez le avec votre moteur de recherche favori. Vous pouvez vérifier les choses comme ceci :

```
# /sbin/ip -f inet6 addr
```

Si vous arrivez à avoir `radvd` tournant sur votre passerelle IPv6 et que vous démarrez une machine avec IPv6 sur votre réseau local, vous serez ravi de voir les bénéfices de l'autoconfiguration IPv6 :

```
# /sbin/ip -f inet6 addr
1: lo: <LOOPBACK,UP> mtu 3924 qdisc noqueue inet6 ::1/128 scope host

3: eth0: <BROADCAST,MULTICAST,UP> mtu 1500 qdisc pfifo_fast qlen 100
   inet6 3ffe:604:6:8:5054:4cff:fe01:e3d6/64 scope global dynamic
   valid_lft forever preferred_lft 604646sec inet6 fe80::5054:4cff:fe01:e3d6/10
   scope link
```

Vous pouvez maintenant configurer votre serveur de noms pour les adresses IPv6. Le type A a un équivalent pour IPv6 : AAAA. L'équivalent de `in-addr.arpa` est : `ip6.int`. Il y a beaucoup d'informations disponibles sur ce sujet.

Il y a un nombre croissant d'applications IPv6 disponibles, comme le shell sécurisé, `telnet`, `inetd`, le navigateur Mozilla, le serveur web Apache et beaucoup d'autres. Mais ceci est en dehors du sujet de ce document de routage ;-).

Du côté Cisco, la configuration ressemblera à ceci :

```
!
interface Tunnel1
description IPv6 tunnel
no ip address
no ip directed-broadcast
ipv6 address 3FFE:604:6:7::1/96
tunnel source Serial0
tunnel destination 145.100.24.181
tunnel mode ipv6ip
!
ipv6 route 3FFE:604:6:8::/64 Tunnel1
```

Si vous n'avez pas un Cisco à votre disposition, essayez un des prestataires de tunnel IPv6 disponible sur Internet. Ils sont prêts à configurer leur Cisco avec un tunnel supplémentaire pour vous, le plus souvent au moyen d'une agréable interface web. Cherchez *ipv6 tunnel broker* avec votre moteur de recherche favori.

A ce jour, deux versions d'IPSEC sont disponibles pour Linux. FreeS/WAN, qui fût la première implémentation majeure, existe pour les noyaux Linux 2.2 et 2.4. Ce projet a [un site officiel](#)¹ et également [un site non officiel](#)², qui est bien maintenu. FreeS/WAN n'a jamais été intégré dans le noyau pour un certain nombre de raisons. Celle qui est la plus souvent mentionnée concerne un problème "politique" avec les américains travaillant sur la cryptographie qui freinent son exportabilité. De plus, la mise en place de FreeS/WAN dans le noyau Linux est délicate, ce qui n'en fait pas un bon candidat pour une réelle intégration.

De plus, [des](#)³ personnes [se sont inquiétées](#)⁴ de la qualité du code. Pour configurer FreeS/WAN, de nombreuses [documentations](#)⁵ sont [disponibles](#)⁶.

Une implémentation native d'IPSEC est présente dans le noyau à partir de la version Linux 2.5.47. Elle a été écrite par Alexey Kuznetsov et Dave Miller, qui se sont inspirés des travaux du groupe USAGI IPv6. Avec cette fusion, les CryptoAPI de James Morris deviennent également une partie du noyau, qui fait ainsi vraiment du cryptage.

Ce HOWTO ne documente que la version 2.5 d'IPSEC. FreeS/WAN est recommandé pour l'instant pour les utilisateurs de Linux 2.4. Faites cependant attention, dans la mesure où sa configuration est différente de l'IPSEC natif. Il y a maintenant une [mise à jour](#)⁷ qui permet au code FreeS/WAN de l'espace utilisateur de fonctionner avec l'IPSEC natif de Linux.

A partir du noyau 2.5.49, IPSEC fonctionne sans l'ajout de mises à jour supplémentaires.



Note

Les outils de l'espace utilisateur sont disponibles [ici](#)⁸. Il y a plusieurs programmes disponibles ; celui qui est proposé dans le lien est basé sur Racoon.

Lors de la compilation du noyau, soyez sûr d'activer 'PF_KEY', 'AH' et tous les éléments de CryptoAPI !



Avertissement

L'auteur de ce chapitre est un complet nigaud en ce qui concerne IPSEC ! Si vous trouvez les inévitables erreurs, envoyez un courrier électronique à Bert Hubert <ahu@ds9a.nl>.

Tout d'abord, nous montrerons comment configurer manuellement une communication sécurisée entre deux hôtes. Une grande partie de ce processus peut être automatisée, mais nous le ferons ici à la main afin de comprendre ce qui se passe "sous le capot".

Passez à la section suivante si la seule gestion automatique des clés vous intéresse. Soyez cependant conscient que la compréhension de la gestion manuelle des clés est utile.

7.1. Introduction sur la gestion manuelle des clés

IPSEC est un sujet compliqué. De nombreuses informations sont disponibles en ligne. Ce HOWTO se concentrera sur la mise en place et à l'explication des principes de base. Tous les exemples sont basés sur Racoon dont le lien est donné au-dessus.



Note

Certaines configurations iptables rejettent les paquets IPSEC ! Pour transmettre IPSEC, utilisez : **iptables -A xxx -p 50 -j ACCEPT** et **'iptables -A xxx -p 51 -j ACCEPT**.

IPSEC offre une version sécurisée de la couche IP (Internet Protocol). La sécurité dans ce contexte prend deux formes : l'encryptage et l'authentification. Une vision naïve de la sécurité ne propose que le cryptage. On peut cependant montrer facilement que c'est insuffisant : il se peut que vous ayez une communication cryptée, mais vous n'avez aucune garantie que l'hôte distant est bien celui auquel vous pensez.

IPSEC supporte 'Encapsulated Security Payload' (Encapsulation Sécurisée de la Charge utile) (ESP) pour le cryptage et 'Authentication Header' (Entête d'Authentification) (AH) pour authentifier le partenaire distant. Vous pouvez configurer les deux, ou décider de ne faire que l'un des deux.

ESP et AH s'appuient tous les deux sur des Associations de Sécurité (Security Associations (SA)). Une Association de Sécurité (SA) consiste en une source, une destination et une instruction. Un exemple simple d'Association de Sécurité (SA) pour l'authentification peut ressembler à ceci :

```
add 10.0.0.11 10.0.0.216 ah 15700 -A hmac-md5 "1234567890123456";
```

Ceci indique que le trafic allant de 10.0.0.11 vers 10.0.0.216 a besoin d'un En-tête d'Authentification (AH) qui peut être signé en utilisant HMAC-MD et le secret 1234567890123456. Cette instruction est repérée

¹ <http://www.freeswan.org/>

² <http://www.freeswan.ca>

³ <http://www.edlug.ed.ac.uk/archive/Sep2002/msg00244.html>

⁴ <http://lists.freeswan.org/pipermail/design/2002-November/003901.html>

⁵ <http://www.freeswan.ca/code/old/freeswan-Snapshot/doc/index.html>

⁶ <http://www.freeswan.org/doc.html>

⁷ <http://gondor.apana.org.au/~herbert/freeswan/>

⁸ <http://sourceforge.net/projects/ipsec-tools>

par l'identificateur SPI (*Security Parameter Index*) 15700, dont nous parlerons plus par la suite. Le point intéressant à propos des Associations de Sécurité (SA) est qu'elles sont symétriques. Les deux cotés de la conversation partagent exactement la même Association de Sécurité (SA), qui n'est pas recopiée sur l'hôte distant. Notez cependant qu'il n'y a pas de règles "d'inversion automatique". Cette Association de Sécurité (SA) décrit une authentification possible de 10.0.0.11 vers 10.0.0.216. Pour un trafic bidirectionnel, deux Associations de Sécurité (SA) sont nécessaires.

Un exemple d'Association de Sécurité (SA) pour le cryptage ESP :

```
add 10.0.0.11 10.0.0.216 esp 15701 -E 3des-cbc "123456789012123456789012";
```

Ceci signifie que le trafic allant de 10.0.0.11 vers 10.0.0.216 est chiffré en utilisant 3des-cbc avec la clé 123456789012123456789012. L'identificateur SPI est 15701.

Jusqu'ici, nous avons vu que les Associations de Sécurité (SA) décrivent les instructions possibles, mais pas la politique qui indique quand ces SA doivent être utilisées. En fait, il pourrait y avoir un nombre arbitraire de SA presque identiques ne se différenciant que par les identificateurs SPI. Entre parenthèses, SPI signifie *Security Parameter Index*, ou Index du Paramètre de Sécurité en français. Pour faire vraiment du cryptage, nous devons décrire une politique. Cette politique peut inclure des choses comme "utiliser ipsec s'il est disponible" ou "rejeter le trafic à moins que vous ayez ipsec".

Une "Politique de Sécurité" (Security Policy (SP)) typique ressemble à ceci :

```
spdadd 10.0.0.216 10.0.0.11 any -P out ipsec
    esp/transport//require
    ah/transport//require;
```

Si cette configuration est appliquée sur l'hôte 10.0.0.216, cela signifie que tout le trafic allant vers 10.0.0.11 doit être encrypté et encapsulé dans un en-tête d'authentification AH. Notez que ceci ne décrit pas quelle SA sera utilisée. Cette détermination est un exercice laissé à la charge du noyau.

En d'autres termes, une Politique de Sécurité spécifique CE QUE nous voulons ; une Association de Sécurité décrit COMMENT nous le voulons.

Les paquets sortants sont étiquetés avec le SPI SA ('le comment') que le noyau utilise pour l'encryptage et l'authentification et l'hôte distant peut consulter les instructions de vérification et de décryptage correspondantes.

Ce qui suit est une configuration très simple permettant le dialogue de l'hôte 10.0.0.216 vers l'hôte 10.0.0.11 en utilisant l'encryptage et l'authentification. Notez que le trafic de retour de cette première version est en clair et que cette configuration ne doit pas être déployée.

Sur l'hôte 10.0.0.216 :

```
#!/sbin/setkey -f
add 10.0.0.216 10.0.0.11 ah 24500 -A hmac-md5 "1234567890123456";
add 10.0.0.216 10.0.0.11 esp 24501 -E 3des-cbc "123456789012123456789012";

spdadd 10.0.0.216 10.0.0.11 any -P out ipsec
    esp/transport//require
    ah/transport//require;
```

Sur l'hôte 10.0.0.11, nous donnons les mêmes Associations de Sécurité (SA). Nous ne donnons pas de Politique de Sécurité :

```
#!/sbin/setkey -f
add 10.0.0.216 10.0.0.11 ah 24500 -A hmac-md5 "1234567890123456";
add 10.0.0.216 10.0.0.11 esp 24501 -E 3des-cbc "123456789012123456789012";
```

Avec la mise en place de la configuration ci-dessus (ces fichiers peuvent être exécutés si 'setkey' est installé dans /sbin), la commande **ping 10.0.0.11** exécutée sur 10.0.0.216 va donner la sortie suivante avec tcpdump :

```
22:37:52 10.0.0.216 > 10.0.0.11: AH(spi=0x00005fb4,seq=0xa): ESP(spi=0x00005fb5,seq=0xa) (DF)
22:37:52 10.0.0.11 > 10.0.0.216: icmp: echo reply
```

Notez que le paquet de retour provenant de 10.0.0.11 est en effet complètement visible. Le paquet ping émis par 10.0.0.216 ne peut évidemment pas être lu par tcpdump, mais celui-ci montre l'Index du Paramètre de Sécurité (SPI) de l'AH, ainsi que l'ESP, qui indique à 10.0.0.11 comment vérifier l'authenticité de notre paquet et comment le décrypter.

Quelques éléments doivent être mentionnés. La configuration ci-dessus est proposée dans de nombreux exemples d'IPSEC, mais elle est très dangereuse. Le problème est qu'elle contient la politique indiquant à 10.0.0.216 comment traiter les paquets allant vers 10.0.0.11 et comment 10.0.0.11 doit traiter ces paquets, mais ceci n'INDIQUE pas à 10.0.0.11 de rejeter le trafic non authentifié et non encrypté !

N'importe qui peut maintenant insérer des données "spooquées" (NdT : usurpées) et entièrement non cryptées que 10.0.0.1 acceptera. Pour remédier à ceci, nous devons avoir sur 10.0.0.11 une Politique de Sécurité pour le trafic entrant :

```
#!/sbin/setkey -f
spdadd 10.0.0.216 10.0.0.11 any -P IN ipsec
    esp/transport//require
    ah/transport//require;
```

Ceci indique à 10.0.0.11 que tout le trafic venant de 10.0.0.216 nécessite d'avoir un ESP et AH valide.

Maintenant, pour compléter cette configuration, nous devons également renvoyer un trafic encrypté et authentifié. La configuration complète sur 10.0.0.216 est la suivante :

```
#!/sbin/setkey -f
flush;
spdflush;

# AH
add 10.0.0.11 10.0.0.216 ah 15700 -A hmac-md5 "1234567890123456";
add 10.0.0.216 10.0.0.11 ah 24500 -A hmac-md5 "1234567890123456";

# ESP
add 10.0.0.11 10.0.0.216 esp 15701 -E 3des-cbc "123456789012123456789012";
add 10.0.0.216 10.0.0.11 esp 24501 -E 3des-cbc "123456789012123456789012";

spdadd 10.0.0.216 10.0.0.11 any -P out ipsec
        esp/transport//require
        ah/transport//require;

spdadd 10.0.0.11 10.0.0.216 any -P in ipsec
        esp/transport//require
        ah/transport//require;
```

Et sur 10.0.0.11 :

```
#!/sbin/setkey -f
flush;
spdflush;

# AH
add 10.0.0.11 10.0.0.216 ah 15700 -A hmac-md5 "1234567890123456";
add 10.0.0.216 10.0.0.11 ah 24500 -A hmac-md5 "1234567890123456";

# ESP
add 10.0.0.11 10.0.0.216 esp 15701 -E 3des-cbc "123456789012123456789012";
add 10.0.0.216 10.0.0.11 esp 24501 -E 3des-cbc "123456789012123456789012";

spdadd 10.0.0.11 10.0.0.216 any -P out ipsec
        esp/transport//require
        ah/transport//require;

spdadd 10.0.0.216 10.0.0.11 any -P in ipsec
        esp/transport//require
        ah/transport//require;
```

Notez que, dans cet exemple, nous avons utilisé des clés identiques pour les deux directions du trafic. Ceci n'est cependant en aucun cas exigé.

Pour examiner la configuration que nous venons de créer, exécuter **setkey -D**, qui montre les SA ou **setkey -DP** qui montre les politiques configurées.

7.2. Gestion automatique des clés

Dans la section précédente, l'encryptage était configuré pour utiliser simplement le partage de secrets. En d'autres termes, pour rester sécurisé, nous devons transférer la configuration de notre encryptage à travers un tunnel sécurisé. Si nous avons configuré l'hôte distant par telnet, n'importe quel tiers pourrait avoir pris connaissance de notre secret partagé et, ainsi, notre configuration ne serait plus sûre.

De plus, puisque le secret est partagé, ce n'est pas un secret. L'hôte distant ne peut pas en faire grand chose, mais nous devons être sûrs d'utiliser un secret différent pour les communications avec tous nos partenaires. Ceci nécessite un grand nombre de clés. Pour 10 partenaires, nous devrions avoir au moins 50 secrets différents.

En plus du problème des clés symétriques, le renouvellement des clés est également nécessaire. Si un tiers écoute suffisamment le trafic, il peut être en position de retrouver la clé par rétro ingénierie. On peut s'en prémunir en modifiant la clé de temps en temps, mais ce processus a besoin d'être automatisé.

Un autre problème est que la gestion manuelle des clés décrite au-dessus impose de définir précisément les algorithmes et les longueurs de clés utilisées, ce qui nécessite une grande coordination avec l'hôte distant. Il serait préférable d'avoir la capacité à décrire une politique des clés plus large comme par exemple "Nous pouvons faire du 3DES et du Blowfish avec les longueurs de clés suivantes".

Pour résoudre ces problèmes, IPSEC fournit l'Echange de Clé sur Internet (Internet Key Exchange (IKE)) permettant d'automatiser l'échange de clés générées aléatoirement. Ces clés sont transmises en utilisant une technologie d'encryptage asymétrique négociée.

L'implémentation IPSEC de Linux 2.5 fonctionne avec le démon IKE "KAME racoon". Depuis le 9 novembre, la version de racoon présente la distribution iptools d'Alexey peut être compilée en supprimant, au préalable `#include <net/route.h>` dans deux fichiers. Je fournis une [version précompilée](http://ds9a.nl/ipsec/racoon.bz2)⁹.



Note

L'Echange de Clé sur Internet (IKE) doit avoir accès au port UDP 500. Soyez sûr que iptables ne le bloque pas.

⁹ <http://ds9a.nl/ipsec/racoon.bz2>

7.2.1. Théorie

Comme expliqué avant, la gestion automatique des clés réalise beaucoup d'opérations pour nous. Spécifiquement, il crée à la volée les Associations de Sécurité. Il ne configure cependant pas la politique pour nous, ce qui est le fonctionnement attendu.

Donc, pour bénéficier de IKE, configurez une politique, mais ne fournissez aucune Association de Sécurité. Si le noyau découvre qu'il y a une politique IPSEC, mais pas d'Association de Sécurité, il va le notifier au démon IKE qui va chercher à en négocier une.

De nouveau, rappelons que la Politique de Sécurité spécifie CE QUE nous voulons tandis que l'Association de Sécurité décrit COMMENT nous le voulons. L'utilisation de la gestion automatique des clés nous permet de ne spécifier que ce que nous voulons.

7.2.2. Exemple

Kame racoon possède un grand nombre d'options dont la plupart des valeurs par défaut sont corrects ; nous n'avons donc pas besoin de les modifier. Comme nous l'avons dit auparavant, l'opérateur doit définir une Politique de Sécurité, mais pas d'Associations de Sécurité. Nous laissons cette négociation au démon IKE.

Dans cet exemple, 10.0.0.1 et 10.0.0.216 sont encore une fois sur le point d'établir des communications sécurisées mais, cette fois, avec l'aide du démon racoon. Par soucis de simplification, cette configuration utilisera des clés pré-partagées, les redoutés 'secrets partagés'. Nous discuterons des certificats X.509 dans une section à part. Voir [Section 7.2.3, « Gestion automatique des clés en utilisant les certificats X.509 »](#).

Nous allons à peu près rester fidèle à la configuration par défaut, qui est identique sur les deux hôtes :

```
path pre_shared_key "/usr/local/etc/racoon/psk.txt";

remote anonymous
{
    exchange_mode aggressive,main;
    doi ipsec_doi;
    situation identity_only;

    my_identifier address;

    lifetime time 2 min; # sec,min,hour
    initial_contact on;
    proposal_check obey; # obey, strict or claim

    proposal {
        encryption_algorithm 3des;
        hash_algorithm sha1;
        authentication_method pre_shared_key;
        dh_group 2 ;
    }
}

sainfo anonymous
{
    pfs_group 1;
    lifetime time 2 min;
    encryption_algorithm 3des ;
    authentication_algorithm hmac_sha1;
    compression_algorithm deflate ;
}
```

Beaucoup de paramètres. Je pense que l'on peut encore en supprimer pour se rapprocher de la configuration par défaut. Remarquons ici quelques éléments notables. Nous avons configuré deux sections "anonymous", ce qui convient pour tous les hôtes distants. Ceci va ainsi faciliter les configurations supplémentaires. Il n'est pas nécessaire d'avoir de sections spécifiques à une machine particulière, à moins que vous ne le vouliez vraiment.

De plus, la configuration précise que nous nous identifions grâce à notre adresse IP ('my_identifier address') et que nous pouvons faire du 3des, sha1 et que nous utiliserons une clé "pré-partagée" se trouvant dans psk.txt.

Dans le fichier psk.txt, nous avons configuré deux entrées qui sont différentes suivant les hôtes. Sur 10.0.0.11 :

```
10.0.0.216 password2
```

Sur 10.0.0.216 :

```
10.0.0.11 password2
```

Soyez sûr que ces fichiers sont la propriété de root, et qu'ils ont le mode 0600. Dans le cas contraire, racoon ne pourra faire confiance à leur contenu. Notez que ces fichiers sont symétriques l'un de l'autre.

Nous sommes maintenant prêt à configurer notre politique qui est assez simple. Sur l'hôte 10.0.0.216 :

```
#!/sbin/setkey -f
flush;
spdflush;
```

```

spdadd 10.0.0.216 10.0.0.11 any -P out ipsec
    esp/transport//require;

spdadd 10.0.0.11 10.0.0.216 any -P in ipsec
    esp/transport//require;

```

Et sur 10.0.0.11 :

```

#!/sbin/setkey -f
flush;
spdf flush;

spdadd 10.0.0.11 10.0.0.216 any -P out ipsec
    esp/transport//require;

spdadd 10.0.0.216 10.0.0.11 any -P in ipsec
    esp/transport//require;

```

Noter que ces politiques sont encore une fois symétriques.

Nous sommes maintenant prêt à lancer racoon ! Une fois lancé, au moment où nous essayons une connexion un telnet depuis 10.0.0.11 vers 10.0.0.216, ou l'inverse, racoon aura démarré la négociation :

```

12:18:44: INFO: isakmp.c:1689:isakmp_post_acquire(): IPsec-SA
    request for 10.0.0.11 queued due to no phase found.
12:18:44: INFO: isakmp.c:794:isakmp_phlbegin_i(): initiate new
    phase 1 negotiation: 10.0.0.216[500]<=>10.0.0.11[500]
12:18:44: INFO: isakmp.c:799:isakmp_phlbegin_i(): begin Aggressive mode.
12:18:44: INFO: vendorid.c:128:check_vendorid(): received Vendor ID:
    KAME/racoon
12:18:44: NOTIFY: oakley.c:2037:oakley_skeyid(): couldn't find
    the proper pskey, try to get one by the peer's address.
12:18:44: INFO: isakmp.c:2417:log_phleestablished(): ISAKMP-SA
    established 10.0.0.216[500]-10.0.0.11[500] spi:044d25dede78a4d1:ff01e5b4804f0680
12:18:45: INFO: isakmp.c:938:isakmp_ph2begin_i(): initiate new phase 2
    negotiation: 10.0.0.216[0]<=>10.0.0.11[0]
12:18:45: INFO: pfkey.c:1106:pk_recvupdate(): IPsec-SA established:
    ESP/Transport 10.0.0.11->10.0.0.216 spi=44556347(0x2a7e03b)
12:18:45: INFO: pfkey.c:1318:pk_recvadd(): IPsec-SA established:
    ESP/Transport 10.0.0.216->10.0.0.11 spi=15863890(0xf21052)

```

L'exécution de la commande **setkey -D**, qui nous montre les Associations de Sécurité, nous indique qu'elles sont en effet présentes :

```

10.0.0.216 10.0.0.11
esp mode=transport spi=224162611(0x0d5c7333) reqxml:id=0(0x00000000)
E: 3des-cbc 5d421c1b d33b2a9f 4e9055e3 857db9fc 211d9c95 ebaead04
A: hmac-sha1 c5537d66 f3c5d869 bd736ae2 08d22133 27f7aa99
seq=0x00000000 replay=4 flags=0x00000000 state=mature
created: Nov 11 12:28:45 2002 current: Nov 11 12:29:16 2002
diff: 31(s) hard: 600(s) soft: 480(s)
last: Nov 11 12:29:12 2002 hard: 0(s) soft: 0(s)
current: 304(bytes) hard: 0(bytes) soft: 0(bytes)
allocated: 3 hard: 0 soft: 0
sadb_seq=1 pxml:id=17112 refcnt=0
10.0.0.11 10.0.0.216
esp mode=transport spi=165123736(0x09d79698) reqxml:id=0(0x00000000)
E: 3des-cbc d7af8466 acd4f14c 872c5443 ec45a719 d4b3fde1 8d239d6a
A: hmac-sha1 41ccc388 4568ac49 19e4e024 628e240c 141ffe2f
seq=0x00000000 replay=4 flags=0x00000000 state=mature
created: Nov 11 12:28:45 2002 current: Nov 11 12:29:16 2002
diff: 31(s) hard: 600(s) soft: 480(s)
last:
current: 231(bytes) hard: 0(bytes) soft: 0(bytes)
allocated: 2 hard: 0 soft: 0
sadb_seq=0 pxml:id=17112 refcnt=0

```

Nous avons les Politiques de Sécurité que nous avons nous-même configurées :

```

10.0.0.11[any] 10.0.0.216[any] tcp
in ipsec
esp/transport//require
created:Nov 11 12:28:28 2002 lastused:Nov 11 12:29:12 2002
lifetime:0(s) validtime:0(s)
spxml:id=3616 seq=5 pxml:id=17134
refcnt=3
10.0.0.216[any] 10.0.0.11[any] tcp
out ipsec
esp/transport//require
created:Nov 11 12:28:28 2002 lastused:Nov 11 12:28:44 2002
lifetime:0(s) validtime:0(s)
spxml:id=3609 seq=4 pxml:id=17134
refcnt=3

```

7.2.2.1. Problèmes et défauts connus

Si cela ne marche pas, vérifiez que tous les fichiers de configuration sont la propriété de root et qu'ils ne peuvent être lus que par celui-ci. Pour démarrer racoon en avant-plan, utilisez '-F'. Pour le forcer à lire un fichier de configuration à la place de celui précisé lors de la compilation, utilisez '-f'. Pour obtenir de nombreux détails, ajouter l'option 'log debug' dans le fichier racoon.conf.

7.2.3. Gestion automatique des clés en utilisant les certificats X.509

Comme nous l'avons dit avant, l'utilisation de secrets partagés est compliquée car ils ne peuvent pas être facilement partagés et, une fois qu'ils le sont, ils ne sont plus secrets. Heureusement, nous avons la technologie d'encryptage asymétrique pour nous aider à résoudre ce problème.

Si chaque participant d'une liaison IPSEC crée une clé publique et privée, des communications sécurisées peuvent être mises en place par les deux parties en publiant leur clé publique et en configurant leur politique.

Créer une clé est relativement facile, bien que cela exige un peu de travail. Ce qui suit est basé sur l'outil 'openssl'.

7.2.3.1. Construire un certificat X.509 pour votre hôte

OpenSSL dispose d'une importante infrastructure de gestion des clés, capable de gérer des clés signées ou non par une autorité de certification. Pour l'instant, nous avons besoin de court-circuiter toute cette infrastructure et de mettre en place une sécurité de charlatan, et de travailler sans autorité de certification.

Nous allons tout d'abord créer une requête de certificat (certificate request) pour notre hôte, appelé 'laptop' :

```
$ openssl req -new -nodes -newkey rsa:1024 -sha1 -keyform PEM -keyout \
laptop.private -outform PEM -out request.pem
```

Des questions nous sont posées :

```
Country Name (2 letter code) [AU]:NL
State or Province Name (full name) [Some-State]:.
Locality Name (eg, city) []:Delft
Organization Name (eg, company) [Internet Widgits Pty Ltd]:Linux Advanced
Routing & Traffic Control
Organizational Unit Name (eg, section) []:laptop
Common Name (eg, YOUR name) []:bert hubert
Email Address []:ahu@ds9a.nl
```

```
Please enter the following 'extra' attributes
to be sent with your certificate request
A challenge password []:
An optional company name []:
```

Vous avez toute liberté quant aux réponses. Vous pouvez ou non mettre le nom d'hôte, en fonction de vos besoins de sécurité. C'est ce que nous avons fait dans cet exemple.

Nous allons maintenant "auto signer" cette requête :

```
$ openssl x509 -req -in request.pem -signkey laptop.private -out \
laptop.public
Signature ok
subject=/C=NL/L=Delft/O=Linux Advanced Routing & Traffic \
Control/OU=laptop/CN=bert hubert/Email=ahu@ds9a.nl
Getting Private key
```

Le fichier "request.pem" peut maintenant être éliminé.

Répétez cette procédure pour tous les hôtes qui ont besoin d'une clé. Vous pouvez distribuer le fichier '.public' en toute impunité, mais garder le fichier '.private' privé !

7.2.3.2. Configuration et lancement

Une fois que nous avons les clés publiques et privées pour nos hôtes, nous pouvons indiquer à racoon de les utiliser.

Reprenons notre configuration précédente et les deux hôtes 10.0.0.11 ('upstairs') et 10.0.0.216 ('laptop').

Dans le fichier racoon.conf présent sur 10.0.0.11, nous ajoutons :

```
path certificate "/usr/local/etc/racoon/certs";

remote 10.0.0.216
{
    exchange_mode aggressive,main;
    my_identifiant asn1dn;
    peers_identifiant asn1dn;

    certificate_type x509 "upstairs.public" "upstairs.private";

    peers_certificate "laptop.public";
    proposal {
        encryption_algorithm 3des;
        hash_algorithm sha1;
        authentication_method rsasig;
        dh_group 2 ;
    }
}
```

Ceci indique à racoon que les certificats se trouvent dans /usr/local/etc/racoon/certs/. De plus, il contient des éléments spécifiques pour l'hôte distant 10.0.0.216.

La ligne 'asn1dn' indique à racoon que l'identification pour l'hôte local et distant doit être extraite des clés publiques. Ceci correspond à la ligne 'subject=/C=NL/L=Delft/O=Linux Advanced Routing & Traffic Control/OU=laptop/CN=bert hubert/Email=ahu@ds9a.nl' donné au-dessus.

La ligne **certificate_type** précise l'emplacement des clés publiques et privées locales. La déclaration **peers_certfile** précise à racoon que la clé publique de l'hôte distant se trouve dans le fichier `laptop.public`.

La section **proposal** reste inchangée par rapport à ce que nous avons vu plus tôt, à l'exception de **authentication_method** qui est maintenant **rsasig**, ce qui indique l'utilisation de clé RSA publique/privée pour l'authentification.

La configuration ajoutée sur 10.0.0.216 est presque identique, exception faite de l'habituelle symétrie :

```
path certificate "/usr/local/etc/racoon/certs";

remote 10.0.0.11
{
  exchange_mode aggressive,main;
  my_identifiant asn1dn;
  peers_identifiant asn1dn;

  certificate_type x509 "laptop.public" "laptop.private";

  peers_certfile "upstairs.public";

  proposal {
    encryption_algorithm 3des;
    hash_algorithm sha1;
    authentication_method rsasig;
    dh_group 2 ;
  }
}
```

Maintenant que nous avons ajouté ces éléments sur les deux hôtes, la seule chose qui reste à faire est de mettre en place les fichiers contenant les clés. La machine 'upstairs' doit avoir les fichiers `upstairs.private`, `upstairs.public`, et `laptop.public` placés dans `/usr/local/etc/racoon/certs`. Soyez sûr que le répertoire est la propriété de root et qu'il possède les droits 0700. Dans le cas contraire, racoon pourrait refuser de lire le contenu de ce répertoire.

La machine 'laptop' doit avoir les fichiers `upstairs.private`, `upstairs.public`, et `laptop.public` placés dans `/usr/local/etc/racoon/certs`. Autrement dit, chaque hôte doit avoir ses propres clés publique et privée et, de plus, la clé publique de l'hôte distant.

Vérifiez que la Politique de Sécurité est en place (exécutez la commande 'spdadd' vue dans [Section 7.2.2](#), « Exemple »). Lancez alors racoon et tout devrait fonctionner.

7.2.3.3. Comment configurer des tunnels sécurisés

Pour configurer des communications sécurisées avec un hôte distant, nous devons échanger des clés publiques. Bien qu'il ne soit pas nécessaire que la clé publique reste secrète, il est important d'être sûr que cette clé n'a pas été modifiée. En d'autres termes, vous devez être certain qu'il n'y a pas de 'man in the middle'. [NdT : 'man in the middle' est le nom d'une attaque qui consiste à se placer entre l'hôte émetteur et l'hôte de destination]

Pour faciliter ceci, OpenSSL propose la commande 'digest' :

```
$ openssl dgst upstairs.public
MD5(upstairs.public)= 78a3bddafb4d681c1ca8ed4d23da4ff1
```

La seule chose que nous devons faire est de vérifier que notre partenaire distant voit la même empreinte. Ceci peut être effectué en se rencontrant physiquement, ou par téléphone, en s'assurant que le numéro de téléphone de l'hôte distant n'a pas été envoyé dans le même courrier électronique que celui qui contenait la clé !

Une autre manière de faire ceci est d'utiliser un tiers de confiance qui exécute le service d'autorité de certification (*Certificate Authority*). Cette autorité de certification (CA) peut alors signer votre clé ; celle que nous avons nous-même créé au-dessus.

7.3. tunnels IPSEC

Jusqu'ici, nous n'avons seulement considéré IPSEC dans le mode appelé 'transport' où les points terminaux comprennent directement IPSEC. Comme ceci n'est pas souvent le cas, il peut être nécessaire d'avoir des routeurs qui, eux seuls, comprennent IPSEC et qui réalisent le travail pour les hôtes se trouvant derrière eux. Ceci est appelé le mode tunnel.

Configurer ceci est très rapide. Pour tunneler tout le trafic vers 130.161.0.0/16 à partir de 10.0.0.216 via 10.0.0.11, nous éditons ce qui suit sur 10.0.0.216 :

```
#!/sbin/setkey -f
flush;
spdflush;

add 10.0.0.216 10.0.0.11 esp 34501
-m tunnel
-E 3des-cbc "123456789012123456789012";

spdadd 10.0.0.0/24 130.161.0.0/16 any -P out ipsec
      esp/tunnel/10.0.0.216-10.0.0.11/require;
```

Notez que l'option '-m tunnel' est vitale ! Ceci configure tout d'abord une Association de Sécurité ESP entre les points terminaux de notre tunnel, à savoir 10.0.0.216 et 10.0.0.11.

Nous allons ensuite réellement configurer le tunnel. On doit indiquer au noyau d'encrypter tout le trafic de 10.0.0.0/24 vers 130.161.0.0. De plus, ce trafic doit être envoyé vers 10.0.0.11.

10.0.0.11 a également besoin d'être configuré :

```
#!/sbin/setkey -f
flush;
spdf flush;

add 10.0.0.216 10.0.0.11 esp 34501
-m tunnel
-E 3des-cbc "123456789012123456789012";

spdadd 10.0.0.0/24 130.161.0.0/16 any -P in ipsec
      esp/tunnel/10.0.0.216-10.0.0.11/require;
```

Notez que ceci est exactement identique, à l'exception du changement de '-P out' en '-P in'. Les exemples précédents n'ont configuré le trafic que dans un seul sens. Il est laissé comme exercice au lecteur le soin de compléter l'autre moitié du tunnel.

Le nom de 'proxy ESP' est également donné pour cette configuration, ce qui est un peu plus clair.



Note

Le tunnel IPSEC a besoin d'avoir la transmission IP activée dans le noyau !

7.4. Autre logiciel IPSEC

Thomas Walpuski précise qu'il a écrit une mise à jour pour que OpenBSD isakmpd puisse fonctionner avec Linux 2.5 IPSEC. De plus, la repository principale CVS de isakmpd contient maintenant le code ! Des notes sont disponibles [sur cette page](#)¹⁰.

isakmpd est différent de racoon mentionné au-dessus, mais de nombreuses personnes l'apprécient. Il peut être trouvé [ici](#)¹¹. D'autres éléments de lecture sur le CVS d'OpenBSD [ici](#)¹². Thomas a également créé un [tarball](#)¹³ pour ceux qui ne sont pas habitués à CVS ou patch.

De plus, des mises à jour sont disponibles pour permettre aux outils FreeS/WAN de l'espace utilisateur de fonctionner avec l'IPSEC natif de Linux 2.5. Vous pourrez les trouver [ici](#)¹⁴.

7.5. Interopérabilité d'IPSEC avec d'autres systèmes

FIXME: Ecrire ceci

7.5.1. Windows

Andreas Jellinghaus <aj@dungeon.inka.de> rapporte : "win2k: cela marche. pré-partage de clé et l'adresse ip pour l'authentification (je ne pense pas que windows supporte fdqn ou userfdqn). Les certificats devraient également marcher, mais cela n'a pas été essayé.

7.5.2. Check Point VPN-1 NG

Peter Bieringer rapporte :

```
Voici des résultats (seul le mode tunnel a été testé,
auth=SHA1) :
DES:      ok
3DES:     ok
AES-128:  ok
AES-192:  non supporté par CP VPN-1
AES-256:  ok
CAST* :   non supporté par le noyau Linux utilisé

Version Testée : FP4 aka R54 aka w/AI
```

Plus d'informations [ici](#)¹⁵.

¹⁰ <http://bender.thinknerd.de/~thomas/IPsec/isakmpd-linux.html>

¹¹ <http://www.openbsd.org/cgi-bin/cvsweb/src/sbin/isakmpd/>

¹² <http://www.openbsd.org/anoncvns.html>

¹³ <http://bender.thinknerd.de/~thomas/IPsec/isakmpd.tgz>

¹⁴ <http://gondor.apana.org.au/~herbert/freeswan/>

¹⁵ <http://www.fw-1.de/aerasesc/ng/vpn-racoon/CP-VPN1-NG-Linux-racoon.html>

FIXME: Pas de rédacteur !

Le Multicast-HOWTO est (relativement) ancien. De ce fait, il peut être imprécis ou induire en erreur à certains endroits.

Avant que vous ne puissiez faire du routage multidistribution, le noyau Linux a besoin d'être configuré pour supporter le type de routage multidistribution que vous voulez faire. Ceci, à son tour, exige une décision quant au choix du protocole de routage multidistribution que vous vous préparez à utiliser. Il y a essentiellement quatre types « communs » de protocoles : DVMRP (la version multidistribution du protocole RIP unicast), MOSPF (la même chose, mais pour OSPF), PIM-SM (*Protocol Independent Multicasting - Sparse Mode*) qui suppose que les utilisateurs de n'importe quel groupe de multidistribution sont dispersés plutôt que regroupés) et PIM-DM (le même, mais *Dense Mode*) qui suppose qu'il y aura un regroupement significatif des utilisateurs d'un même groupe de multidistribution.

On pourra noter que ces options n'apparaissent pas dans le noyau Linux. Ceci s'explique par le fait que le protocole lui-même est géré par une application de routage, comme Zebra, mrouterd ou pimd. Cependant, vous devez avoir une bonne idée de ce que vous allez utiliser, de manière à sélectionner les bonnes options dans le noyau.

Pour tout routage multidistribution, vous avez forcément besoin de sélectionner les options `multicasting` et `multicasting routing`. Ceci est suffisant pour DVMRP et MOSPF. Dans le cas de PIM, vous devez également valider les options `PIMv1` ou `PIMv2` suivant que le réseau que vous connectez utilise la version 1 ou 2 du protocole PIM.

Une fois que tout cela a été réalisé, et que votre nouveau noyau a été compilé, vous verrez au démarrage que IGMP est inclus dans la liste des protocoles IP. Celui-ci est un protocole permettant de gérer les groupes multidistribution. Au moment de la rédaction, Linux ne supportait que les versions 1 et 2 de IGMP, bien que la version 3 existe et ait été documentée. Ceci ne va pas vraiment nous affecter dans la mesure où IGMPv3 est encore trop récent pour que ses fonctionnalités supplémentaires soient largement utilisées. Puisque IGMP s'occupe des groupes, seules les fonctionnalités présentes dans la plus simple version de IGMP gérant un groupe entier seront utilisées. IGMPv2 sera utilisé dans la plupart des cas, bien que IGMPv1 puisse encore être rencontré.

Jusqu'à-là, c'est bon. Nous avons activé la multidistribution. Nous devons dire au noyau de l'utiliser concrètement. Nous allons donc démarrer le routage. Ceci signifie que nous ajoutons un réseau virtuel de multidistribution à la table du routeur :

```
ip route add 224.0.0.0/4 dev eth0
```

(En supposant bien sûr, que vous diffusez à travers eth0 ! Remplacez-le par le périphérique de votre choix, si nécessaire.)

Maintenant, dire à Linux de transmettre les paquets...

```
echo 1 > /proc/sys/net/ipv4/ip_forward
```

Arrivé ici, il se peut que vous vous demandiez si ceci va faire quelque chose. Donc, pour tester notre connexion, nous pinguons le groupe par défaut, 224.0.0.1, pour voir si des machines sont présentes. Toutes les machines du réseau local avec la multidistribution activée *DEVRAIENT* répondre, et aucune autre. Vous remarquerez qu'aucune des machines qui répondent ne le fait avec l'adresse IP 224.0.0.1. Quelle surprise ! :) Ceci est une adresse de groupe (une « diffusion » pour les abonnés) et tous les membres du groupe répondront avec leur propre adresse, et non celle du groupe.

```
ping -c 2 224.0.0.1
```

Maintenant, vous êtes prêt à faire du vrai routage multidistribution. Bien, en supposant que vous ayez deux réseaux à router l'un vers l'autre.

(A continuer !)

Quand je l'ai découvert, cela m'a *VRAIMENT* soufflé. Linux 2.2 contient toutes les fonctionnalités pour la gestion de la bande passante, de manière comparable à un système dédié de haut niveau.

Linux dépasse même ce que l'ATM et le Frame peuvent fournir.

Afin d'éviter toute confusion, voici les règles utilisées par **tc** pour la spécification de la bande passante :

```
mbps = 1024 kbps = 1024 * 1024 bps => byte/s (octets/s)
mbit = 1024 kbit => kilo bit/s.
mb = 1024 kb = 1024 * 1024 b => byte (octet)
mbit = 1024 kbit => kilo bit.
```

En interne, les nombres sont stockés en bps (octet/s) et b (octet).

Mais **tc** utilise l'unité suivante lors de l'affichage des débits :

```
1Mbit = 1024 Kbit = 1024 * 1024 bps => octets/s
```

9.1. Explication sur les files d'attente et la gestion de la mise en file d'attente

Avec la mise en file d'attente, nous déterminons la manière dont les données sont *ENVOYEES*. Il est important de comprendre que nous ne pouvons mettre en forme que les données que nous transmettons.

Avec la manière dont Internet travaille, nous n'avons pas de contrôle direct sur ce que les personnes nous envoient. C'est un peu comme votre boîte aux lettres (physique !) chez vous. Il n'y a pas de façon d'influencer le nombre de lettres que vous recevez, à moins de contacter tout le monde.

Cependant, l'Internet est principalement basé sur TCP/IP qui possède quelques fonctionnalités qui vont pouvoir nous aider. TCP/IP n'a pas d'aptitude à connaître les performances d'un réseau entre deux hôtes. Il envoie donc simplement des paquets de plus en plus rapidement (« *slow start* ») et quand des paquets commencent à se perdre, il ralentit car il n'a plus la possibilité de les envoyer. En fait, c'est un peu plus élégant que cela, mais nous en dirons plus par la suite.

C'est comme si vous ne lisiez que la moitié de votre courrier en espérant que vos correspondants arrêteront de vous en envoyer. À la différence que ça marche sur Internet :-)

Si vous avez un routeur et que vous souhaitez éviter que certains hôtes de votre réseau aient des vitesses de téléchargement trop grandes, vous aurez besoin de mettre en place de la mise en forme de trafic sur l'interface *INTERNE* de votre routeur, celle qui envoie les données vers vos propres ordinateurs.

Vous devez également être sûr que vous contrôlez le goulot d'étranglement de la liaison. Si vous avez une carte réseau à 100Mbit et un routeur avec un lien à 256kbit, vous devez vous assurer que vous n'envoyez pas plus de données que ce que le routeur peut manipuler. Autrement, ce sera le routeur qui contrôlera le lien et qui mettra en forme la bande passante disponible. Nous devons pour ainsi dire « être le propriétaire de la file d'attente » et être le lien le plus lent de la chaîne. Heureusement, c'est facilement réalisable.

9.2. Gestionnaires de mise en file d'attente simples, sans classes

Comme nous l'avons déjà dit, la gestion de mise en file d'attente permet de modifier la façon dont les données sont envoyées. Les gestionnaires de mise en file d'attente sans classes sont ceux qui, en gros, acceptent les données et qui ne font que les réordonner, les retarder ou les jeter.

Ils peuvent être utilisés pour mettre en forme le trafic d'une interface sans aucune subdivision. Il est primordial que vous compreniez cet aspect de la mise en file d'attente avant de continuer sur les gestionnaires de mise en files d'attente basés sur des classes contenant d'autres gestionnaires de mise en file d'attente.

Le gestionnaire le plus largement utilisé est de loin `pfifo_fast`, qui est celui par défaut. Ceci explique aussi pourquoi ces fonctionnalités avancées sont si robustes. Elles ne sont rien de plus « qu'une autre file d'attente ».

Chacune de ces files d'attente a ses forces et ses faiblesses. Toutes n'ont peut-être pas été bien testées.

9.2.1. `pfifo_fast`

Cette file d'attente, comme son nom l'indique : premier entré, premier sorti (*First In First Out*), signifie que les paquets ne subissent pas de traitements spéciaux. En fait, ce n'est pas tout à fait vrai. Cette file d'attente a trois « bandes ». A l'intérieur de chacune de ces bandes, des règles FIFO s'appliquent. Cependant, tant qu'il y a un paquet en attente dans la bande 0, la bande 1 ne sera pas traitée. Il en va de même pour la bande 1 et la bande 2.

Le noyau prend en compte la valeur du champ Type de Service des paquets et prend soin d'insérer dans la bande 0 les paquets ayant le bit « délai minimum » activé.

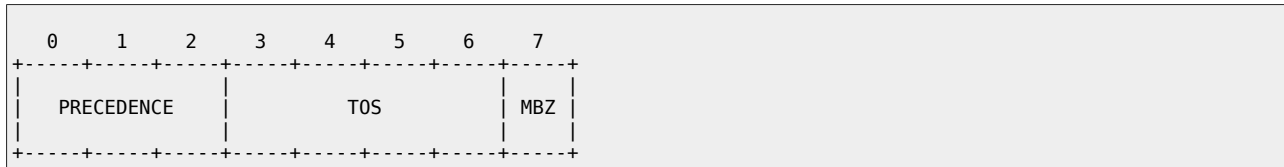
Ne pas confondre ce gestionnaire de mise en file d'attente sans classes avec celui basé sur des classes *PRIO* ! Bien qu'ils aient des comportements similaires, `pfifo_fast` ne possède pas de classes et vous ne pourrez pas y ajouter de nouveaux gestionnaires avec la commande **tc**.

9.2.1.1. Paramètres & usage

Vous ne pouvez pas configurer le gestionnaire `pfifo_fast`, dans la mesure où c'est celui par défaut. Voici sa configuration par défaut :

prionmap

Détermine comment les priorités des paquets sont reliées aux bandes, telles que définies par le noyau. La relation est établie en se basant sur l'octet TOS du paquet, qui ressemble à ceci :



Les quatre bits TOS (le champ TOS) sont définis comme suit :

Binaire	Décimal	Signification
1000	8	Minimise le Délai (Minimize delay) (md)
0100	4	Maximalise le Débit (Maximize throughput) (mt)
0010	2	Maximalise la Fiabilité (Maximize reliability) (mr)
0001	1	Minimalise le Coût Monétaire (Minimize monetary cost) (mmc)
0000	0	Service Normal

Comme il y a 1 bit sur la droite de ces quatre bits, la valeur réelle du champ TOS est le double de la valeur des bits TOS. `tcpdump -v -v` fournit la valeur de tout le champ TOS, et non pas seulement la valeur des quatre bits. C'est la valeur que l'on peut voir dans la première colonne du tableau suivant :

TOS	Bits	Signification	Priorité Linux	Bande
0x0	0	Service Normal	0 Best Effort	1
0x2	1	Minimise le Coût Monétaire (mmc)	1 Filler	2
0x4	2	Maximalise la Fiabilité (mr)	0 Best Effort	1
0x6	3	mmc+mr	0 Best Effort	1
0x8	4	Maximalise le Débit (mt)	2 Masse	2
0xa	5	mmc+mt	2 Masse	2
0xc	6	mr+mt	2 Masse	2
0xe	7	mmc+mr+mt	2 Masse	2
0x10	8	Minimise le Délai (md)	6 Interactive	0
0x12	9	mmc+md	6 Interactive	0
0x14	10	mr+md	6 Interactive	0
0x16	11	mmc+mr+md	6 Interactive	0
0x18	12	mt+md	4 Int. Masse	1
0x1a	13	mmc+mt+md	4 Int. Masse	1
0x1c	14	mr+mt+md	4 Int. Masse	1
0x1e	15	mmc+mr+mt+md	4 Int. Masse	1

[NdT : par flux de masse (*bulk flow*), il faut entendre « gros flot de données transmises en continu » comme un transfert FTP. A l'opposé, un flux interactif (*interactive flow*), correspond à celui généré par des requêtes SSH].

Beaucoup de nombres. La seconde colonne contient la valeur correspondante des quatre bits TOS, suivi de leur signification. Par exemple, 15 représente un paquet voulant un coût monétaire minimal, une fiabilité maximum, un débit maximum *ET* un délai minimum. J'appellerai ceci un « paquet Hollandais ».

La quatrième colonne liste la manière dont le noyau Linux interprète les bits TOS, en indiquant à quelle priorité ils sont reliés.

La dernière colonne montre la carte des priorités par défaut. Sur la ligne de commande, la carte des priorités ressemble à ceci :

```
1, 2, 2, 2, 1, 2, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1
```

Ceci signifie, par exemple, que la priorité 4 sera reliée à la bande numéro 1. La carte des priorités vous permet également de lister des priorités plus grandes (> 7) qui ne correspondent pas à une relation avec le champ TOS, mais qui sont configurées par d'autres moyens.

Le tableau suivant provenant de la RFC 1349 (à lire pour plus de détails) indique comment les applications devraient configurer leurs bits TOS pour fonctionner correctement :

TELNET		1000	(minimise le délai)
FTP	Contrôle	1000	(minimise le délai)
	Données	0100	(maximalise le débit)
TFTP		1000	(minimise le délai)
SMTP	phase de commande	1000	(minimise le délai)
	phase DATA	0100	(maximalise le débit)
Domain Name Service	requête UDP	1000	(minimise le délai)
	requête TCP	0000	
	Transfert de Zone	0100	(maximalise le débit)

NNTP	0001	(minimise le coût monétaire)
ICMP		
Erreurs	0000	
Requêtes	0000	(presque)
Réponses	<même chose que requête>	(presque)

txqueuelen

La longueur de cette file d'attente est fournie par la configuration de l'interface, que vous pouvez voir et configurer avec **ifconfig** et **ip**. Pour configurer la longueur de la file d'attente à 10, exécuter : **ifconfig eth0 txqueuelen 10**

Vous ne pouvez pas configurer ce paramètre avec **tc** !

9.2.2. Filtre à seau de jetons (*Token Bucket Filter*)

Le *Token Bucket Filter* (TBF) est un gestionnaire de mise en file d'attente simple. Il ne fait que laisser passer les paquets entrants avec un débit n'excédant pas une limite fixée administrativement. L'envoi de courtes rafales de données avec un débit dépassant cette limite est cependant possible.

TBF est très précis, et peu gourmand du point de vue réseau et processeur. Considérez-le en premier si vous voulez simplement ralentir une interface !

L'implémentation TBF consiste en un tampon (seau), constamment rempli par des éléments virtuels d'information appelés jetons, avec un débit spécifique (débit de jeton). Le paramètre le plus important du tampon est sa taille, qui correspond au nombre de jetons qu'il peut stocker.

Chaque jeton entrant laisse sortir un paquet de données de la file d'attente de données et ce jeton est alors supprimé du seau. L'association de cet algorithme avec les deux flux de jetons et de données, nous conduit à trois scénarios possibles :

- Les données arrivent dans TBF avec un débit *EGAL* au débit des jetons entrants. Dans ce cas, chaque paquet entrant a son jeton correspondant et passe la file d'attente sans délai.
- Les données arrivent dans TBF avec un débit *PLUS PETIT* que le débit des jetons. Seule une partie des jetons est supprimée au moment où les paquets de données sortent de la file d'attente, de sorte que les jetons s'accumulent jusqu'à atteindre la taille du tampon. Les jetons libres peuvent être utilisés pour envoyer des données avec un débit supérieur au débit des jetons standard, si de courtes rafales de données arrivent.
- Les données arrivent dans TBF avec un débit *PLUS GRAND* que le débit des jetons. Ceci signifie que le seau sera bientôt dépourvu de jetons, ce qui provoque l'arrêt de TBF pendant un moment. Ceci s'appelle « une situation de dépassement de limite » (*overlimit situation*). Si les paquets continuent à arriver, ils commenceront à être éliminés.

Le dernier scénario est très important, car il autorise la mise en forme administrative de la bande passante disponible pour les données traversant le filtre.

L'accumulation de jetons autorise l'émission de courtes rafales de données sans perte en situation de dépassement de limite, mais toute surcharge prolongée causera systématiquement le retard des paquets, puis leur rejet.

Notez que, dans l'implémentation réelle, les jetons correspondent à des octets, et non des paquets.

9.2.2.1. Paramètres & usage

Même si vous n'aurez probablement pas besoin de les changer, TBF a des paramètres. D'abord, ceux toujours disponibles sont :

limit or latency

Limit est le nombre d'octets qui peuvent être mis en file d'attente en attendant la disponibilité de jetons. Vous pouvez également indiquer ceci d'une autre manière en configurant le paramètre latency, qui spécifie le temps maximal pendant lequel un paquet peut rester dans TBF. Ce dernier paramètre prend en compte la taille du seau, le débit, et s'il est configuré, le débit de crête (peakrate).

burst/buffer/maxburst

Taille du seau, en octets. C'est la quantité maximale, en octets, de jetons dont on disposera simultanément. En général, plus les débits de mise en forme sont importants, plus le tampon doit être grand. Pour 10 Mbit/s sur plateforme Intel, vous avez besoin d'un tampon d'au moins 10 kilo-octets si vous voulez atteindre la limitation configurée !

Si votre tampon est trop petit, les paquets pourront être rejetés car il arrive plus de jetons par top d'horloge que ne peut en contenir le tampon.

mpu

Un paquet de taille nulle n'utilise pas une bande passante nulle. Pour ethernet, la taille minimale d'un paquet est de 64 octets. L'Unité Minimale de Paquet (*Minimum Packet Unit*) détermine le nombre minimal de jetons à utiliser pour un paquet.

rate

Le paramètre de la vitesse. Voir les remarques au-dessus à propos des limites !

Si le seau contient des jetons et qu'il est autorisé à se vider, alors il le fait par défaut avec une vitesse infinie. Si ceci vous semble inacceptable, utilisez les paramètres suivants :

peakrate

Si des jetons sont disponibles et que des paquets arrivent, ils sont immédiatement envoyés par défaut ; et pour ainsi dire à « la vitesse de la lumière ». Cela peut ne pas vous convenir, spécialement si vous avez un grand seau.

Le débit de crête (*peak rate*) peut être utilisé pour spécifier la vitesse à laquelle le seau est autorisé à se vider. Si tout se passe comme écrit dans les livres, ceci est réalisé en libérant un paquet, puis en attendant suffisamment longtemps, pour libérer le paquet suivant. Le temps d'attente est calculé de manière à obtenir un débit égal au débit de crête.

Cependant, étant donné que la résolution du minuteur (*timer*) d'UNIX est de 10 ms et que les paquets ont une taille moyenne de 10 000 bits, nous sommes limités à un débit de crête de 1mbit/s !

mtu/minburst

Le débit de crête de 1Mb/s ne sert pas à grand chose si votre débit habituel est supérieur à cette valeur. Un débit de crête plus élevé peut être atteint en émettant davantage de paquets par top du minuteur, ce qui a pour effet de créer un second seau.

Ce second *bucket* ne prend par défaut qu'un seul paquet, et n'est donc en aucun cas un seau.

Pour calculer le débit de crête maximum, multipliez le *mtu* que vous avez configuré par 100 (ou plus exactement par HZ, qui est égal à 100 sur Intel et à 1024 sur Alpha).

9.2.2.2. Configuration simple

Voici une configuration simple, mais *très* utile :

```
# tc qdisc add dev ppp0 root tbf rate 220kbit latency 50ms burst 1540
```

Pourquoi est-ce utile ? Si vous avez un périphérique réseau avec une grande file d'attente, comme un modem DSL ou un modem câble, et que le dialogue se fasse à travers une interface rapide, comme une interface ethernet, vous observerez que télécharger vers l'amont (*uploading*) dégrade complètement l'interactivité.

[NdT : *uploading* désigne une opération qui consiste à transférer des données ou des programmes stockés dans un ordinateur local vers un ordinateur distant à travers un réseau. La traduction officielle pour ce terme est « téléchargement vers l'amont ». On parle alors de voie montante. Le *downloading* désigne l'opération inverse (transfert d'un hôte distant vers l'ordinateur local) et est traduit par « téléchargement » ou « téléchargement vers l'aval ». On parle alors de la voie descendante.]

Le téléchargement vers l'amont va en effet remplir la file d'attente du modem. Celle-ci est probablement *ENORME* car cela aide vraiment à obtenir de bon débit de téléchargement vers l'amont. Cependant, ceci n'est pas forcément ce que voulez. Vous ne voulez pas forcément avoir une file d'attente importante de manière à garder l'interactivité et pouvoir encore faire des choses pendant que vous envoyez des données.

La ligne de commande au-dessus ralentit l'envoi de données à un débit qui ne conduit pas à une mise en file d'attente dans le modem. La file d'attente réside dans le noyau Linux, où nous pouvons lui imposer une taille limite.

Modifier la valeur 220kbit avec votre vitesse de lien *REELLE* moins un petit pourcentage. Si vous avez un modem vraiment rapide, augmenter un peu le paramètre burst.

9.2.3. Mise en file d'attente stochastiquement équitable (*Stochastic Fairness Queueing*)

Stochastic Fairness Queueing (SFQ) est une implémentation simple de la famille des algorithmes de mise en file d'attente équitable. Cette implémentation est moins précise que les autres, mais elle nécessite aussi moins de calculs tout en étant presque parfaitement équitable.

Le mot clé dans SFQ est conversation (ou flux), qui correspond principalement à une session TCP ou un flux UDP. Le trafic est alors divisé en un grand nombre de jolies files d'attente FIFO : une par conversation. Le trafic est alors envoyé dans un tourniquet, donnant une chance à chaque session d'envoyer leurs données tour à tour.

Ceci conduit à un comportement très équitable et empêche qu'une seule conversation étouffe les autres. SFQ est appelé « Stochastic » car il n'alloue pas vraiment une file d'attente par session, mais a un algorithme qui divise le trafic à travers un nombre limité de files d'attente en utilisant un algorithme de hachage.

A cause de ce hachage, plusieurs sessions peuvent finir dans le même seau, ce qui peut réduire de moitié les chances d'une session d'envoyer un paquet et donc réduire de moitié la vitesse effective disponible. Pour empêcher que cette situation ne devienne importante, SFQ change très souvent son algorithme de hachage pour que deux sessions entrantes en collision ne le fassent que pendant un nombre réduit de secondes.

Il est important de noter que SFQ n'est seulement utile que dans le cas où votre interface de sortie est vraiment saturée ! Si ce n'est pas le cas, il n'y aura pas de files d'attente sur votre machine Linux et donc,

pas d'effets. Plus tard, nous décrirons comment combiner SFQ avec d'autres gestionnaires de mise en files d'attente pour obtenir le meilleur des deux mondes.

Configurer spécialement SFQ sur l'interface ethernet qui est en relation avec votre modem câble ou votre routeur DSL est vain sans d'autres mises en forme du trafic !

9.2.3.1. Paramètres & usage

SFQ est presque configuré de base :

`perturb`

Reconfigure le hachage une fois toutes les `perturb` secondes. S'il n'est pas indiqué, le hachage se sera jamais reconfiguré. Ce n'est pas recommandé. 10 secondes est probablement une bonne valeur.

`quantum`

Nombre d'octets qu'un flux est autorisé à retirer de la file d'attente avant que la prochaine file d'attente ne prenne son tour. Par défaut, égal à la taille maximum d'un paquet (MTU). Ne le configurez pas en dessous du MTU !

9.2.3.2. Configuration simple

Si vous avez un périphérique qui a une vitesse identique à celle du lien et un débit réel disponible, comme un modem téléphonique, cette configuration aidera à promouvoir l'équité :

```
# tc qdisc add dev ppp0 root sfq perturb 10
# tc -s -d qdisc ls
qdisc sfq 800c: dev ppp0 quantum 1514b limit 128p flows 128/1024 perturb 10sec
Sent 4812 bytes 62 pkts (dropped 0, overlimits 0)
```

Le nombre `800c` est un descripteur (*handle*) automatiquement assigné et `limit` signifie que 128 paquets peuvent attendre dans la file d'attente. Il y a 1024 « seaux de hachage » disponibles pour la comptabilité, 128 pouvant être actifs à la fois (pas plus de paquets ne conviennent dans la file d'attente). Le hachage est reconfiguré toutes les 10 secondes.

9.3. Conseils pour le choix de la file d'attente

Pour résumer, ces files d'attente simples gèrent le trafic en réordonnant, en ralentissant ou en supprimant les paquets.

Les astuces suivantes peuvent vous aider à choisir la file d'attente à utiliser. Elles mentionnent certaines files d'attente décrites dans le chapitre *Gestionnaires de mise en file d'attente avancés*.

- Pour simplement ralentir le trafic sortant, utilisez le *Token Bucket Filter*. Il convient bien pour les énormes bandes passantes, si vous paramétrez en conséquence le seau.
- Si votre lien est vraiment saturé et que vous voulez être sûr qu'aucune session ne va accaparer la bande passante vers l'extérieur, utilisez le *Stochastic Fairness Queueing*.
- Si vous avez une grande dorsale et que vous voulez savoir ce que vous faites, considérez *Random Early Drop* (voir le chapitre *Gestionnaires de mise en file d'attente avancés*).
- Pour « mettre en forme » le trafic entrant qui n'est pas transmis, utilisez la réglementation Ingress (*Ingress Policier*). La mise en forme du flux entrant est appelée « réglementation » (*policing*) et non « mise en forme » (*shaping*).
- Si vous transmettez le trafic, utilisez TBF sur l'interface vers laquelle vous transmettez les données. Si vous voulez mettre en forme un trafic pouvant sortir par plusieurs interfaces, alors le seul facteur commun est l'interface entrante. Dans ce cas, utilisez la réglementation Ingress.
- Si vous ne voulez pas mettre en forme le trafic, mais que vous voulez voir si votre interface est tellement chargée qu'elle a dû mettre en file d'attente les données, utilisez la file d'attente `pfifo` (pas `pfifo_fast`). Elle n'a pas de bandes internes, mais assure le comptage de la taille de son accumulateur.
- Finalement, vous pouvez aussi faire de la « mise en forme sociale ». La technologie n'est pas toujours capable de réaliser ce que vous voulez. Les utilisateurs sont hostiles aux contraintes techniques. Un mot aimable peut également vous aider à avoir votre bande passante correctement divisée !

9.4. terminologie

Pour comprendre correctement des configurations plus compliquées, il est d'abord nécessaire d'expliquer quelques concepts. A cause de la complexité et de la relative jeunesse du sujet, beaucoup de mots différents sont utilisés par les personnes mais ils signifient en fait la même chose.

Ce qui suit est lâchement inspiré du texte `draft-ietf-diffserv-model-06.txt`, *An Informal Management Model for Diffserv Routers*. Il peut être trouvé à l'adresse <http://www.ietf.org/internet-drafts/draft-ietf-diffserv-model-04.txt>¹.

¹ <http://www.ietf.org/internet-drafts/draft-ietf-diffserv-model-04.txt>

Lisez-le pour les définitions strictes des termes utilisés.

Gestionnaire de mise en file d'attente (qdisc) (*Queueing Discipline*)

Un algorithme qui gère la file d'attente d'un périphérique, soit pour les données entrantes (*ingress*), soit pour les données sortantes (*egress*).

Gestionnaire de mise en file d'attente sans classes (*Classless qdisc*)

Un gestionnaire de mise en file d'attente qui n'a pas de subdivisions internes configurables.

Gestionnaire de mise en file d'attente basé sur des classes (*Classful qdisc*)

Un gestionnaire de mise en file d'attente basé sur des classes contient de multiples classes. Certaines de ces classes contiennent un gestionnaire de mise en file d'attente supplémentaire, qui peut encore être basé sur des classes, mais ce n'est pas obligatoire. Si l'on s'en tient à la définition stricte, `pfifo_fast EST` basé sur des classes, dans la mesure où il contient trois bandes, qui sont en fait des classes. Cependant, d'un point de vue des perspectives de configuration pour l'utilisateur, il est sans classes dans la mesure où ces classes ne peuvent être modifiées avec l'outil `tc`.

Classes

Un gestionnaire de mise en file d'attente basé sur les classes peut avoir beaucoup de classes, chacune d'elles étant internes au gestionnaire. Une classe peut à son tour se voir ajouter plusieurs classes. Une classe peut donc avoir comme parent soit un gestionnaire de mise en file d'attente, soit une autre classe. Une classe terminale est une classe qui ne possède de classes enfants. Seul 1 gestionnaire de mise en file d'attente est attaché à cette classe. Ce gestionnaire est responsable de l'envoi des données de cette classe. Quand vous créez une classe, un gestionnaire de mise en file d'attente fifo est créé. Quand vous ajoutez une classe enfant, ce gestionnaire est supprimé. Le gestionnaire fifo d'une classe terminale peut être remplacé par un autre gestionnaire plus adapté. Vous pouvez même remplacer ce gestionnaire fifo par un gestionnaire de mise en file d'attente basé sur des classes de sorte que vous pourrez rajouter des classes supplémentaires.

Classificateur (*Classifier*)

Chaque gestionnaire de mise en file d'attente basé sur des classes a besoin de déterminer vers quelles classes il doit envoyer un paquet. Ceci est réalisé en utilisant le classificateur.

Filtre (*Filter*)

La classification peut être réalisée en utilisant des filtres. Un filtre est composé d'un certain nombre de conditions qui, si elles sont toutes vérifiées, satisfait le filtre.

Ordonnancement (*Scheduling*)

Un gestionnaire de mise en file d'attente peut, avec l'aide d'un classificateur, décider que des paquets doivent sortir plus tôt que d'autres. Ce processus est appelé ordonnancement (*scheduling*), et est réalisé par exemple par le gestionnaire `pfifo_fast` mentionné plus tôt. L'ordonnancement est aussi appelé « reclassement » (*reordering*), ce qui peut prêter à confusion.

Mise en forme (*Shaping*)

Le processus qui consiste à retarder l'émission des paquets sortants pour avoir un trafic conforme à un débit maximum configuré. La mise en forme est réalisée sur *egress*. Familièrement, rejeter des paquets pour ralentir le trafic est également souvent appelé Mise en forme.

Réglementation (*Policing*)

Retarder ou jeter des paquets dans le but d'avoir un trafic restant en dessous d'une bande passante configurée. Dans Linux, la réglementation ne peut que jeter un paquet, et non le retarder dans la mesure où il n'y a pas de « file d'attente d'entrée » (*ingress queue*).

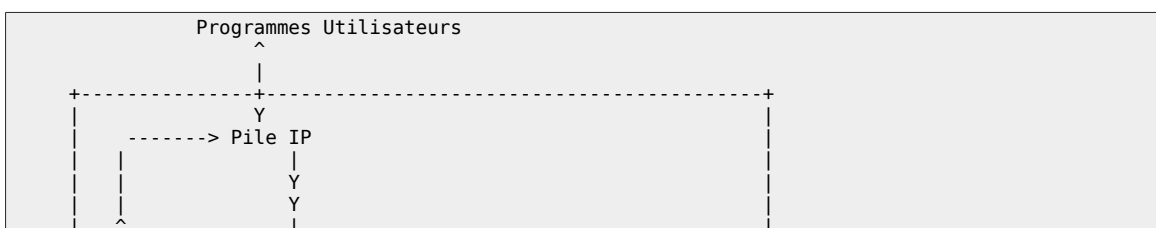
Work-Conserving

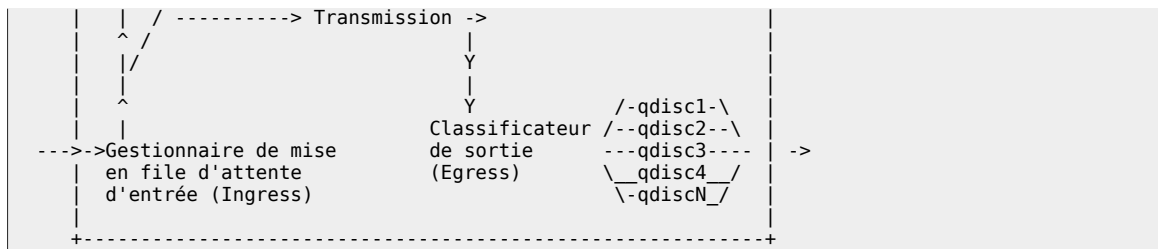
Un gestionnaire de mise en file d'attente *work-conserving* délivre toujours un paquet s'il y en a un de disponible. En d'autres termes, il ne retarde jamais un paquet si l'adaptateur réseau est prêt à l'envoyer (dans le cas du gestionnaire *egress*).

non-Work-Conserving

Quelques gestionnaire de mise en files d'attente, comme par exemple le *Token Bucket Filter*, peuvent avoir besoin de maintenir un paquet pendant un certain temps pour limiter la bande passante. Ceci signifie qu'ils refusent parfois de libérer un paquet, bien qu'ils en aient un de disponible.

Maintenant que nous avons défini notre terminologie, voyons où tous ces éléments sont situés.





Merci à Jamal Hadi Salim pour cette représentation ASCII.

Le grand rectangle représente le noyau. La flèche la plus à gauche représente le trafic du réseau entrant dans votre machine. Celui-ci alimente alors le gestionnaire de mise en file d'attente Ingress qui peut appliquer des filtres à un paquet, et décider de le supprimer. Ceci est appelé « réglementation » (*Policing*).

Ce processus a lieu très tôt, avant d'avoir beaucoup parcouru le noyau. C'est par conséquent un très bon endroit pour rejeter au plus tôt du trafic, sans pour autant consommer beaucoup de ressources CPU.

Si le paquet est autorisé à continuer, il peut être destiné à une application locale et, dans ce cas, il entre dans la couche IP pour être traité et délivré à un programme utilisateur. Le paquet peut également être transmis sans entrer dans une application et dans ce cas, être destiné à *egress*. Les programmes utilisateurs peuvent également délivrer des données, qui sont alors transmises et examinées par le classificateur *Egress*.

Là, il est examiné et mis en file d'attente vers un certain nombre de gestionnaire de mise en file d'attente. Par défaut, il n'y a qu'un seul gestionnaire *egress* installé, *pfifo_fast*, qui reçoit tous les paquets. Ceci correspond à « la mise en file d'attente » (*enqueueing*).

Le paquet réside maintenant dans le gestionnaire de mise en file d'attente, attendant que le noyau le réclame pour le transmettre à travers l'interface réseau. Ceci correspond au « retrait de la file d'attente » (*dequeueing*).

Le schéma ne montre que le cas d'un seul adaptateur réseau. Les flèches entrantes et sortantes du noyau ne doivent pas être trop prises au pied de la lettre. Chaque adaptateur réseau a un gestionnaire d'entrée et de sortie.

9.5. Gestionnaires de file d'attente basés sur les classes

Les gestionnaires de mise en file d'attente basés sur des classes sont très utiles si vous avez différentes sortes de trafic qui doivent être traités différemment. L'un d'entre eux est appelé CBQ, pour *Class Based Queueing*. Il est si souvent mentionné que les personnes identifient les gestionnaires de mise en file d'attente basés sur des classes uniquement à CBQ, ce qui n'est pas le cas.

CBQ est le mécanisme le plus ancien, ainsi que le plus compliqué. Il n'aura pas forcément les effets que vous recherchez. Ceci surprendra peut-être ceux qui sont sous l'emprise de « l'effet Sendmail », qui nous enseigne qu'une technologie complexe, non documentée est forcément meilleure que toute autre.

Nous évoquerons bientôt, plus à propos, CBQ et ses alternatives.

9.5.1. Flux à l'intérieur des gestionnaires basés sur des classes & à l'intérieur des classes

Quand le trafic entre dans un gestionnaire de mise en file d'attente basé sur des classes, il doit être envoyé vers l'une de ses classes ; il doit être « classifié ». Pour déterminer que faire d'un paquet, les éléments appelés « filtres » sont consultés. Il est important de savoir que les filtres sont appelés de l'intérieur d'un gestionnaire, et pas autrement !

Les filtres attachés à ce gestionnaire renvoient alors une décision que le gestionnaire utilise pour mettre en file d'attente le paquet vers l'une des classes. Chaque sous-classe peut essayer d'autres filtres pour voir si de nouvelles instructions s'appliquent. Si ce n'est pas le cas, la classe met le paquet en file d'attente dans le gestionnaire de mise en file d'attente qu'elle contient.

En plus de contenir d'autres gestionnaires, la plupart des gestionnaires de mise en file d'attente basés sur des classes réalisent également de la mise en forme. Ceci est utile pour réaliser à la fois l'ordonnancement (avec SFQ, par exemple) et le contrôle de débit. Vous avez besoin de ceci dans les cas où vous avez une interface à haut débit (ethernet, par exemple) connectée à un périphérique plus lent (un modem câble).

Si vous n'utilisez que SFQ, rien ne devait se passer, dans la mesure où les paquets entrent et sortent du routeur sans délai : l'interface de sortie est de loin beaucoup plus rapide que la vitesse réelle de votre liaison ; il n'y a alors pas de files d'attente à réordonnancer.

9.5.2. La famille des gestionnaires de mise en file d'attente : racines, descendants, et parents

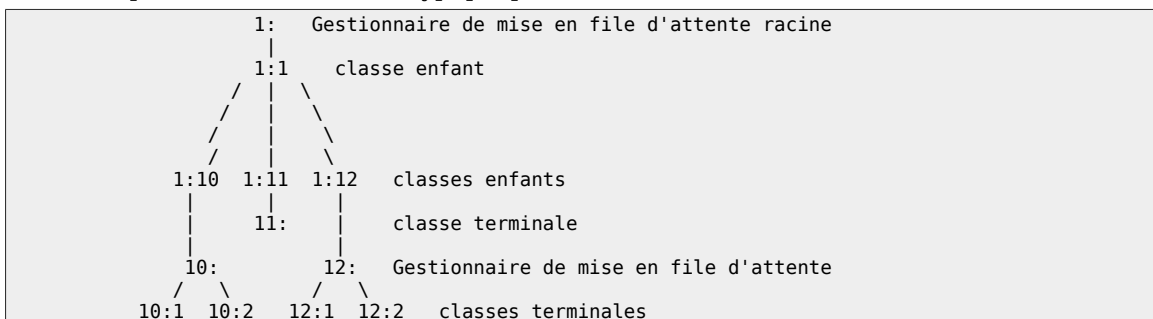
Chaque interface à « un gestionnaire de mise en file d'attente racine » de sortie (*egress root qdisc*). Par défaut, le gestionnaire de mise en file d'attente sans classes mentionné plus tôt *pfifo_fast*. Chaque gestionnaire et classe est repéré par un descripteur (*handle*), qui pourra être utilisé par les prochaines déclarations de configuration pour se référer à ce gestionnaire. En plus du gestionnaire de sortie, une interface peut également avoir un gestionnaire d'entrée (*ingress*), qui régleme le trafic entrant.

Ces descripteurs sont constitués de deux parties : un nombre majeur et un nombre mineur : <major>:<minor>. Il est habituel de nommer le gestionnaire racine 1:, ce qui est équivalent à 1:0. Le nombre mineur d'un gestionnaire de mise en file d'attente est toujours 0.

Les classes doivent avoir le même nombre majeur que leur parent. Le nombre majeur doit être unique à l'intérieur d'une configuration egress ou ingress. Le nombre mineur doit être unique à l'intérieur d'un gestionnaire de mise en file d'attente et de ses classes.

9.5.2.1. Comment les filtres sont utilisés pour classifier le trafic

Pour récapituler, une hiérarchie typique pourrait ressembler à ceci :



Mais ne laissez pas cet arbre vous abuser ! Vous ne devriez *pas* imaginer le noyau être au sommet de l'arbre et le réseau en dessous, ce qui n'est justement pas le cas. Les paquets sont mis et retirés de la file d'attente à la racine du gestionnaire, qui est le seul élément avec lequel le noyau dialogue.

Un paquet pourrait être classifié à travers une chaîne suivante :

```
1: -> 1:1 -> 12: -> 12:2
```

Le paquet réside maintenant dans la file d'attente du gestionnaire attaché à la classe 12:2. Dans cet exemple, un filtre a été attaché à chaque nœud de l'arbre, chacun choisissant la prochaine branche à prendre. Cela est réalisable. Cependant, ceci est également possible :

```
1: -> 12:2
```

Dans ce cas, un filtre attaché à la racine a décidé d'envoyer le paquet directement à 12:2.

9.5.2.2. Comment les paquets sont retirés de la file d'attente et envoyés vers le matériel

Quand le noyau décide qu'il doit extraire des paquets pour les envoyer vers l'interface, le gestionnaire racine 1: reçoit une requête de queue, qui est transmise à 1:1 et qui, à son tour, est passée à 10:, 11: et 12:, chacune interrogeant leurs descendances qui essaient de retirer les paquets de leur file d'attente. Dans ce cas, le noyau doit parcourir l'ensemble de l'arbre, car seul 12:2 contient un paquet.

En résumé, les classes « emboîtées » parlent *uniquement* à leur gestionnaire de mise en file d'attente parent ; jamais à une interface. Seul la file d'attente du gestionnaire racine est vidée par le noyau !

Ceci a pour résultat que les classes ne retirent jamais les paquets d'une file d'attente plus vite que ce que leur parent autorise. Et c'est exactement ce que nous voulons : de cette manière, nous pouvons avoir SFQ dans une classe interne qui ne fait pas de mise en forme, mais seulement de l'ordonnancement, et avoir un gestionnaire de mise en file d'attente extérieur qui met en forme le trafic.

9.5.3. Le gestionnaire de mise en file d'attente PRIO

Le gestionnaire de mise en file d'attente ne met pas vraiment en forme le trafic ; il ne fait que le subdiviser en se basant sur la manière dont vous avez configuré vos filtres. Vous pouvez considérer les gestionnaires PRIO comme une sorte de super `pfifo_fast` dopé, où chaque bande est une classe séparée au lieu d'une simple FIFO.

Quand un paquet est mis en file d'attente dans le gestionnaire PRIO, une classe est choisie en fonction des filtres que vous avez donnés. Par défaut, trois classes sont créées. Ces classes contiennent par défaut de purs gestionnaires de mise en file d'attente FIFO sans structure interne, mais vous pouvez les remplacer par n'importe quels gestionnaires disponibles.

Chaque fois qu'un paquet doit être retiré d'une file d'attente, la classe :1 est d'abord testée. Les classes plus élevées ne sont utilisées que si aucune des bandes plus faibles n'a pas fourni de paquets.

Cette file d'attente est très utile dans le cas où vous voulez donner la priorité à certains trafics en utilisant toute la puissance des filtres `tc` et en ne se limitant pas seulement aux options du champ TOS. Vous pouvez également ajouter un autre gestionnaire de mise en file d'attente aux trois classes prédéfinies, tandis que `pfifo_fast` est limité aux simples gestionnaires FIFO.

Puisqu'il ne met pas vraiment en forme, on applique le même avertissement que pour SFQ. Utilisez PRIO seulement si votre lien physique est vraiment saturé ou intégrez-le à l'intérieur d'un gestionnaire de mise en file d'attente basé sur des classes qui réalisent la mise en forme. Ce dernier cas est valable pour pratiquement tous les modems-câbles et les périphériques DSL.

En termes formels, le gestionnaire de mise en file d'attente PRIO est un ordonnanceur *Work-Conserving*.

9.5.3.1. Paramètres PRIO & usage

Les paramètres suivants sont reconnus par **tc** :

bands

Nombre de bandes à créer. Chaque bande est en fait une classe. Si vous changez ce nombre, vous devez également changer :

priomap

Si vous ne fournissez pas de filtres **tc** pour classer le trafic, le gestionnaire PRIO regarde la priorité TC_PRIO pour décider comment mettre en file d'attente le trafic.

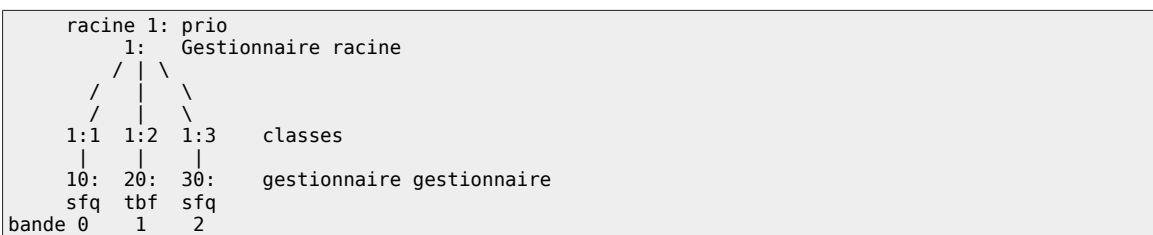
Ceci fonctionne comme le gestionnaire de mise en file d'attente `pfifo_fast` mentionné plus tôt. Voir la section correspondante pour plus de détails.

Les bandes sont des classes et sont appelées par défaut majeur:1 à majeur:3. Donc, si votre gestionnaire de mise en file d'attente est appelé 12:, **tc** filtre le trafic vers 12:1 pour lui accorder une plus grande priorité.

Par itération, la bande 0 correspond au nombre mineur 1, la bande 1 au nombre mineur 2, etc ...

9.5.3.2. Configuration simple

Nous allons créer cet arbre :



Le trafic de masse ira vers 30: tandis que le trafic interactif ira vers 20: ou 10:.

Les lignes de commande :

```

# tc qdisc add dev eth0 root handle 1: prio
## Ceci crée *instantanément* les classes 1:1, 1:2, 1:3

# tc qdisc add dev eth0 parent 1:1 handle 10: sfq
# tc qdisc add dev eth0 parent 1:2 handle 20: tbf rate 20kbit buffer 1600 limit 3000
# tc qdisc add dev eth0 parent 1:3 handle 30: sfq
  
```

Regardons maintenant ce que nous avons créé :

```

# tc -s qdisc ls dev eth0
qdisc sfq 30: quantum 1514b
  Sent 0 bytes 0 pkts (dropped 0, overlimits 0)

qdisc tbf 20: rate 20Kbit burst 1599b lat 667.6ms
  Sent 0 bytes 0 pkts (dropped 0, overlimits 0)

qdisc sfq 10: quantum 1514b
  Sent 132 bytes 2 pkts (dropped 0, overlimits 0)

qdisc prio 1: bands 3 priomap 1 2 2 2 1 2 0 0 1 1 1 1 1 1 1 1
  Sent 174 bytes 3 pkts (dropped 0, overlimits 0)
  
```

Comme vous pouvez le voir, la bande 0 a déjà reçu du trafic, et un paquet a été envoyé pendant l'exécution de cette commande !

Nous allons maintenant générer du trafic de masse avec un outil qui configure correctement les options TOS, et regarder de nouveau :

```

# scp tc ahu@10.0.0.11:./
ahu@10.0.0.11's password:
tc 100% |*****| 353 KB 00:00
# tc -s qdisc ls dev eth0
qdisc sfq 30: quantum 1514b
  Sent 384228 bytes 274 pkts (dropped 0, overlimits 0)

qdisc tbf 20: rate 20Kbit burst 1599b lat 667.6ms
  Sent 2640 bytes 20 pkts (dropped 0, overlimits 0)

qdisc sfq 10: quantum 1514b
  Sent 2230 bytes 31 pkts (dropped 0, overlimits 0)

qdisc prio 1: bands 3 priomap 1 2 2 2 1 2 0 0 1 1 1 1 1 1 1 1
  Sent 389140 bytes 326 pkts (dropped 0, overlimits 0)
  
```

Comme vous pouvez le voir, tout le trafic a été envoyé comme prévu vers le descripteur 30:, qui est la bande de plus faible priorité. Maintenant, pour vérifier que le trafic interactif va vers les bandes de plus grande priorité, nous générons du trafic interactif :

```
# tc -s qdisc ls dev eth0
qdisc sfq 30: quantum 1514b
Sent 384228 bytes 274 pkts (dropped 0, overlimits 0)

qdisc tbf 20: rate 20Kbit burst 1599b lat 667.6ms
Sent 2640 bytes 20 pkts (dropped 0, overlimits 0)

qdisc sfq 10: quantum 1514b
Sent 14926 bytes 193 pkts (dropped 0, overlimits 0)

qdisc prio 1: bands 3 priomap 1 2 2 2 1 2 0 0 1 1 1 1 1 1 1
Sent 401836 bytes 488 pkts (dropped 0, overlimits 0)
```

Ca a marché. Tout le trafic supplémentaire a été vers 10:, qui est notre gestionnaire de plus grande priorité. Aucun trafic n'a été envoyé vers les priorités les plus faibles, qui avaient reçu au préalable tout le trafic venant de notre **scp**.

9.5.4. Le célèbre gestionnaire de mise en file d'attente CBQ

Comme dit avant, CBQ est le gestionnaire de mise en file d'attente disponible le plus complexe, celui qui a eu le plus de publicité, qui est le moins compris et qui est probablement le plus farceur lors de sa mise au point. Ce n'est pas parce que les auteurs sont mauvais ou incompetents, loin de là, mais l'algorithme CBQ n'est pas remarquablement précis et il ne correspond pas vraiment à la façon dont Linux fonctionne.

En plus d'être basé sur des classes, CBQ sert également à la mise en forme de trafic et c'est sur cet aspect qu'il ne fonctionne pas très bien. Il travaille comme ceci : si vous essayez de mettre en forme une connexion de 10mbit/s à 1mbits/s, le lien doit être inactif 90% du temps. Si ce n'est pas le cas, nous devons limiter le taux de sorte qu'il *soit* inactif 90% du temps.

Ceci est assez dur à mesurer et c'est pour cette raison que CBQ déduit le temps d'inactivité du nombre de microsecondes qui s'écoulent entre les requêtes de la couche matérielle pour avoir plus de données. Cette combinaison peut être utilisée pour évaluer si le lien est chargé ou non.

Ceci est plutôt léger et l'on arrive pas toujours à des résultats convenables. Par exemple, qu'en est-il de la vitesse de liaison réelle d'une interface qui n'est pas capable de transmettre pleinement les données à 100mbit/s, peut-être à cause d'un mauvais pilote de périphérique ? Une carte réseau PCMCIA ne pourra jamais atteindre 100mbit/s à cause de la conception du bus. De nouveau, comment calculons-nous le temps d'inactivité ?

Cela devient même pire quand on considère un périphérique réseau "pas-vraiment-réel" comme *PPP Over Ethernet* ou *PPTP over TCP/IP*. La largeur de bande effective est dans ce cas probablement déterminée par l'efficacité des tubes vers l'espace utilisateur, qui est énorme.

Les personnes qui ont effectué des mesures ont découvert que CBQ n'est pas toujours très exact, et parfois même, très éloigné de la configuration.

Cependant, il marche bien dans de nombreuses circonstances. Avec la documentation fournie ici, vous devriez être capable de le configurer pour qu'il fonctionne bien dans la plupart des cas.

9.5.4.1. Mise en forme CBQ en détail

Comme dit précédemment, CBQ fonctionne en s'assurant que le lien est inactif juste assez longtemps pour abaisser la bande passante réelle au débit configuré. Pour réaliser cela, il calcule le temps qui devrait s'écouler entre des paquets de taille moyenne.

En cours de fonctionnement, le temps d'inactivité effectif (*the effective idletime*) est mesuré en utilisant l'algorithme EWMA (*Exponential Weighted Moving Average*), qui considère que les paquets récents sont exponentiellement plus nombreux que ceux passés. La charge moyenne UNIX (*UNIX loadaverage*) est calculée de la même manière.

Le temps d'inactivité calculé est soustrait à celui mesuré par EWMA et le nombre résultant est appelé *avgidle*. Un lien parfaitement chargé a un *avgidle* nul : un paquet arrive à chaque intervalle calculé.

Une liaison surchargée a un *avgidle* négatif et s'il devient trop négatif, CBQ s'arrête un moment et se place alors en dépassement de limite (*overlimit*).

Inversement, un lien inutilisé peut accumuler un *avgidle* énorme, qui autoriserait alors des bandes passantes infinies après quelques heures d'inactivité. Pour éviter cela, *avgidle* est borné à *maxidle*.

En situation de dépassement de limite, CBQ peut en théorie bloquer le débit pour une durée équivalente au temps qui doit s'écouler entre deux paquets moyens, puis laisser passer un paquet et bloquer de nouveau le débit. Regardez cependant le paramètre *minburst* ci-dessous.

Voici les paramètres que vous pouvez spécifier pour configurer la mise en forme :

avpkt

Taille moyenne d'un paquet mesurée en octets. Nécessaire pour calculer *maxidle*, qui dérive de *maxburst*, qui est spécifié en paquets.

bandwidth

La bande passante physique de votre périphérique nécessaire pour les calculs du temps de non utilisation (*idle time*).

cell

La durée de transmission d'un paquet n'augmente pas nécessairement de manière linéaire en fonction de sa taille. Par exemple, un paquet de 800 octets peut être transmis en exactement autant de temps qu'un paquet de 806 octets. Ceci détermine la granularité. Cette valeur est généralement positionnée à 8, et doit être une puissance de deux.

maxburst

Ce nombre de paquets est utilisé pour calculer `maxidle` de telle sorte que quand `avgidle` est égal à `maxidle`, le nombre de paquets moyens peut être envoyé en rafale avant que `avgidle` ne retombe à 0. Augmentez-le pour être plus tolérant vis à vis des rafales de données. Vous ne pouvez pas configurer `maxidle` directement, mais seulement via ce paramètre.

minburst

Comme nous l'avons déjà indiqué, CBQ doit bloquer le débit dans le cas d'un dépassement de limite. La solution idéale est de le faire pendant exactement le temps d'inutilisation calculé, puis de laisser passer un paquet. Cependant, les noyaux UNIX ont généralement du mal à prévoir des événements plus courts que 10 ms, il vaut donc mieux limiter le débit pendant une période plus longue, puis envoyer `minburst` paquets d'un seul coup et dormir pendant une durée de `minburst`.

Le temps d'attente est appelé *offtime*. De plus grandes valeurs de `minburst` mènent à une mise en forme plus précise dans le long terme, mais provoquent de plus grandes rafales de données pendant des périodes de quelques millisecondes.

minidle

Si `avgidle` est inférieur à 0, nous sommes en dépassement de limite et nous devons attendre jusqu'à ce que `avgidle` devienne suffisamment important pour envoyer un paquet. Pour éviter qu'une brusque rafale de données n'empêche le lien de fonctionner pendant une durée prolongée, `avgidle` est remis à `minidle` s'il atteint une valeur trop basse.

La valeur `minidle` est spécifiée en microsecondes négatives : 10 signifie alors que `avgidle` est borné à -10µs.

mpu

Taille minimum d'un paquet. Nécessaire car même un paquet de taille nulle est encapsulé par 64 octets sur ethernet et il faut donc un certain temps pour le transmettre. CBQ doit connaître ce paramètre pour calculer précisément le temps d'inutilisation.

rate

Débit du trafic sortant du gestionnaire. Ceci est le « paramètre de vitesse » !

En interne, CBQ est finement optimisé. Par exemple, les classes qui sont connues pour ne pas avoir de données présentes dans leur file d'attente ne sont pas interrogées. Les classes en situation de dépassement de limite sont pénalisées par la diminution de leur priorité effective. Tout ceci est très habile et compliqué.

9.5.4.2. Le comportement CBQ classful

En plus de la mise en forme, en utilisant les approximations `idletime` mentionnées ci-dessus, CBQ peut également agir comme une file d'attente PRIO dans le sens où les classes peuvent avoir différentes priorités. Les priorités de plus faible valeur seront examinées avant celles de valeurs plus élevées.

Chaque fois qu'un paquet est demandé par la couche matérielle pour être envoyé sur le réseau, un processus *weighted round robin* (WRR) démarre en commençant par les classes de plus faibles numéros.

Celles-ci sont regroupées et interrogées si elles ont des données disponibles. Après qu'une classe ait été autorisée à retirer de la file d'attente un nombre d'octets, la classe de priorité suivante est consultée.

Les paramètres suivants contrôlent le processus WRR :

allot

Quand le CBQ racine reçoit une demande d'envoi de paquets sur une interface, il va essayer tous les gestionnaires internes (dans les classes) tour à tour suivant l'ordre du paramètre `priority`. A chaque passage, une classe ne peut envoyer qu'une quantité limitée de données. Le paramètre `allot` est l'unité de base de cette quantité. Voir le paramètre `weight` pour plus d'informations.

prio

CBQ peut également agir comme un périphérique PRIO. Les classes internes avec les priorités les plus élevées sont consultées en premier et, aussi longtemps qu'elles ont du trafic, les autres classes ne sont pas examinées.

weight

Le paramètre `weight` assiste le processus *Weighted Round Robin*. Chaque classe a tour à tour la possibilité d'envoyer ses données. Si vous avez des classes avec des bandes passantes significativement plus importantes, il est logique de les autoriser à envoyer plus de données à chaque tour que les autres.

Vous pouvez utiliser des nombres arbitraires dans la mesure où CBQ additionne tous les paramètres `weight` présents sous une classe et les normalise. La règle empirique qui consiste à prendre `rate/10` semble fonctionner correctement. Le paramètre `weight` normalisé est multiplié par le paramètre `allot` pour déterminer la quantité de données à envoyer à chaque tour.

Notez, s'il vous plaît, que toutes les classes à l'intérieur d'une hiérarchie CBQ doivent avoir le même nombre majeur !

9.5.4.3. Paramètres CBQ qui déterminent le partage & le prêt du lien

En plus de purement limiter certains trafics, il est également possible de spécifier quelles classes peuvent emprunter de la bande passante aux autres classes ou, réciproquement, prêter sa bande passante.

`isolated/ sharing`

Une classe qui est configurée avec `isolated` ne prêtera pas sa bande passante à ses classes soeurs. Utilisez ceci si vous avez sur votre lien deux agences concurrentes ou qui ne s'apprécient pas et qui ne veulent pas se prêter gratuitement de la bande passante.

Le programme de contrôle `tc` connaît également `sharing`, qui agit à l'inverse du paramètre `isolated`.

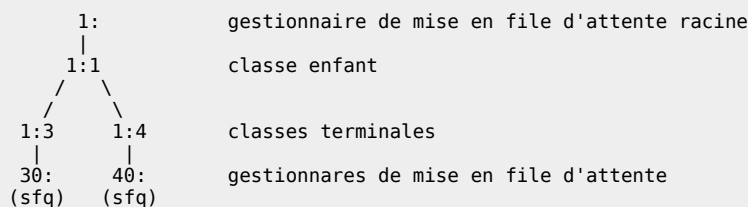
`bounded/ borrow`

Une classe peut aussi être bornée (`bounded`), ce qui signifie qu'elle n'essaiera pas d'emprunter de la bande passante à ses classes enfants. `tc` connaît également `borrow`, qui agit à l'inverse de `bounded`.

Une situation typique pourrait être le cas où vous avez deux agences présentes sur votre lien qui sont à la fois `isolated` et `bounded`. Ceci signifie qu'elles sont strictement limitées à leur débit et qu'elles ne prêteront pas aux autres leur bande passante.

A l'intérieur de ces classes d'agence, il pourrait y avoir d'autres classes qui sont autorisées à échanger leur bande passante.

9.5.4.4. Configuration simple



Cette configuration limite le trafic d'un serveur web à 5 mbit et le trafic SMTP à 3 mbit. Il est souhaitable qu'ils n'occupent pas plus de 6 mbit à eux deux. Nous avons une carte réseau à 100 mbit et les classes peuvent s'emprunter mutuellement de la bande passante.

```
# tc qdisc add dev eth0 root handle 1:0 cbq bandwidth 100Mbit \
  avpkt 1000 cell 8
# tc class add dev eth0 parent 1:0 classid 1:1 cbq bandwidth 100Mbit \
  rate 6Mbit weight 0.6Mbit prio 8 allot 1514 cell 8 maxburst 20 \
  avpkt 1000 bounded
```

Cette partie installe la racine et la classe 1:1 habituelle. La classe 1:1 est bornée, la bande passante totale ne pourra donc pas excéder 6 mbit.

Comme dit avant, CBQ a besoin de *NOMBREUX* paramètres. Tous ces paramètres sont cependant expliqués au-dessus. La configuration HTB correspondante est beaucoup plus simple.

```
# tc class add dev eth0 parent 1:1 classid 1:3 cbq bandwidth 100Mbit \
  rate 5Mbit weight 0.5Mbit prio 5 allot 1514 cell 8 maxburst 20 \
  avpkt 1000
# tc class add dev eth0 parent 1:1 classid 1:4 cbq bandwidth 100Mbit \
  rate 3Mbit weight 0.3Mbit prio 5 allot 1514 cell 8 maxburst 20 \
  avpkt 1000
```

Ce sont nos deux classes. Notez comment nous avons configuré la valeur du paramètre `weight` en fonction du paramètre `rate`. La bande passante de l'ensemble des deux classes ne pourra jamais dépasser 6 mbit. En fait, les identifiants de classe (`classid`) doivent avoir le même numéro majeur que le gestionnaire de mise en file d'attente parent !

```
# tc qdisc add dev eth0 parent 1:3 handle 30: sfq
# tc qdisc add dev eth0 parent 1:4 handle 40: sfq
```

Les deux classes ont par défaut un gestionnaire de mise en file d'attente FIFO. Nous les remplaçons par une file d'attente SFQ de telle sorte que chaque flux de données soit traité de manière égale.

```
# tc filter add dev eth0 parent 1:0 protocol ip prio 1 u32 match ip \
  sport 80 0xffff flowid 1:3
# tc filter add dev eth0 parent 1:0 protocol ip prio 1 u32 match ip \
  sport 25 0xffff flowid 1:4
```

Ces commandes, directement attachées à la racine, envoient le trafic vers le bon gestionnaire de mise en file d'attente.

Notez que nous utilisons `tc class add` pour *CREER* des classes à l'intérieur d'un gestionnaire de mise en file d'attente, et que nous utilisons `tc qdisc add` pour véritablement configurer ces classes.

Vous vous demandez peut-être ce qui arrive au trafic qui n'est classifié par aucune des deux règles. Dans ce cas, les données seront traitées à l'intérieur de 1:0, et le débit ne sera pas limité.

Si le trafic SMTP+web tente de dépasser la limite de 6 mbit/s, la bande passante sera divisée selon le paramètre `weight`, donnant 5/8 du trafic au serveur web et 3/8 au serveur smtp.

Avec cette configuration, vous pouvez également dire que le trafic du serveur web sera au minimum de $5/8 * 6 \text{ mbit} = 3.75 \text{ mbit}$.

9.5.4.5. D'autres paramètres CBQ : split & defmap

Comme précisé avant, un gestionnaire de mise en file d'attente basé sur des classes doit appeler des filtres pour déterminer dans quelle classe un paquet sera mis en file d'attente.

En plus d'appeler les filtres, CBQ offre d'autres options : `defmap` & `split`. C'est plutôt compliqué à comprendre et, de plus, ce n'est pas vital. Mais, étant donné que ceci est le seul endroit connu où `defmap` & `split` sont correctement expliqués, je vais faire de mon mieux.

Étant donné que nous voulons le plus souvent réaliser le filtrage en ne considérant que le champ TOS, une syntaxe spéciale est fournie. Chaque fois que CBQ doit trouver où le paquet doit être mis en file d'attente, il vérifie si le nœud est un nœud d'aiguillage (*split node*). Si c'est le cas, un de ses sous-gestionnaires a indiqué son souhait de recevoir tous les paquets configurés avec une certaine priorité. Celle-ci peut être dérivée du champ TOS ou des options des sockets positionnées par les applications.

Les bits de priorités des paquets subissent un ET logique avec le champ `defmap` pour voir si une correspondance existe. En d'autres termes, c'est un moyen pratique de créer un filtre très rapide, qui ne sera actif que pour certaines priorités. Un `defmap` de `ff` (en hexadécimal) vérifiera tout tandis qu'une valeur de `0` ne vérifiera rien. Une configuration simple aidera peut-être à rendre les choses plus claires :

```
# tc qdisc add dev eth1 root handle 1: cbq bandwidth 10Mbit allot 1514 \
  cell 8 avpkt 1000 mpu 64

# tc class add dev eth1 parent 1:0 classid 1:1 cbq bandwidth 10Mbit \
  rate 10Mbit allot 1514 cell 8 weight 1Mbit prio 8 maxburst 20 \
  avpkt 1000
```

Préambule standard de CBQ. Je n'ai jamais pris l'habitude de la quantité de nombres nécessaires !

Le paramètre `defmap` se réfère aux bits `TC_PRIO` qui sont définis comme suit :

TC_PRIO..	Num	Correspond à TOS
BESTEFFORT	0	Maximalise la Fiabilité
FILLER	1	Minimalise le Coût
BULK	2	Maximalise le Débit (0x8)
INTERACTIVE_BULK	4	
INTERACTIVE	6	Minimise le Délai (0x10)
CONTROL	7	

Les nombres `TC_PRIO..` correspondent aux bits comptés à partir de la droite. Voir la section `pfifo_fast` pour plus de détails sur la façon dont les bits TOS sont convertis en priorités.

Maintenant, les classes interactive et de masse :

```
# tc class add dev eth1 parent 1:1 classid 1:2 cbq bandwidth 10Mbit \
  rate 1Mbit allot 1514 cell 8 weight 100Kbit prio 3 maxburst 20 \
  avpkt 1000 split 1:0 defmap c0

# tc class add dev eth1 parent 1:1 classid 1:3 cbq bandwidth 10Mbit \
  rate 8Mbit allot 1514 cell 8 weight 800Kbit prio 7 maxburst 20 \
  avpkt 1000 split 1:0 defmap 3f
```

La gestion de mise en file d'attente d'aiguillage (*split qdisc*) est `1:0` et c'est à ce niveau que le choix sera fait. `c0` correspond au nombre binaire `11000000` et `3f` au nombre binaire `00111111`. Ces valeurs sont choisies de telle sorte qu'à elles deux, elles vérifient tous les bits. La première classe correspond aux bits 6 & 7, ce qui est équivalent aux trafics « interactif » et de « contrôle ». La seconde classe correspond au reste.

Le nœud `1:0` possède maintenant la table suivante :

```
priorité envoyer à
0 1:3
1 1:3
2 1:3
3 1:3
4 1:3
5 1:3
6 1:2
7 1:2
```

Pour d'autres amusements, vous pouvez également donner un « masque de changement » qui indique exactement les priorités que vous souhaitez changer. N'utilisez ceci qu'avec la commande `tc class change`. Par exemple, pour ajouter le trafic *best effort* à la classe `1:2`, nous devons exécuter ceci :

```
# tc class change dev eth1 classid 1:2 cbq defmap 01/01
```

La carte des priorités au niveau de `1:0` ressemble maintenant à ceci :

```
priorité envoyer à
0 1:2
1 1:3
2 1:3
```

```
3 1:3
4 1:3
5 1:3
6 1:2
7 1:2
```

FIXME: `tc class change` n'a pas été testé, mais simplement vu dans les sources.

9.5.5. Seau de jetons à contrôle hiérarchique (*Hierarchical Token Bucket*)

Martin Devera (<devik>) réalisa à juste titre que CBQ est complexe et qu'il ne semble pas optimisé pour de nombreuses situations classiques. Son approche hiérarchique est bien adaptée dans le cas de configurations où il y a une largeur de bande passante fixée à diviser entre différents éléments. Chacun de ces éléments aura une bande passante garantie, avec la possibilité de spécifier la quantité de bande passante qui pourra être empruntée.

HTB travaille juste comme CBQ, mais il n'a pas recourt à des calculs de temps d'inoccupation pour la mise en forme. A la place, c'est un *Token Bucket Filter* basé sur des classes, d'où son nom. Il n'a que quelques paramètres, qui sont bien documentés sur ce [site](#)².

Au fur et à mesure que votre configuration HTB se complexifie, votre configuration s'adapte bien. Avec CBQ, elle est déjà complexe même dans les cas simples ! HTB3 (voir [sa page principale](#)³ pour les détails des versions HTB) fait maintenant parti des sources officielles du noyau (à partir des versions 2.4.20-pre1 et 2.5.31 et supérieures). Il est encore cependant possible que vous soyez obligé de récupérer la version mise à jour de 'tc' pour HTB3. Les programmes de l'espace utilisateur et la partie HTB du noyau doivent avoir le même numéro majeur. Sans cela, 'tc' ne marchera pas avec HTB.

Si vous avez déjà un noyau récent ou si vous êtes sur le point de mettre à jour votre noyau, considérez HTB coûte que coûte.

9.5.5.1. Configuration simple

Fonctionnellement presque identique à la configuration simple CBQ présentée ci-dessus :

```
# tc qdisc add dev eth0 root handle 1: htb default 30
# tc class add dev eth0 parent 1: classid 1:1 htb rate 6mbit burst 15k
# tc class add dev eth0 parent 1:1 classid 1:10 htb rate 5mbit burst 15k
# tc class add dev eth0 parent 1:1 classid 1:20 htb rate 3mbit ceil 6mbit burst 15k
# tc class add dev eth0 parent 1:1 classid 1:30 htb rate 1kbit ceil 6mbit burst 15k
```

L'auteur recommande SFQ sous ces classes :

```
# tc qdisc add dev eth0 parent 1:10 handle 10: sfq perturb 10
# tc qdisc add dev eth0 parent 1:20 handle 20: sfq perturb 10
# tc qdisc add dev eth0 parent 1:30 handle 30: sfq perturb 10
```

Ajouter les filtres qui dirigent le trafic vers les bonnes classes :

```
# U32="tc filter add dev eth0 protocol ip parent 1:0 prio 1 u32"
# $U32 match ip dport 80 0xffff flowid 1:10
# $U32 match ip sport 25 0xffff flowid 1:20
```

Et, c'est tout. Pas de vilains nombres non expliqués, pas de paramètres non documentés.

HTB semble vraiment merveilleux. Si 10: et 20: ont atteint tous les deux leur bande passante garantie et qu'il en reste à partager, ils l'empruntent avec un rapport de 5:3, comme attendu.

Le trafic non classifié est acheminé vers 30:, qui a une petite bande passante, mais qui peut emprunter tout ce qui est laissé libre. Puisque nous avons choisi SFQ en interne, on hérite naturellement de l'équité.

9.6. Classifier des paquets avec des filtres

Pour déterminer quelle classe traitera un paquet, la « chaîne de classificateurs » est appelée chaque fois qu'un choix a besoin d'être fait. Cette chaîne est constituée de tous les filtres attachés aux gestionnaires de mise en file d'attente basés sur des classes qui doivent prendre une décision.

On reprend l'arbre qui n'est pas un arbre :



Quand un paquet est mis en file d'attente, l'instruction appropriée de la chaîne de filtre est consultée à chaque branche. Une configuration typique devrait avoir un filtre en 1:1 qui dirige le paquet vers 12: et un filtre en 12: qui l'envoie vers 12:2.

² <http://luxik.cdi.cz/~devik/qos/htb/>

³ <http://luxik.cdi.cz/~devik/qos/htb/>

Vous pourriez également avoir ce dernier filtre en 1:1, mais vous pouvez gagner en efficacité en ayant des tests plus spécifiques plus bas dans la chaîne.

A ce propos, vous ne pouvez pas filtrer un paquet « vers le haut ». Donc, avec HTB, vous devrez attacher tous les filtres à la racine !

Encore une fois, les paquets ne sont mis en file d'attente que vers le bas ! Quand ils sont retirés de la file d'attente, ils montent de nouveau, vers l'interface. Ils ne tombent PAS vers l'extrémité de l'arbre en direction de l'adaptateur réseau !

9.6.1. Quelques exemples simples de filtrage

Comme expliqué dans le chapitre *Filtres avancés pour la classification des paquets*, vous pouvez vraiment analyser n'importe quoi en utilisant une syntaxe très compliquée. Pour commencer, nous allons montrer comment réaliser les choses évidentes, ce qui heureusement est plutôt facile.

Disons que nous avons un gestionnaire de mise en file d'attente PRIO appelé 10: qui contient trois classes, et que nous voulons assigner à la bande de plus haute priorité tout le trafic allant et venant du port 22. Les filtres seraient les suivants :

```
# tc filter add dev eth0 protocol ip parent 10: prio 1 u32 match \
ip dport 22 0xffff flowid 10:1
# tc filter add dev eth0 protocol ip parent 10: prio 1 u32 match \
ip sport 80 0xffff flowid 10:1
# tc filter add dev eth0 protocol ip parent 10: prio 2 flowid 10:2
```

Qu'est-ce que cela signifie ? Cela dit : attacher à eth0, au nœud 10: un filtre u32 de priorité 1 qui analyse le port de destination ip 22 et qui l'envoie vers la bande 10:1. La même chose est répétée avec le port source 80. La dernière commande indique que si aucune correspondance n'est trouvée, alors le trafic devra aller vers la bande 10:2, la plus grande priorité suivante.

Vous devez ajouter eth0 ou n'importe laquelle de vos interfaces, car chaque interface possède un espace de nommage de ses descripteurs qui lui est propre.

Pour sélectionner une adresse IP, utilisez ceci :

```
# tc filter add dev eth0 parent 10:0 protocol ip prio 1 u32 \
match ip dst 4.3.2.1/32 flowid 10:1
# tc filter add dev eth0 parent 10:0 protocol ip prio 1 u32 \
match ip src 1.2.3.4/32 flowid 10:1
# tc filter add dev eth0 protocol ip parent 10: prio 2 \
flowid 10:2
```

Ceci dirige le trafic allant vers 4.3.2.1 et venant de 1.2.3.4 vers la file d'attente de plus haute priorité, tandis que le reste ira vers la prochaine plus haute priorité.

Vous pouvez rassembler ces deux vérifications pour récupérer le trafic venant de 1.2.3.4 avec le port source 80 :

```
# tc filter add dev eth0 parent 10:0 protocol ip prio 1 u32 match ip src 4.3.2.1/32
match ip sport 80 0xffff flowid 10:1
```

9.6.2. Toutes les commandes de filtres dont vous aurez normalement besoin

La plupart des commandes présentées ici commencent avec le préambule suivant :

```
# tc filter add dev eth0 parent 1:0 protocol ip prio 1 u32 ..
```

Ils sont appelés filtres u32 et analysent *N'IMPORTE QUELLE* partie d'un paquet.

Sur l'adresse source/destination

Masque pour la source `match ip src 1.2.3.0/24` et masque pour la destination `match ip dst 4.3.2.0/24`. Pour analyser un hôte simple, employez /32 ou omettez le masque.

Sur le port source/destination, tous les protocoles IP

Source: `match ip sport 80 0xffff` et destination : `match ip dport ?? 0xffff`

Sur le protocole ip (tcp, udp, icmp, gre, ipsec)

Utilisez les nombres définis dans `/etc/protocols`, par exemple 1 pour icmp : `match ip protocol 1 0xff`.

Sur fwmark

Vous pouvez marquer les paquets avec ipchains ou iptables et voir cette marque préservée lors du routage à travers les interfaces. Ceci est vraiment utile pour mettre uniquement en forme le trafic sur eth1 et venant de eth0, par exemple. La syntaxe est la suivante :

```
# tc filter add dev eth1 protocol ip parent 1:0 prio 1 handle 6 fw flowid 1:1
```

Notez que ce n'est pas une correspondance u32 !

Vous pouvez positionner une marque comme ceci :

```
# iptables -A PREROUTING -t mangle -i eth0 -j MARK --set-mark 6
```

Le nombre 6 est arbitraire.

Si vous ne voulez pas assimiler la syntaxe complète de **tc filter**, utilisez juste **iptables** et apprenez seulement la sélection basée sur fwmark.

Sur le champ TOS

Pour sélectionner le trafic interactif, délai minimum :

```
# tc filter add dev ppp0 parent 1:0 protocol ip prio 10 u32 \  
  match ip tos 0x10 0xff \  
  flowid 1:4
```

Utilisez 0x08 0xff pour le trafic de masse.

Pour plus de commandes de filtrage, voir le chapitre *Filtres avancés pour la classification des paquets*.

9.7. Le périphérique de file d'attente intermédiaire (The Intermediate queueing device (IMQ))

Le périphérique IMQ n'est pas un gestionnaire de mise en file d'attente mais son utilisation est fortement liée à ceux-ci. Au coeur de Linux, les gestionnaires de mise en file d'attente sont attachés aux périphériques réseaux et tout ce qui est mis en file d'attente dans ce périphérique l'est d'abord dans le gestionnaire. Avec ce concept, il existe deux limitations :

1. Seule la mise en forme du trafic sortant est possible (un gestionnaire d'entrée existe, mais ses possibilités sont très limitées en comparaison des gestionnaires de mise en file basés sur les classes).
2. Un gestionnaire de mise en file d'attente ne voit le trafic que d'une interface, et des limitations globales ne peuvent pas être mises en place.

IMQ est ici pour aider à résoudre ces deux limitations. En résumé, vous pouvez mettre tout ce que vous voulez dans un gestionnaire de mise en file d'attente. Les paquets spécialement marqués sont interceptés par les points d'accroche netfilter NF_IP_PRE_ROUTING et NF_IP_POST_ROUTING et sont transférés vers le gestionnaire attaché au périphérique imq. Une cible iptables est utilisée pour le marquage des paquets.

Ceci vous permet de réaliser de la mise en forme d'entrée étant donné que vous pouvez marquer les paquets entrant par un périphérique quelconque et/ou traiter les interfaces comme des classes pour configurer des limites globales. Vous pouvez également réaliser de nombreuses autres choses comme simplement mettre votre trafic http dans un gestionnaire, mettre les requêtes de nouvelles connexions dans un gestionnaire, ...

9.7.1. Configuration simple

La première chose qui devrait vous venir à l'esprit est d'utiliser la mise en forme du trafic entrant pour vous garantir une grande passante. ;) La configuration se fait comme avec n'importe quelle autre interface :

```
tc qdisc add dev imq0 root handle 1: htb default 20  
  
tc class add dev imq0 parent 1: classid 1:1 htb rate 2mbit burst 15k  
  
tc class add dev imq0 parent 1:1 classid 1:10 htb rate 1mbit  
tc class add dev imq0 parent 1:1 classid 1:20 htb rate 1mbit  
  
tc qdisc add dev imq0 parent 1:10 handle 10: pfifo  
tc qdisc add dev imq0 parent 1:20 handle 20: sfq  
  
tc filter add dev imq0 parent 10:0 protocol ip prio 1 u32 match \  
  ip dst 10.0.0.230/32 flowid 1:10
```

Dans cet exemple, u32 est utilisé pour la classification. Les autres classificateurs devraient marcher tout aussi bien. Le trafic doit ensuite être sélectionné et marqué pour être mis en file d'attente vers imq0.

```
iptables -t mangle -A PREROUTING -i eth0 -j IMQ --todev 0  
  
ip link set imq0 up
```

Les cibles iptables IMQ sont valides dans les chaînes PREROUTING et POSTROUTING de la table mangle. La syntaxe est la suivante :

```
IMQ [ --todev n ] n : numéro du périphérique imq
```

Il existe aussi une cible ip6tables.

Notez que le trafic n'est pas mis en file d'attente quand la cible est activée, mais après. La localisation exacte de l'entrée du trafic dans le périphérique imq dépend de la direction de ce trafic (entrant/sortant). Ces entrées sont les points d'accroche prédéfinis de netfilter et utilisés par iptables :

```
enum nf_ip_hook_priorities {  
  NF_IP_PRI_FIRST = INT_MIN,  
  NF_IP_PRI_CONNTRACK = -200,  
  NF_IP_PRI_MANGLE = -150,  
  NF_IP_PRI_NAT_DST = -100,  
  NF_IP_PRI_FILTER = 0,  
  NF_IP_PRI_NAT_SRC = 100,  
  NF_IP_PRI_LAST = INT_MAX,  
};
```

Pour le trafic entrant, imq se déclare avec la priorité NF_IP_PRI_MANGLE + 1, ce qui signifie que les paquets entrent dans le périphérique imq juste après la chaîne PREROUTING de la table mangle.

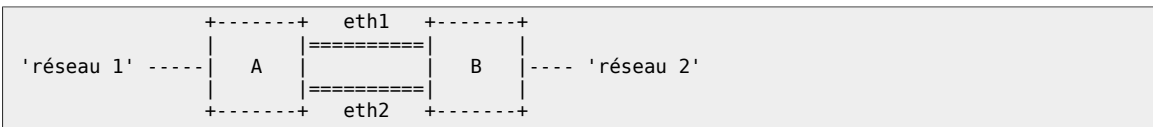
Pour le trafic sortant, `imq` utilise `NF_IP_PRI_LAST` qui honore le fait que les paquets rejetés par la table `filter` n'occuperont pas de bande passante.

Les mises à jour et de plus amples informations peuvent être trouvées sur le [site `imq`](http://luxik.cdi.cz/~patrick/imq/)⁴.

⁴ <http://luxik.cdi.cz/~patrick/imq/>

Il existe plusieurs manières pour le faire. Une des plus faciles et des plus directes est TEQL (*True (or Trivial) Link Equalizer*). Comme la plupart des éléments en relation avec la gestion de file d'attente, l'équilibrage de charge est bidirectionnel. Les deux équipements terminaux du lien ont besoin de participer pour obtenir une efficacité optimale.

Imaginez la situation suivante :



A et B sont des routeurs dont nous supposons qu'ils fonctionnent avec Linux pour le moment. Si le trafic va du réseau 1 vers le réseau 2, le routeur A a besoin de distribuer les paquets sur les deux liens allant vers B. Le routeur B a besoin d'être configuré pour l'accepter. On retrouve la même chose dans le sens inverse, pour les paquets allant du réseau 2 vers le réseau 1. Le routeur B a besoin d'envoyer les paquets à la fois sur eth1 et eth2.

La répartition est faite par un périphérique TEQL, comme ceci (cela ne pouvait pas être plus simple) :

```
# tc qdisc add dev eth1 root teql0
# tc qdisc add dev eth2 root teql0
# ip link set dev teql0 up
```

N'oubliez pas la commande **ip link set up** !

Ceci a besoin d'être fait sur les deux hôtes. Le périphérique `teql0` est basiquement un distributeur tourniquet au-dessus de `eth1` et `eth2` pour l'envoi des paquets. Aucune donnée n'arrive jamais à travers un périphérique `teql`, mais les données apparaissent sur `eth1` et `eth2`.

Nous n'avons pour le moment que les périphériques et nous avons également besoin d'un routage correct. L'une des possibilités pour réaliser cela est d'assigner un réseau `/31` sur chacun des liens, ainsi que sur le périphérique `teql0` :

FIXME: Avons nous besoin de quelque chose comme `nobroadcast` ? Un `/31` est trop petit pour contenir une adresse réseau et une adresse de diffusion. Si cela ne marche pas comme prévu, essayez un `/30`, et ajustez les adresses IP. Vous pouvez même essayer sans attribuer d'adresses à `eth1` et `eth2`.

Sur le routeur A:

```
# ip addr add dev eth1 10.0.0.0/31
# ip addr add dev eth2 10.0.0.2/31
# ip addr add dev teql0 10.0.0.4/31
```

Sur le routeur B:

```
# ip addr add dev eth1 10.0.0.1/31
# ip addr add dev eth2 10.0.0.3/31
# ip addr add dev teql0 10.0.0.5/31
```

Le routeur A devrait maintenant être capable de lancer un **ping** vers `10.0.0.1`, `10.0.0.3` et `10.0.0.5` à travers les deux liens physiques et le périphérique « égalisé ». Le routeur B devrait maintenant être capable de lancer un **ping** vers `10.0.0.0`, `10.0.0.2` et `10.0.0.4` à travers les liens.

Si cela marche, le routeur A peut prendre `10.0.0.5` comme route vers le réseau 2 et le routeur B `10.0.0.4` comme route vers le réseau 1. Pour le cas particulier où le réseau 1 est votre réseau personnel et où le réseau 2 est l'Internet, le routeur A peut prendre `10.0.0.5` comme passerelle par défaut.

10.1. Avertissement

Rien n'est aussi simple qu'il y paraît. Les interfaces `eth1` et `eth2` sur les deux routeurs A et B ne doivent pas avoir la fonction de filtrage par chemin inverse activée. Dans le cas contraire, ils rejettent les paquets destinés à des adresses autres que les leurs :

```
# echo 0 > /proc/sys/net/ipv4/conf/eth1/rp_filter
# echo 0 > /proc/sys/net/ipv4/conf/eth2/rp_filter
```

Il y a un sérieux problème avec le réordonnement des paquets. Supposons que six paquets aient besoin d'être envoyés de A vers B. Par exemple, `eth1` peut traiter les paquets 1, 3 et 5 et `eth2` les paquets 2, 4 et 6. Dans un monde idéal, le routeur B devrait recevoir ces paquets dans l'ordre 1, 2, 3, 4, 5, 6. Mais il est plus probable que le noyau les recevra comme ceci : 2, 1, 4, 3, 6, 5. Ce problème va perturber TCP/IP. Alors qu'il n'y a pas de problèmes pour les liens transportant différentes sessions TCP/IP, vous ne serez pas capable de regrouper plusieurs liens et obtenir par ftp un simple fichier beaucoup plus rapidement, à moins que le système d'exploitation envoyant ou recevant ne soit Linux. En effet, celui-ci n'est pas facilement perturbé par de simples réordonnements.

Cependant, l'équilibrage de charge est une bonne idée pour de nombreuses applications.

Jusqu'à maintenant, nous avons vu comment iproute travaille, et netfilter a été mentionné plusieurs fois. Vous ne perdrez pas votre temps à consulter [The netfilter/iptables HOWTO's](#)¹. Le logiciel Netfilter peut être trouvé [ici](#)².

Netfilter nous permet de filtrer les paquets ou de désosser leurs en-têtes. Une de ses fonctionnalités particulières est de pouvoir marquer un paquet avec un nombre, grâce à l'option `--set-mark`.

Comme exemple, la commande suivante marque tous les paquets destinés au port 25, en l'occurrence le courrier sortant.

```
# iptables -A PREROUTING -i eth0 -t mangle -p tcp --dport 25 \
-j MARK --set-mark 1
```

Disons que nous avons plusieurs connexions, une qui est rapide (et chère au mégaoctet) et une qui est plus lente, mais avec un tarif moins élevé. Nous souhaiterions que le courrier passe par la route la moins chère.

Nous avons déjà marqué le paquet avec un "1" et nous allons maintenant renseigner la base de données de la politique de routage pour qu'elle agisse sur ces paquets marqués.

```
# echo 201 mail.out >> /etc/iproute2/rt_tables
# ip rule add fwmark 1 table mail.out
# ip rule ls
0: from all lookup local
32764: from all fwmark 1 lookup mail.out
32766: from all lookup main
32767: from all lookup default
```

Nous allons maintenant générer la table mail.out avec une route vers la ligne lente, mais peu coûteuse.

```
# /sbin/ip route add default via 195.96.98.253 dev ppp0 table mail.out
```

Voilà qui est fait. Il se peut que nous voulions mettre en place des exceptions, et il existe de nombreux moyens pour le faire. Nous pouvons modifier la configuration de netfilter pour exclure certains hôtes ou nous pouvons insérer une règle avec une priorité plus faible qui pointe sur la table principale pour nos hôtes faisant exception.

Nous pouvons aussi utiliser cette fonctionnalité pour nous conformer aux bits TOS en marquant les paquets avec différents types de service et les nombres correspondants. On crée ensuite les règles qui agissent sur ces types de service. De cette façon, on peut dédier une ligne RNIS aux connexions interactives.

Inutile de le dire, cela marche parfaitement sur un hôte qui fait de la traduction d'adresse (NAT), autrement dit du *masquerading*.

IMPORTANT : Nous avons reçu une information selon laquelle MASQ et SNAT entrent en conflit avec le marquage de paquets. Rusty Russell l'explique dans [ce courrier](#)³.

Désactivez le filtrage de chemin inverse pour que cela fonctionne correctement.

Note : pour marquer les paquets, vous aurez besoin de valider quelques options du noyau :

```
IP: advanced router (CONFIG_IP_ADVANCED_ROUTER) [Y/n/?]
IP: policy routing (CONFIG_IP_MULTIPLE_TABLES) [Y/n/?]
IP: use netfilter MARK value as routing key (CONFIG_IP_ROUTE_FWMARK) [Y/n/?]
```

Voir aussi [Section 15.5, « Cache web transparent utilisant netfilter, iproute2, ipchains et squid »](#) dans le chapitre [Recettes de cuisine](#).

¹ <http://netfilter.org/documentation/>

² <http://netfilter.org/>

³ <http://lists.netfilter.org/pipermail/netfilter/2000-November/006089.html>

Comme expliqué dans la section sur les gestionnaires de mise en file d'attente basés sur des classes, les filtres sont nécessaires pour classifier les paquets dans n'importe laquelle des sous-files d'attente. Ces filtres sont appelés à l'intérieur des gestionnaires de mise en file d'attente basés sur des classes.

Voici une liste incomplète des classificateurs disponibles :

fw

Base la décision sur la façon dont le pare-feu a marqué les paquets. Ceci peut être un passage facile si vous ne voulez pas apprendre la syntaxe **tc** liée aux filtres. Voir le chapitre sur les gestionnaires de mise en file d'attente pour plus de détails.

u32

Base la décision sur les champs à l'intérieur du paquet (c'est-à-dire l'adresse IP source, etc.)

route

Base la décision sur la route que va emprunter le paquet.

rsvp, rsvp6

Route les paquets en se basant sur **RSVP**¹. Seulement utile sur les réseaux que vous contrôlez. Internet ne respecte pas RSVP.

tcindex

Utilisé par le gestionnaire de file d'attente **DSMARK**. Voir la section **DSMARK**.

Notez qu'il y a généralement plusieurs manières de classifier un paquet. Cela dépend du système de classification que vous souhaitez utiliser.

Les classificateurs acceptent en général quelques arguments communs. Ils sont listés ici pour des raisons pratiques :

protocol

Le protocole que ce classificateur acceptera. Généralement, on n'acceptera que le trafic IP. Exigé.

parent

Le descripteur auquel ce classificateur est attaché. Ce descripteur doit être une classe déjà existante. Exigé.

prio

La priorité de ce classificateur. Les plus petits nombres seront testés en premier.

handle

Cette référence a plusieurs significations suivant les différents filtres.

Toutes les sections suivantes supposeront que vous essayez de mettre en forme le trafic allant vers `HostA`. Ces sections supposeront que la classe racine a été configurée sur `1:` et que la classe vers laquelle vous voulez envoyer le trafic sélectionné est `1:1`.

12.1. Le classificateur `u32`

Le filtre `u32` est le filtre le plus avancé dans l'implémentation courante. Il est entièrement basé sur des tables de hachage, ce qui le rend robuste quand il y a beaucoup de règles de filtrage.

Dans sa forme la plus simple, le filtre `u32` est une liste d'enregistrements, chacun consistant en deux champs : un sélecteur et une action. Les sélecteurs, décrits ci-dessous, sont comparés avec le paquet IP traité jusqu'à la première correspondance, et l'action associée est réalisée. Le type d'action le plus simple serait de diriger le paquet vers une classe CBQ définie.

La ligne de commande du programme **tc filter**, utilisée pour configurer le filtre, consiste en trois parties : la spécification du filtre, un sélecteur et une action. La spécification du filtre peut être définie comme :

```
tc filter add dev IF [ protocol PROTO ]
                    [ (preference|priority) PRIO ]
                    [ parent CBQ ]
```

Le champ `protocol` décrit le protocole sur lequel le filtre sera appliqué. Nous ne discuterons que du cas du protocole `ip`. Le champ `preference` (`priority` peut être utilisé comme alternative) fixe la priorité du filtre que l'on définit. C'est important dans la mesure où vous pouvez avoir plusieurs filtres (listes de règles) avec des priorités différentes. Chaque liste sera scrutée dans l'ordre d'ajout des règles. Alors, la liste avec la priorité la plus faible (celle qui a le numéro de préférence le plus élevé) sera traitée. Le champ `parent` définit le sommet de l'arbre CBQ (par ex. `1:0`) auquel le filtre doit être attaché.

Les options décrites s'appliquent à tous les filtres, pas seulement à `u32`.

12.1.1. Le sélecteur `U32`

Le sélecteur `U32` contient la définition d'un modèle, qui sera comparé au paquet traité. Plus précisément, il définit quels bits doivent correspondre dans l'en-tête du paquet, et rien de plus, mais cette méthode simple

¹ <http://www.isi.edu/div7/rsvp/overview.html>

est très puissante. Jetons un oeil sur l'exemple suivant, directement tiré d'un filtre assez complexe réellement existant :

```
# tc filter parent 1: protocol ip pref 10 u32 fh 800::800 order 2048 key ht 800 bkt 0 flowid 1:3 \
  match 00100000/00ff0000 at 0
```

Pour l'instant, laissons de côté la première ligne ; tous ces paramètres décrivent les tables de hachage du filtre. Focalisons-nous sur la ligne de sélection contenant le mot-clé `match`. Ce sélecteur fera correspondre les en-têtes IP dont le second octet sera `0x10` (`0010`). Comme nous pouvons le deviner, le nombre `00ff` est le masque de correspondance, disant au filtre quels bits il doit regarder. Ici, c'est `0xff`, donc l'octet correspondra si c'est exactement `0x10`. Le mot-clé `at` signifie que la correspondance doit démarrer au décalage spécifié (en octets) - dans notre cas, c'est au début du paquet. Traduisons tout cela en langage humain : le paquet correspondra si son champ Type de Service (TOS) a le bit « faible délai » positionné. Analysons une autre règle :

```
# tc filter parent 1: protocol ip pref 10 u32 fh 800::803 order 2051 key ht 800 bkt 0 flowid 1:3 \
  match 00000016/0000ffff at nexthdr+0
```

L'option `nexthdr` désigne l'en-tête suivant encapsulé dans le paquet IP, c'est à dire celui du protocole de la couche supérieure. La correspondance commencera également au début du prochain en-tête. Elle devrait avoir lieu dans le deuxième mot de 32 bits de l'en-tête. Dans les protocoles TCP et UDP, ce champ contient le port de destination du paquet. Le nombre est donné dans le format big-endian, c'est-à-dire les bits les plus significatifs en premier. Il faut donc lire `0x0016` comme 22 en décimal, qui correspond au service SSH dans le cas de TCP. Comme vous le devinez, cette correspondance est ambiguë sans un contexte, et nous en discuterons plus loin.

Ayant compris tout cela, nous trouverons le sélecteur suivant très facile à lire : `match c0a80100/ffffff00 at 16`. Ce que nous avons ici, c'est une correspondance de trois octets au 17ème octet, en comptant à partir du début de l'en-tête IP. Cela correspond aux paquets qui ont une adresse de destination quelconque dans le réseau `192.168.1/24`. Après avoir analysé les exemples, nous pouvons résumer ce que nous avons appris.

12.1.2. Sélecteurs généraux

Les sélecteurs généraux définissent le modèle, le masque et le décalage qui seront comparés au contenu du paquet. En utilisant les sélecteurs généraux, vous pouvez rechercher des correspondances sur n'importe quel bit de l'en-tête IP (ou des couches supérieures). Ils sont quand même plus difficiles à écrire et à lire que les sélecteurs spécifiques décrits ci-dessus. La syntaxe générale des sélecteurs est :

```
match [ u32 | u16 | u8 ] PATTERN MASK [ at OFFSET | nexthdr+OFFSET]
```

Un des mots-clés `u32`, `u16` ou `u8` doit spécifier la longueur du modèle en bits. `PATTERN` et `MASK` se rapportent à la longueur définie par ce mot-clé. Le paramètre `OFFSET` est le décalage, en octets, pour le démarrage de la recherche de correspondance. Si le mot-clé `nexthdr+` est présent, le décalage sera relatif à l'en-tête de la couche réseau supérieure.

Quelques exemples :

```
# tc filter add dev ppp14 parent 1:0 prio 10 u32 \
  match u8 64 0xff at 8 \
  flowid 1:4
```

Un paquet correspondra à cette règle si sa « durée de vie » (TTL) est de 64. TTL est le champ démarrant juste après le 8ème octet de l'en-tête IP.

Correspond à tous les paquets TCP ayant le bit ACK activé :

```
# tc filter add dev ppp14 parent 1:0 prio 10 u32 \
  match ip protocol 6 0xff \
  match u8 0x10 0xff at nexthdr+13 \
  flowid 1:3
```

Utilisez ceci pour déterminer la présence du bit ACK sur les paquets d'une longueur inférieure à 64 octets :

```
## Vérifie la présence d'un ACK,
## protocol IP 6,
## longueur de l'en-tête IP 0x5(mots de 32 bits),
## longueur total IP 0x34 (ACK + 12 octets d'options TCP)
## TCP ack actif (bit 5, offset 33)
# tc filter add dev ppp14 parent 1:0 protocol ip prio 10 u32 \
  match ip protocol 6 0xff \
  match u8 0x05 0x0f at 0 \
  match u16 0x0000 0xffc0 at 2 \
  match u8 0x10 0xff at 33 \
  flowid 1:3
```

Seuls les paquets TCP sans charge utile et avec le bit ACK positionné vérifieront cette règle. Ici, nous pouvons voir un exemple d'utilisation de deux sélecteurs, le résultat final étant un ET logique de leur résultat. Si nous jetons un coup d'oeil sur un schéma de l'en-tête TCP, nous pouvons voir que le bit ACK est le second bit (`0x10`) du 14ème octet de l'en-tête TCP (`at nexthdr+13`). Comme second sélecteur, si nous voulons nous compliquer la vie, nous pouvons écrire `match u8 0x06 0xff at 9` à la place du sélecteur spécifique `protocol tcp`, puisque 6 est le numéro du protocole TCP, spécifié au 10ème octet de l'en-tête IP. D'un autre côté, dans cet exemple, nous ne pourrions pas utiliser de sélecteur spécifique pour la première correspondance, simplement parce qu'il n'y a pas de sélecteur spécifique pour désigner les bits TCP ACK.

Le filtre ci-dessous est une version modifiée du filtre présenté au-dessus. La différence est qu'il ne vérifie pas la longueur de l'en-tête ip. Pourquoi ? Car le filtre au-dessus ne marche que sur les systèmes 32 bits.

```
tc filter add dev ppp14 parent 1:0 protocol ip prio 10 u32 \
  match ip protocol 6 0xff \
  match u8 0x10 0xff at nexthdr+13 \
  match u16 0x0000 0xffc0 at 2 \
  flowid 1:3
```

12.1.3. Les sélecteurs spécifiques

La table suivante contient la liste de tous les sélecteurs spécifiques que les auteurs de cette section ont trouvés dans le code source du programme **tc**. Ils rendent simplement la vie plus facile en accroissant la lisibilité de la configuration du filtre.

FIXME: emplacement de la table - la table est dans un fichier séparé "selector.html"

FIXME: C'est encore en Polonais :-)

FIXME: doit être "sgmlisé"

Quelques exemples :

```
# tc filter add dev ppp0 parent 1:0 prio 10 u32 \
  match ip tos 0x10 0xff \
  flowid 1:4
```

FIXME: tcp dport match ne fonctionne pas comme décrit ci-dessous :

La règle ci-dessus correspondra à des paquets qui ont le champ TOS égal à 0x10. Le champ TOS commence au deuxième octet du paquet et occupe 1 octet, ce qui nous permet d'écrire un sélecteur général équivalent : `match u8 0x10 0xff at 1`. Cela nous donne une indication sur l'implémentation du filtre u32. Les règles spécifiques sont toujours traduites en règles générales, et c'est sous cette forme qu'elles sont stockées en mémoire par le noyau. Cela amène à une autre conclusion : les sélecteurs `tcp` et `udp` sont exactement les mêmes et c'est la raison pour laquelle vous ne pouvez pas utiliser un simple sélecteur `match tcp dport 53 0xffff` pour désigner un paquet TCP envoyé sur un port donné. Ce sélecteur désigne aussi les paquets UDP envoyés sur ce port. Vous devez également spécifier le protocole avec la règle suivante :

```
# tc filter add dev ppp0 parent 1:0 prio 10 u32 \
  match tcp dport 53 0xffff \
  match ip protocol 0x6 0xff \
  flowid 1:2
```

12.2. Le classificateur `route`

Ce classificateur filtre en se basant sur les informations des tables de routage. Quand un paquet passant à travers les classes et en atteint une qui est marquée avec le filtre `route`, il divise le paquet en se basant sur l'information de la table de routage.

```
# tc filter add dev eth1 parent 1:0 protocol ip prio 100 route
```

Ici, nous ajoutons un classificateur `route` sur le nœud parent 1:0, avec la priorité 100. Quand un paquet atteint ce nœud (ce qui arrive immédiatement, puisqu'il est racine), il consulte la table de routage et si une entrée de la table correspond, il envoie le paquet vers la classe donnée et lui donne une priorité de 100. Ensuite, vous ajoutez l'entrée de routage appropriée pour finalement activer les choses.

L'astuce ici est de définir `realm` en se basant soit sur la destination, soit sur la source. Voici la façon de procéder :

```
# ip route add Host/Network via Gateway dev Device realm RealmNumber
```

Par exemple, nous pouvons définir notre réseau de destination 192.168.10.0 avec le nombre `realm` égal à 10 :

```
# ip route add 192.168.10.0/24 via 192.168.10.1 dev eth1 realm 10
```

Quand on ajoute des filtres `route`, on peut utiliser les nombres `realm` pour représenter les réseaux ou les hôtes et spécifier quelle est la correspondance entre les routes et les filtres.

```
# tc filter add dev eth1 parent 1:0 protocol ip prio 100 \
  route to 10 classid 1:10
```

La règle ci-dessus indique que les paquets allant vers le réseau 192.168.10.0 correspondent à la classe 1:10.

Le filtre `route` peut aussi être utilisé avec les routes sources. Par exemple, il y a un sous-réseau attaché à notre routeur Linux sur `eth2`.

```
# ip route add 192.168.2.0/24 dev eth2 realm 2
# tc filter add dev eth1 parent 1:0 protocol ip prio 100 \
  route from 2 classid 1:2
```

Ici, le filtre spécifie que les paquets venant du réseau 192.168.2.0 (`realm 2`) correspondront à la classe 1:2.

12.3. Les filtres de réglementation (*Policing filters*)

Pour réaliser des configurations encore plus compliquées, vous pouvez avoir des filtres qui analysent le trafic à hauteur d'une certaine bande passante. Vous pouvez configurer un filtre pour qu'il cesse complètement l'analyse de tout le trafic au-dessus d'un certain débit ou pour qu'il n'analyse pas la bande passante dépassant un certain débit.

Ainsi, si vous décidez de régler à 4mbit/s, mais qu'un trafic de 5mbit/s est présent, vous pouvez cesser d'analyser l'ensemble des 5mbit/s ou seulement cesser d'analyser le 1 mbit/s supplémentaire et envoyer 4 mbit/s à la classe correspondante.

Si la bande passante dépasse le débit configuré, vous pouvez rejeter un paquet, le reclassifier ou voir si un autre filtre y correspond.

12.3.1. Techniques de réglementation

Il y a essentiellement deux façons de régler. Si vous avez compilé le noyau avec *Estimators*, celui-ci peut mesurer plus ou moins pour chaque filtre le trafic qui est passé. Ces estimations ne sont pas coûteuses en temps CPU, étant donné qu'il ne compte que 25 fois par seconde le nombre de données qui sont passées, et qu'il calcule le débit à partir de là.

L'autre manière utilise encore le *Token Bucket Filter* qui réside à l'intérieur du filtre cette fois. Le TBF analyse seulement le trafic *A HAUTEUR* de la bande passante que vous avez configurée. Si cette bande passante est dépassée, seul l'excès est traité par l'action de dépassement de limite configurée.

12.3.1.1. Avec l'estimateur du noyau

Ceci est très simple et il n'y a qu'un seul paramètre : *avrate*. Soit le flux demeure sous *avrate* et le filtre classe le trafic vers la classe appropriée, soit votre débit le dépasse et l'action indiquée par défaut, la « reclassification », est réalisée dans ce cas.

Le noyau utilise l'algorithme EWMA pour votre bande passante, ce qui la rend moins sensible aux courtes rafales de données.

12.3.1.2. Avec le *Token Bucket Filter*

Utilisez les paramètres suivants :

- *buffer/maxburst*
- *mtu/minburst*
- *mpu*
- *rate*

Ceux-ci se comportent la plupart du temps de manière identique à ceux décrits dans la section *Filtre à seuil de jetons*. Notez cependant que si vous configurez le *mtu* du filtre de réglementation TBF trop bas, aucun paquet ne passera et le gestionnaire de mise en file d'attente de sortie TBF ne fera que les ralentir.

Une autre différence est que la réglementation ne peut que laisser passer ou jeter un paquet. Il ne peut pas le retenir dans le but de le retarder.

12.3.2. Actions de dépassement de limite (*Overlimit actions*)

Si votre filtre décide qu'un dépassement de limite est atteint, il peut mettre en oeuvre des « actions ». Actuellement, trois actions sont disponibles :

continue

Provoque l'arrêt de l'analyse du filtre, bien que d'autres filtres aient la possibilité de le faire.

drop

Ceci est une option très féroce qui supprime simplement le trafic excédant un certain débit. Elle est souvent employée dans le *Ingress policer* et a des utilisations limitées. Par exemple, si vous avez un serveur de noms qui s'écroule s'il traite plus de 5mbit/s de paquets, alors, vous pourrez dans ce cas utiliser un filtre d'entrée pour être sûr qu'il ne traitera jamais plus de 5mbit/s.

Pass/OK

Transmettre le trafic. Peut être utilisé pour mettre hors service un filtre compliqué, tout en le laissant en place.

reclassify

Permet le plus souvent une reclassification vers *Best Effort*. Ceci est l'action par défaut.

12.3.3. Exemples

Le seul vrai exemple connu est mentionné dans la section *Protéger votre machine des inondations SYN*.

FIXME: Si vous avez déjà utilisé ceci, partagez s'il vous plaît votre expérience avec nous.

12.4. Filtres hachés pour un filtrage massif très rapide

Si vous avez besoin de milliers de règles, par exemple, dans le cas où vous avez beaucoup de clients ou d'ordinateurs, tous avec des spécifications QoS différentes, vous pourrez constater que le noyau passe beaucoup de temps à analyser toutes ces règles.

Par défaut, tous les filtres résident dans une grande chaîne qui est analysée par ordre décroissant des priorités. Si vous avez 1000 règles, 1000 contrôles peuvent être nécessaires pour déterminer ce qu'il faut faire d'un paquet.

La vérification irait plus vite s'il y avait 256 chaînes avec chacune quatre règles et si vous pouviez répartir les paquets sur ces 256 chaînes, afin que la bonne règle soit présente.

Ceci est rendu possible par le hachage. Imaginons que vous ayez sur votre réseau 1024 clients avec des modems câble, avec des adresses IP allant de 1.2.0.0 à 1.2.3.255, et que chacun doit avoir un classement particulier, par exemple « pauvre », « moyen » et « bourrage ». Cela vous ferait alors 1024 règles, dans le genre :

```
# tc filter add dev eth1 parent 1:0 protocol ip prio 100 match ip src \
  1.2.0.0 classid 1:1
# tc filter add dev eth1 parent 1:0 protocol ip prio 100 match ip src \
  1.2.0.1 classid 1:1
...
# tc filter add dev eth1 parent 1:0 protocol ip prio 100 match ip src \
  1.2.3.254 classid 1:3
# tc filter add dev eth1 parent 1:0 protocol ip prio 100 match ip src \
  1.2.3.255 classid 1:2
```

Pour aller plus vite, nous pouvons utiliser la dernière partie de l'adresse IP comme « clé de hachage ». Nous obtenons alors 256 tables, la première ressemblant à ceci :

```
# tc filter add dev eth1 parent 1:0 protocol ip prio 100 match ip src \
  1.2.0.0 classid 1:1
# tc filter add dev eth1 parent 1:0 protocol ip prio 100 match ip src \
  1.2.1.0 classid 1:1
# tc filter add dev eth1 parent 1:0 protocol ip prio 100 match ip src \
  1.2.2.0 classid 1:3
# tc filter add dev eth1 parent 1:0 protocol ip prio 100 match ip src \
  1.2.3.0 classid 1:2
```

La suivante commence comme ceci :

```
# tc filter add dev eth1 parent 1:0 protocol ip prio 100 match ip src \
  1.2.0.1 classid 1:1
...
```

De cette manière, seules quatre recherches au plus sont nécessaires et deux en moyenne.

La configuration est plutôt compliquée, mais elle en vaut vraiment la peine du fait des nombreuses règles. Nous créons d'abord un filtre racine, puis une table avec 256 entrées :

```
# tc filter add dev eth1 parent 1:0 prio 5 protocol ip u32
# tc filter add dev eth1 parent 1:0 prio 5 handle 2: u32 divisor 256
```

Nous ajoutons maintenant des règles dans la table précédemment créée :

```
# tc filter add dev eth1 protocol ip parent 1:0 prio 5 u32 ht 2:7b: \
  match ip src 1.2.0.123 flowid 1:1
# tc filter add dev eth1 protocol ip parent 1:0 prio 5 u32 ht 2:7b: \
  match ip src 1.2.1.123 flowid 1:2
# tc filter add dev eth1 protocol ip parent 1:0 prio 5 u32 ht 2:7b: \
  match ip src 1.2.3.123 flowid 1:3
# tc filter add dev eth1 protocol ip parent 1:0 prio 5 u32 ht 2:7b: \
  match ip src 1.2.4.123 flowid 1:2
```

Ceci est l'entrée 123, qui contient les correspondances pour 1.2.0.13, 1.2.1.123, 1.2.2.123 et 1.2.3.123 qui les envoient respectivement vers 1:1, 1:2, 1:3 et 1:2. Notez que nous devons spécifier notre seau de hachage en hexadécimal, 0x7b pour 123.

Nous créons ensuite un « filtre de hachage » qui dirige le trafic vers la bonne entrée de la table de hachage :

```
# tc filter add dev eth1 protocol ip parent 1:0 prio 5 u32 ht 800:: \
  match ip src 1.2.0.0/16 \
  hashkey mask 0x000000ff at 12 \
  link 2:
```

Ok, certains nombres doivent être expliqués. La table de hachage par défaut est appelée 800:: et tous les filtres démarrent de là. Nous sélectionnons alors l'adresse source qui est en position 12, 13, 14 et 15 dans l'en-tête IP, et indiquons que seule la dernière partie nous intéresse. Ceci est envoyé vers la table de hachage 2: qui a été créée plus tôt.

C'est plutôt compliqué, mais cela marche en pratique et les performances seront époustouflantes. Notez que cet exemple pourrait être amélioré pour que chaque chaîne contienne un filtre, ce qui représenterait le cas idéal !

12.5. Filtrer le trafic IPv6

12.5.1. Comment se fait-il que ces filtres tc IPv6 ne fonctionnent pas ?

La base de données des politiques de routage (RPDB) a remplacé le routage IPv4 et la structure d'adressage à l'intérieur du noyau Linux, ce qui a permis les merveilleuses fonctionnalités décrites dans ce HOWTO. Malheureusement, la pile IPv6 à l'intérieur de Linux a été implémentée en dehors de cette structure

principale. Bien qu'ils partagent des fonctionnalités, la structure RPDB de base ne participe pas dans ou avec les structures d'adressage et de routage de IPv6.

Ceci va sûrement changer, nous devons juste attendre un peu plus longtemps.

FIXME : Des idées sur des personnes travaillant sur ce sujet ? Planifications ?

12.5.2. Marquer les paquets IPv6 en utilisant ip6tables

ip6tables est capable de marquer un paquet et de lui assigner un numéro :

```
# ip6tables -A PREROUTING -i eth0 -t mangle -p tcp -j MARK --mark 1
```

Ceci ne va cependant pas nous aider dans la mesure où le paquet ne passera pas par la structure RPDB.

12.5.3. Utiliser le sélecteur u32 pour repérer le paquet IPv6

IPv6 est normalement encapsulé dans un tunnel SIT et transporté à travers les réseaux IPv4. Voir la section sur le tunnel IPv6 pour de plus amples informations quant à la configuration d'un tel tunnel. Ceci nous permet de filtrer les paquets IPv4 en considérant les paquets IPv6 comme la charge utile.

Le filtre suivant repère tous les paquets IPv6 encapsulés dans des paquets IPv4 :

```
# tc filter add dev $DEV parent 10:0 protocol ip prio 10 u32 \
    match ip protocol 41 0xff flowid 42:42
```

Continuons. Supposons que les paquets IPv6 soient envoyés grâce à des paquets IPv4 et que ces paquets n'ont pas d'options. On pourrait utiliser le filtre suivant pour repérer ICMPv6 dans IPv6 dans IPv4 n'ayant aucune option. 0x3a (58) est le type du champ en-tête suivant pour ICMPv6.

```
# tc filter add dev $DEV parent 10:0 protocol ip prio 10 u32 \
    match ip protocol 41 0xff \
    match u8 0x05 0x0f at 0 \
    match u8 0x3a 0xff at 26 \
    flowid 42:42
```

Repérer l'adresse de destination IPv6 nécessite un peu plus de travail. Le filtre suivant repère l'adresse de destination 3ffe:202c:ffff:32:230:4fff:fe08:358d:

```
# tc filter add dev $DEV parent 10:0 protocol ip prio 10 u32 \
    match ip protocol 41 0xff \
    match u8 0x05 0x0f at 0 \
    match u8 0x3f 0xff at 44 \
    match u8 0xfe 0xff at 45 \
    match u8 0x20 0xff at 46 \
    match u8 0x2c 0xff at 47 \
    match u8 0xff 0xff at 48 \
    match u8 0xff 0xff at 49 \
    match u8 0x00 0xff at 50 \
    match u8 0x32 0xff at 51 \
    match u8 0x02 0xff at 52 \
    match u8 0x30 0xff at 53 \
    match u8 0x4f 0xff at 54 \
    match u8 0xff 0xff at 55 \
    match u8 0xfe 0xff at 56 \
    match u8 0x08 0xff at 57 \
    match u8 0x35 0xff at 58 \
    match u8 0x8d 0xff at 59 \
    flowid 10:13
```

La même technique peut être utilisée pour repérer les réseaux. Par exemple 2001::

```
# tc filter add dev $DEV parent 10:0 protocol ip prio 10 u32 \
    match ip protocol 41 0xff \
    match u8 0x05 0x0f at 0 \
    match u8 0x20 0xff at 28 \
    match u8 0x01 0xff at 29 \
    flowid 10:13
```

Le noyau utilise de nombreux paramètres qui peuvent être ajustés en différentes circonstances. Bien que, comme d'habitude, les paramètres par défaut conviennent à 99% des installations, nous ne pourrions pas appeler ce document « HOWTO avancé » sans en dire un mot.

Les éléments intéressants sont dans `/proc/sys/net`, jetez-y un oeil. Tout ne sera pas documenté ici au départ, mais nous y travaillons.

En attendant, vous pouvez jeter un oeil dans les sources du noyau Linux et lire le fichier `Documentation/filesystems/proc.txt`. La plupart des fonctionnalités y sont expliquées.

13.1. Filtrage de Chemin Inverse (*Reverse Path Filtering*)

Par défaut, les routeurs routent tout, même les paquets qui visiblement n'appartiennent pas à votre réseau. Un exemple courant est l'espace des adresses IP privées s'échappant sur Internet. Si vous avez une interface avec une route pour `195.96.96.0/24` dessus, vous ne vous attendrez pas à voir arriver des paquets venant de `212.64.94.1`.

Beaucoup d'utilisateurs veulent désactiver cette fonctionnalité. Les développeurs du noyau ont permis de le faire facilement. Il y a des fichiers dans `/proc` où vous pouvez ordonner au noyau de le faire pour vous. La méthode est appelée « Filtrage par Chemin Inverse » (*Reverse Path Filtering*). Pour faire simple, si la réponse à ce paquet ne sort pas par l'interface par laquelle il est entré, alors c'est un paquet « bogué » et il sera ignoré.

Les instructions suivantes vont activer cela pour toutes les interfaces courantes et futures.

```
# for i in /proc/sys/net/ipv4/conf/*/rp_filter ; do
> echo 2 > $i
> done
```

En reprenant l'exemple du début, si un paquet arrivant sur le routeur Linux par `eth1` prétend venir du réseau Bureau+FAI, il sera éliminé. De même, si un paquet arrivant du réseau Bureau prétend être de quelque part à l'extérieur du pare-feu, il sera également éliminé.

Ce qui est présenté ci-dessus est le filtrage de chemin inverse complet. Le paramétrage par défaut filtre seulement sur les adresses IP des réseaux directement connectés. Ce paramétrage par défaut est utilisé parce que le filtrage complet échoue dans le cas d'un routage asymétrique (où il y a des paquets arrivant par un chemin et ressortant par un autre, comme dans le cas du trafic satellite ou si vous avez des routes dynamiques (bgp, ospf, rip) dans votre réseau. Les données descendent vers la parabole satellite et les réponses repartent par des lignes terrestres normales).

Si cette exception s'applique dans votre cas (vous devriez être au courant), vous pouvez simplement désactiver le `rp_filter` sur l'interface d'arrivée des données satellite. Si vous voulez voir si des paquets sont éliminés, le fichier `log_martians` du même répertoire indiquera au noyau de les enregistrer dans votre syslog.

```
# echo 1 >/proc/sys/net/ipv4/conf/<interfacename>/log_martians
```

FIXME: Est-ce que la configuration des fichiers dans `.../conf/{default,all}` suffit ? - martijn

13.2. Configurations obscures

Bon, il y a beaucoup de paramètres qui peuvent être modifiés. Nous essayons de tous les lister. Voir aussi une documentation partielle dans `Documentation/ip-sysctl.txt`.

Certaines de ces configurations ont des valeurs par défaut différentes suivant que vous répondez Yes ou No à la question `Configure as router and not host` lors de la compilation du noyau.

Oskar Andreasson a une page sur ces options et il apparaît qu'elle soit meilleure que la notre. De ce fait, allez également voir <http://ipsysctl-tutorial.frozentux.net>¹.

13.2.1. ipv4 générique

En remarque générale, les fonctionnalités de limitation de débit ne fonctionnent pas sur l'interface `loopback`. N'essayez donc pas de les tester localement. Les limites sont exprimées en « tic-tac » (*jiffies*) et elles utilisent obligatoirement le *Token Bucket Filter* mentionné plus tôt.

[NdT : le terme *jiffies* désigne un mouvement régulier, faisant référence au « tic-tac » d'une horloge. Dans le noyau lui-même, une variable globale nommée `jiffies` est incrémentée à chaque interruption d'horloge]

Le noyau a une horloge interne qui tourne à HZ impulsions (ou *jiffies*) par seconde. Sur Intel, HZ est la plupart du temps égale à 100. Donc, configurer un fichier `*_rate` à, disons 50, autorise 2 paquets par seconde. Le *Token Bucket Filter* est également configuré pour autoriser une rafale de données de 6 paquets au plus, si suffisamment de jetons ont été gagnés.

Plusieurs éléments de la liste suivante proviennent du fichier `/usr/src/linux/Documentation/networking/ip-sysctl.txt`, écrit par Alexey Kuznetsov et Andi Kleen.

```
/proc/sys/net/ipv4/icmp_destunreach_rate
```

Si le noyau décide qu'il ne peut pas délivrer un paquet, il le rejettera et enverra à la source du paquet un ICMP notifiant ce rejet.

```
/proc/sys/net/ipv4/icmp_echo_ignore_all
```

N'agit en aucun cas comme écho pour les paquets. Ne configurez pas ceci par défaut. Cependant, si vous êtes utilisé comme relais dans une attaque de Déni de Services, cela peut être utile.

¹ <http://ipsysctl-tutorial.frozentux.net/>

`/proc/sys/net/ipv4/icmp_echo_ignore_broadcasts` [Utile]

Si vous pinguez l'adresse de diffusion d'un réseau, tous les hôtes sont sensés répondre. Cela permet de coquettes attaques de déni de service. Mettez cette valeur à 1 pour ignorer ces messages de diffusion.

`/proc/sys/net/ipv4/icmp_echo_reply_rate`

Le débit auquel les réponses echo sont envoyées aux destinataires.

`/proc/sys/net/ipv4/icmp_ignore_bogus_error_responses`

Configurer ceci pour ignorer les erreurs ICMP d'hôtes du réseau réagissant mal aux trames envoyées vers ce qu'ils perçoivent comme l'adresse de diffusion.

`/proc/sys/net/ipv4/icmp_paramprob_rate`

Un message ICMP relativement peu connu, qui est envoyé en réponse à des paquets qui ont des en-têtes IP ou TCP erronés. Avec ce fichier, vous pouvez contrôler le débit auquel il est envoyé.

`/proc/sys/net/ipv4/icmp_timeexceed_rate`

Voici la célèbre cause des « étoiles Solaris » dans traceroute. Limite le nombre de messages ICMP Time Exceeded envoyés.

`/proc/sys/net/ipv4/igmp_max_memberships`

Nombre maximal de sockets igmp (multidistribution) en écoute sur l'hôte. FIXME: Est-ce vrai ?

`/proc/sys/net/ipv4/inet_peer_gc_maxtime`

FIXME : Ajouter une petite explication sur le stockage des partenaires internet (inet peer) ? Intervalle de temps minimum entre deux passages du ramasse-miettes. Cet intervalle est pris en compte lors d'une faible (voire inexistante) utilisation du *pool*. Mesuré en *jiffies*. [NdT : Le *pool* désigne ici la liste des adresses IP des partenaires internet.]

`/proc/sys/net/ipv4/inet_peer_gc_mintime`

Intervalle de temps minimum entre deux passages du ramasse-miettes. Cet intervalle est pris en compte lors d'une utilisation intensive du *pool*. Mesuré en *jiffies*.

`/proc/sys/net/ipv4/inet_peer_maxttl`

Durée de conservation maximale des enregistrements. Les entrées non utilisées expireront au bout de cet intervalle de temps (c'est-à-dire quand le nombre d'entrées dans le *pool* est très petit). Mesuré en *jiffies*.

`/proc/sys/net/ipv4/inet_peer_minttl`

Durée de conservation minimale des enregistrements. Devrait être suffisante pour prendre en compte le temps de vie des fragments sur l'hôte qui doit réassembler les paquets. Cette durée minimale est garantie si le nombre d'éléments dans le *pool* est inférieur au seuil fixé par `inet_peer_threshold`.

`/proc/sys/net/ipv4/inet_peer_threshold`

Taille approximative de l'espace de stockage des partenaires internet. A partir de ce seuil, les entrées sont effacées. Ce seuil détermine la durée de vie des entrées, ainsi que les intervalles de temps entre deux déclenchements du ramasse-miettes. Plus il y a d'entrées, plus le temps de vie est faible et plus l'intervalle du ramasse-miettes est faible.

`/proc/sys/net/ipv4/ip_autoconfig`

Ce fichier contient la valeur 1 si l'hôte a reçu sa configuration IP par RARP, BOOTP, DHCP ou un mécanisme similaire. Autrement, il contient la valeur zéro.

`/proc/sys/net/ipv4/ip_default_ttl`

Durée de vie (TTL) des paquets. Fixer à la valeur sûre de 64. Augmentez-la si vous avez un réseau immense, mais pas « pour s'amuser » : les boucles sans fin d'un mauvais routage sont plus dangereuses si le TTL est élevé. Vous pouvez même envisager de diminuer la valeur dans certaines circonstances.

`/proc/sys/net/ipv4/ip_dynaddr`

Vous aurez besoin de positionner cela si vous utilisez la connexion à la demande avec une adresse d'interface dynamique. Une fois que votre interface a été configurée, toutes les sockets TCP locaux qui n'ont pas eu de paquets de réponse seront retraitées pour avoir la bonne adresse. Cela résout le problème posé par une connexion défectueuse ayant configuré une interface, suivie par une deuxième tentative réussie (avec une adresse IP différente).

`/proc/sys/net/ipv4/ip_forward`

Le noyau doit-il essayer de transmettre les paquets ? Désactivé par défaut.

`/proc/sys/net/ipv4/ip_local_port_range`

Intervalle des ports locaux pour les connexions sortantes. En fait, assez petit par défaut, 1024 à 4999.

`/proc/sys/net/ipv4/ip_no_pmtu_disc`

Configurez ceci si vous voulez désactiver la découverte du MTU de chemin, une technique pour déterminer le plus grand MTU possible sur votre chemin. Voir aussi la section sur la découverte du MTU de chemin dans le chapitre *Recettes de cuisine*.

```
/proc/sys/net/ipv4/ipfrag_high_thresh
```

Mémoire maximum utilisée pour réassembler les fragments IP. Quand `ipfrag_high_thresh` octets de mémoire sont alloués pour cela, le gestionnaire de fragments rejettera les paquets jusqu'à ce que `ipfrag_low_thresh` soit atteint.

```
/proc/sys/net/ipv4/ip_nonlocal_bind
```

Configurez ceci si vous voulez que vos applications soient capables de se lier à une adresse qui n'appartient pas à une interface de votre système. Ceci peut être utile quand votre machine est sur un lien non-permanent (ou même permanent). Vos services sont donc capables de démarrer et de se lier à une adresse spécifique quand votre lien est inactif.

```
/proc/sys/net/ipv4/ipfrag_low_thresh
```

Mémoire minimale utilisée pour réassembler les fragments IP.

```
/proc/sys/net/ipv4/ipfrag_time
```

Temps en secondes du maintien d'un fragment IP en mémoire.

```
/proc/sys/net/ipv4/tcp_abort_on_overflow
```

Une option booléenne contrôlant le comportement dans le cas de nombreuses connexions entrantes. Quand celle-ci est activée, le noyau envoie rapidement des paquets RST quand un service est surchargé.

```
/proc/sys/net/ipv4/tcp_fin_timeout
```

Temps de maintien de l'état `FIN-WAIT-2` pour un socket dans le cas où il a été fermé de notre côté. Le partenaire peut être défectueux et ne jamais avoir fermé son côté ou même mourir de manière inattendue. La valeur par défaut est de 60 secondes. La valeur usuelle utilisée dans le noyau 2.2 était de 180 secondes. Vous pouvez la remettre, mais rappelez vous que si votre machine a un serveur WEB surchargé, vous risquez de dépasser la mémoire avec des kilotonnes de sockets morts. Les sockets `FIN-WAIT2` sont moins dangereux que les sockets `FIN-WAIT1` parce qu'ils consomment au maximum 1,5K de mémoire, mais ils ont tendance à vivre plus longtemps. Cf `tcp_max_orphans`.

```
/proc/sys/net/ipv4/tcp_keepalive_time
```

Durée entre l'envoi de deux messages *keepalive* quand l'option *keepalive* est activée. Par défaut : 2 heures.

```
/proc/sys/net/ipv4/tcp_keepalive_intvl
```

A quelle fréquence les sondes sont retransmises lorsqu'il n'y a pas eu acquittement de sonde. Par défaut : 75 secondes.

```
/proc/sys/net/ipv4/tcp_keepalive_probes
```

Combien de sondes TCP *keepalive* seront envoyées avant de décider que la connexion est brisée. Par défaut : 9. En multipliant par `tcp_keepalive_intvl`, cela donne le temps pendant lequel un lien peut être actif sans donner de réponses après l'envoi d'un *keepalive*.

```
/proc/sys/net/ipv4/tcp_max_orphans
```

Nombre maximum de sockets TCP qui ne sont pas reliés à un descripteur de fichier utilisateur, géré par le système. Si ce nombre est dépassé, les connexions orphelines sont immédiatement réinitialisées et un avertissement est envoyé. Cette limite existe seulement pour prévenir des attaques de déni de services simples. Vous ne devez pas compter sur ceci ou diminuer cette limite artificiellement, mais plutôt l'augmenter (probablement après avoir augmenté la mémoire) si les conditions du réseau réclament plus que cette valeur par défaut et régler vos services réseau pour qu'ils détruisent sans tarder ce type d'état. Laissez-moi vous rappeler encore que chaque orphelin consomme jusqu'à environ 64K de mémoire non *swappable*.

```
/proc/sys/net/ipv4/tcp_orphan_retries
```

Combien d'essais avant de détruire une connexion TCP, fermée par notre côté. La valeur par défaut de 7 correspond à un temps d'environ 50s à 16 min suivant le RTO. Si votre machine supporte un serveur Web, vous pouvez envisager de baisser cette valeur, dans la mesure où de tels sockets peuvent consommer des ressources significatives. Cf `tcp_max_orphans`.

```
/proc/sys/net/ipv4/tcp_max_syn_backlog
```

Nombre maximum de requêtes d'une connexion mémorisée, qui n'avait pas encore reçu d'accusé de réception du client connecté. La valeur par défaut est de 1024 pour des systèmes avec plus de 128 Mo de mémoire et 128 pour des machines avec moins de mémoire. Si un serveur souffre de surcharge, essayez d'augmenter ce nombre. Attention ! Si vous positionnez une valeur supérieure à 1024, il serait préférable de changer `TCP_SYNQ_HSIZE` dans le fichier `include/net/tcp.h` pour garder `TCP_SYNQ_HSIZE*16 <= tcp_max_syn_backlog` et de recompiler de noyau.

```
/proc/sys/net/ipv4/tcp_max_tw_buckets
```

Nombre maximum de sockets `timewait` gérées par le système simultanément. Si ce nombre est dépassé, le socket `timewait` est immédiatement détruit et un message d'avertissement est envoyé. Cette limite n'existe que pour prévenir des attaques de déni de services simples. Vous ne devez pas diminuer cette limite artificiellement, mais plutôt l'augmenter (probablement après avoir augmenté la mémoire) si les conditions du réseau réclament plus que cette valeur par défaut.

`/proc/sys/net/ipv4/tcp_retrans_collapse`
Compatibilité bug à bug avec certaines imprimantes défectueuses. Tentative d'envoi de plus gros paquets lors de la retransmission pour contourner le bug de certaines piles TCP.

`/proc/sys/net/ipv4/tcp_retries1`
Combien d'essais avant de décider que quelque chose est erroné et qu'il est nécessaire d'informer de cette suspicion la couche réseau. La valeur minimale du RFC est de 3. C'est la valeur par défaut ; elle correspond à un temps d'environ 3 sec à 8 min suivant le RTO.

`/proc/sys/net/ipv4/tcp_retries2`
Combien d'essais avant de détruire une connexion TCP active. Le [RFC 1122](#)² précise que la limite ne devrait pas dépasser 100 secondes. C'est un nombre trop petit. La valeur par défaut de 15 correspond à un temps de environ 13 à 30 minutes suivant le RTO.

`/proc/sys/net/ipv4/tcp_rfc1337`
Ce booléen active un rectificatif pour « l'assassinat hasardeux des time-wait dans tcp », décrit dans le RFC 1337. S'il est activé, le noyau rejette les paquets RST pour les sockets à l'état de time-wait. Par défaut : 0

`/proc/sys/net/ipv4/tcp_sack`
Utilise un ACK sélectif qui peut être utilisé pour signifier que des paquets spécifiques sont manquant. Facilite ainsi une récupération rapide.

`/proc/sys/net/ipv4/tcp_stdurg`
Utilise l'interprétation du RFC *Host Requirements* du champ TCP pointeur urgent. La plupart des hôtes utilisent la vieille interprétation BSD. Donc, si vous activez cette option, il se peut que Linux ne communique plus correctement avec eux. Par défaut : FALSE (FAUX)

`/proc/sys/net/ipv4/tcp_syn_retries`
Nombre de paquets SYN que le noyau enverra avant de tenter l'établissement d'une nouvelle connexion.

`/proc/sys/net/ipv4/tcp_synack_retries`
Pour ouvrir l'autre côté de la connexion, le noyau envoie un SYN avec un ACK superposé (*piggyback*), pour accuser réception du SYN précédemment envoyé. C'est la deuxième partie de la poignée de main à trois voies (*three-way handshake*). Cette configuration détermine le nombre de paquets SYN+ACK à envoyer avant que le noyau n'abandonne la connexion.

`/proc/sys/net/ipv4/tcp_timestamps`
Les estampillages horaires sont utilisés, entre autres, pour se protéger du rebouclage des numéros de séquence. On peut concevoir qu'un lien à 1 gigabit puisse de nouveau rencontrer un numéro de séquence précédent avec une valeur hors-ligne parcequ'il était d'une génération précédente. L'estampillage horaire permet de reconnaître cet « ancien paquet ».

`/proc/sys/net/ipv4/tcp_tw_recycle`
Mise en place du recyclage rapide des sockets TIME-WAIT. La valeur par défaut est 1. Celle-ci ne devrait pas être changée sans le conseil/demande d'experts techniques.

`/proc/sys/net/ipv4/tcp_window_scaling`
TCP/IP autorise normalement des fenêtres jusqu'à une taille de 65535 octets. Pour des réseaux vraiment rapides, cela peut ne pas être assez. Les options `windows_scaling` autorisent des fenêtres jusqu'au gigaoctet, ce qui est adapté pour les produits à grande bande passante.

13.2.2. Configuration des périphériques

DEV peut désigner soit une interface réelle, soit `all`, soit `default`. `default` change également les paramètres des interfaces qui seront créées par la suite.

`/proc/sys/net/ipv4/conf/DEV/accept_redirects`
Si un routeur décide que vous l'utilisez à tort (c'est-à-dire qu'il a besoin de ré-envoyer votre paquet sur la même interface), il vous enverra un message ICMP Redirect. Cela présente cependant un petit risque pour la sécurité, et vous pouvez le désactiver ou utiliser les redirections sécurisées.

`/proc/sys/net/ipv4/conf/DEV/accept_source_route`
Plus vraiment utilisé. On l'utilisait pour être capable de donner à un paquet une liste d'adresses IP à visiter. Linux peut être configuré pour satisfaire cette option IP.

`/proc/sys/net/ipv4/conf/DEV/bootp_relay`
Accepte les paquets avec une adresse source 0.b.c.d et des adresses destinations qui ne correspondent ni à cet hôte, ni au réseau local. On suppose qu'un démon de relais BOOTP interceptera et transmettra de tels paquets.

La valeur par défaut est 0, puisque cette fonctionnalité n'est pas encore implémentée (noyau 2.2.12).

`/proc/sys/net/ipv4/conf/DEV/forwarding`
Active ou désactive la transmission IP sur cette interface.

² <http://www.ietf.org/rfc/rfc1122.txt>

/proc/sys/net/ipv4/conf/DEV/log_martians

Voir la section sur le *Filtrage de Chemin Inverse*.

/proc/sys/net/ipv4/conf/DEV/mc_forwarding

Si vous faites de la transmission multidistribution (*multicast*) sur cette interface.

/proc/sys/net/ipv4/conf/DEV/proxy_arp

Si vous configurez ceci à 1, cet interface répondra aux requêtes ARP pour les adresses que le noyau doit router. Peut être très utile si vous mettez en place des « pseudo ponts ip ». Prenez bien garde d'avoir des masques de sous-réseau corrects avant d'activer cette option. Faites également attention que le `rp_filter` agisse aussi sur les requêtes ARP !

/proc/sys/net/ipv4/conf/DEV/rp_filter

Voir la section sur le *Filtrage de Chemin Inverse*.

/proc/sys/net/ipv4/conf/DEV/secure_redirects

Accepte les messages de redirection ICMP seulement pour les passerelles indiquées dans la liste des passerelles par défaut. Activé par défaut.

/proc/sys/net/ipv4/conf/DEV/send_redirects

Active la possibilité d'envoyer les messages de redirections mentionnées ci-dessus.

/proc/sys/net/ipv4/conf/DEV/shared_media

Si cela n'est pas activé, le noyau ne considère pas que différents sous-réseaux peuvent communiquer directement sur cette interface. La configuration par défaut est `Yes`.

/proc/sys/net/ipv4/conf/DEV/tag

FIXME: à remplir

13.2.3. Politique de voisinage

DEV peut désigner soit une interface réelle, soit `all`, soit `default`. `default` change également les paramètres des interfaces qui seront créées par la suite.

/proc/sys/net/ipv4/neighbor/DEV/anycast_delay

Valeur maximum du délai aléatoire de réponse exprimé en *jiffies* (1/100 sec) aux messages de sollicitation des voisins. N'est pas encore implémenté (Linux ne possède pas encore le support *anycast*).

/proc/sys/net/ipv4/neighbor/DEV/app_solicit

Détermine le nombre de requêtes à envoyer au démon ARP de l'espace utilisateur. Utilisez 0 pour désactiver.

/proc/sys/net/ipv4/neighbor/DEV/base_reachable_time

Une valeur de base utilisée pour le calcul du temps aléatoire d'accès comme spécifié dans le RFC2461.

/proc/sys/net/ipv4/neighbor/DEV/delay_first_probe_time

Délai avant de tester pour la première fois si le voisin peut être atteint. (voir `gc_stale_time`)

/proc/sys/net/ipv4/neighbor/DEV/gc_stale_time

Détermine la fréquence à laquelle on doit vérifier les vieilles entrées ARP. Si une entrée est obsolète, elle devra de nouveau être résolue (ce qui est utile quand une adresse IP a été attribuée à une autre machine). Si `ucast_solicit` est supérieur à 0, alors on essaie d'abord d'envoyer un paquet ARP directement à l'hôte connu. Si cela échoue, et que `mcast_solicit` est supérieur à 0, alors une requête ARP est multidiffusée.

/proc/sys/net/ipv4/neighbor/DEV/locktime

Une entrée ARP n'est remplacée par une nouvelle que si l'ancienne est au moins présente depuis `locktime`. Cela évite trop d'écriture dans le cache.

/proc/sys/net/ipv4/neighbor/DEV/mcast_solicit

Nombre maximum d'essais consécutifs pour une sollicitation *multicast*.

/proc/sys/net/ipv4/neighbor/DEV/proxy_delay

Temps maximum (le temps réel est aléatoire et compris entre 0 et `proxytime`) avant de répondre à une requête ARP pour laquelle nous avons une entrée de proxy ARP. Peut être utilisé dans certains cas pour se prémunir des inondations réseaux.

/proc/sys/net/ipv4/neighbor/DEV/proxy_qlen

Longueur maximale de la file d'attente du temporisateur de cache arp en attente (Voir `proxy_delay`).

/proc/sys/net/ipv4/neighbor/DEV/retrans_time

Le temps, exprimé en *jiffies* (1/100 sec), entre deux requêtes ARP. Utilisé pour la résolution d'adresses et pour déterminer si un voisin est inaccessible.

/proc/sys/net/ipv4/neighbor/DEV/ucast_solicit

Nombre maximum de requêtes ARP unicast.

`/proc/sys/net/ipv4/neigh/DEV/unres_qlen`

Longueur maximum de la file d'attente pour la requête ARP en cours : le nombre de paquets qui sont acceptés des autres couches pendant la résolution ARP d'une adresse.

Internet QoS: Architectures and Mechanisms for Quality of Service, Zheng Wang, ISBN 1-55860-608-4

Livre traitant des sujets liés à la qualité de service. Bien pour comprendre les concepts de base.

13.2.4. Configuration du routage

`/proc/sys/net/ipv4/route/error_burst`

Ces paramètres sont utilisés pour limiter le nombre de messages d'avertissement écrits dans le journal du noyau par le code de routage. Plus le paramètre `error_burst` est grand, moins il y aura de messages. `Error_burst` contrôle le moment où les messages seront supprimés. Les configurations par défaut se limitent à un message d'avertissement toutes les cinq secondes.

`/proc/sys/net/ipv4/route/error_cost`

Ces paramètres sont utilisés pour limiter le nombre de messages d'avertissement écrits dans le journal du noyau par le code de routage. Plus le paramètre `error_cost` est grand, moins il y aura de messages. `error_burst` contrôle le moment où les messages seront jetés. Les configurations par défaut se limitent à un message d'avertissement toutes les cinq secondes.

`/proc/sys/net/ipv4/route/flush`

L'écriture dans ce fichier provoque la vidange du cache du routage.

`/proc/sys/net/ipv4/route/gc_elasticity`

Valeurs qui contrôlent la fréquence et le comportement de l'algorithme *garbage collection* du cache de routage. Ceci peut être important en cas de défaut. Au moins `gc_timeout` secondes s'écouleront avant que le noyau ne passe à une autre route si la précédente n'est plus opérationnelle. Configuré par défaut à 300. Diminuez cette valeur si vous voulez passer plus rapidement ce type de problème.

Voir aussi [ce message](#)³ par Ard van Breemen.

`/proc/sys/net/ipv4/route/gc_interval`

Voir `/proc/sys/net/ipv4/route/gc_elasticity`.

`/proc/sys/net/ipv4/route/gc_min_interval`

Voir `/proc/sys/net/ipv4/route/gc_elasticity`.

`/proc/sys/net/ipv4/route/gc_thresh`

Voir `/proc/sys/net/ipv4/route/gc_elasticity`.

`/proc/sys/net/ipv4/route/gc_timeout`

Voir `/proc/sys/net/ipv4/route/gc_elasticity`.

`/proc/sys/net/ipv4/route/max_delay`

Délai d'attente pour la vidange du cache du routage.

`/proc/sys/net/ipv4/route/max_size`

Taille maximum du cache de routage. Les vieilles entrées seront purgées quand le cache aura atteint cette taille.

`/proc/sys/net/ipv4/route/min_adv_mss`

FIXME: à remplir

`/proc/sys/net/ipv4/route/min_delay`

Délai d'attente pour vider le cache de routage.

`/proc/sys/net/ipv4/route/min_pmtu`

FIXME: à remplir

`/proc/sys/net/ipv4/route/mtu_expires`

FIXME: à remplir

`/proc/sys/net/ipv4/route/redirect_load`

Facteurs qui déterminent si plus de redirections ICMP doivent être envoyées à un hôte spécifique. Aucune redirection ne sera envoyée une fois que la limite de charge (*load limit*) ou que le nombre maximum de redirections aura été atteint.

`/proc/sys/net/ipv4/route/redirect_number`

Voir `/proc/sys/net/ipv4/route/redirect_load`.

`/proc/sys/net/ipv4/route/redirect_silence`

Temporisation pour les redirections. Au delà de cette période, les redirections seront de nouveau envoyées, même si elles ont été stoppées parce que la charge ou le nombre limite a été atteint.

³ <http://mailman.ds9a.nl/pipermail/lartc/2002q1/002667.html>

Si vous constatez que vous avez des besoins qui ne sont pas gérés par les files d'attente citées précédemment, le noyau contient quelques autres files d'attente plus spécialisées mentionnées ici.

14.1. bfifo/pfifo

Ces files d'attente sans classes sont plus simples que `pfifo_fast` dans la mesure où elles n'ont pas de bandes internes, tout le trafic étant vraiment équivalent. Elles ont cependant l'avantage important de réaliser des statistiques. Donc, même si vous n'avez pas besoin de mise en forme ou de donner une priorité, vous pouvez employer ce gestionnaire pour déterminer l'arriéré (*backlog*) de votre interface.

`pfifo` mesure en paquets et `bfifo` en octets.

14.1.1. Paramètres & usage

`limit`

Spécifie la taille de la file d'attente. Mesurée en octets pour `bfifo` et en paquets pour `pfifo`. Par défaut, correspond à des paquets de taille égale au paramètre `txqueuelen` de l'interface (voir le chapitre `pfifo_fast`) ou `txqueuelen*mtu` octets pour `bfifo`.

14.2. Algorithme Clark-Shenker-Zhang (CSZ)

Ceci est si théorique que même Alexey (l'auteur principal de CBQ) prétend ne pas le comprendre. De son propre avis :

David D. Clark, Scott Shenker and Lixia Zhang *Supporting Real-Time Applications in an Integrated Services Packet Network: Architecture and Mechanism*.

Comme je le comprends, l'idée principale est de créer des flux WFQ pour chaque service garanti et d'allouer le reste de la bande passante au flux factice, appelé `flow-0`. Le Flux-0 inclut le trafic de service prédictif et le trafic *best-effort*. Il est traité par un ordonnanceur de priorité qui alloue la bande passante de plus grande priorité aux services prédictifs, et le reste aux paquets *best-effort*.

Notez que dans CSZ, les flux ne sont PAS limités à leur bande passante. On suppose que le flux a passé le contrôle d'admission à la frontière du réseau QoS et qu'il n'a pas besoin de mises en forme supplémentaires. N'importe quelles autres tentatives pour améliorer le flux ou pour le mettre en forme grâce à un seau de jetons lors d'étapes intermédiaires introduiront des retards non désirés et augmenteront la gigue.

A l'heure actuelle, CSZ est le seul ordonnanceur qui fournit un véritable service garanti. Les autres implémentations (incluant CBQ) n'assurent pas un délai garanti et rendent la gigue aléatoire.

Ne semble pas actuellement un bon candidat à utiliser, à moins que vous n'ayez lu et compris l'article mentionné.

14.3. DSMARK

Esteve Camps Ce texte est un extrait de ma thèse sur le *support QoS dans Linux*, Septembre 2000.

Documents sources :

- [Draft-almesberger-wajhak-diffserv-linux-01.txt](#)¹.
- Exemples de la distribution `iproute2`.
- [White Paper-QoS protocols and architectures](#)² et [Foires Aux Questions IP QoS](#)³, les deux par *Quality of Service Forum*.

14.3.1. Introduction

Avant tout, il serait préférable de lire les RFC écrits sur ce sujet (RFC2474, RFC2475, RFC2597 et RFC2598) sur le [site web du groupe de travail DiffServ IETF](#)⁴ et sur le [site web de Werner Almesberger](#)⁵ (Il a écrit le code permettant le support des Services Différenciés sous Linux).

14.3.2. A quoi DSMARK est-il relié ?

DSMARK est un gestionnaire de mise en file d'attente qui offre les fonctionnalités dont ont besoin les services différenciés (*Differentiated Services*) (également appelés DiffServ ou tout simplement DS). *DiffServ* est l'une des deux architectures actuelles de la Qualité de Services (QoS : *Quality of Services*) (l'autre est appelée *services intégrés* (*Integrated Services*)). Elle se base sur la valeur du champ DS contenu dans l'en-tête IP du paquet.

¹ <ftp://icaftp.epfl.ch/pub/linux/diffserv/misc/dsid-01.txt.gz>

² http://www.qosforum.com/white-papers/qosprot_v3.pdf

³ <http://www.qosforum.com/docs/faq>

⁴ <http://www.ietf.org/html.charters/diffserv-charter.html>

⁵ <http://diffserv.sf.net/>

Une des premières solutions dans IP pour offrir des niveaux de qualité de services était le champ *Type de Service* (octet TOS) de l'en-tête IP. En modifiant la valeur de ce champ, nous pouvions choisir un niveau élevé/faible du débit, du délai ou de la fiabilité. Cependant, cela ne fournissait pas une flexibilité suffisante pour les besoins de nouveaux services (comme les applications temps réel, les applications interactives et autres). Par la suite, de nouvelles architectures sont apparues. L'une d'elle a été *DiffServ* qui a gardé les bits TOS et les a renommés champ DS.

14.3.3. Guide des services différenciés

Les services différenciés sont orientés groupes. Cela signifie que nous ne savons rien des flux (ce sera le but des services intégrés (*integrated Services*)). Nous connaissons en revanche les agrégations de flux et nous adopterons des comportements différents suivant l'agrégation à laquelle appartient le paquet.

Quand un paquet arrive à un nœud frontalier (nœud d'entrée du domaine DiffServ) et entre dans un domaine DiffServ, nous devons avoir une politique, une mise en forme et/ou un marquage de ces paquets (le marquage fait référence à la mise en place d'une valeur dans le champ DS. Comme on le ferait pour des vaches :-)). Ce sera cette marque/valeur que les nœuds internes de votre domaine DiffServ regarderont pour déterminer quel comportement ou niveau de qualité de service appliquer.

Comme vous pouvez le déduire, les Services Différenciés impliquent un domaine sur lequel toutes les règles DS devront être appliquées. Vous pouvez raisonner de la façon suivante : « Je classifierai tous les paquets entrant dans mon domaine. Une fois qu'ils seront entrés dans mon domaine, ils seront soumis aux règles que ma classification impose et chaque nœud traversé appliquera son niveau de qualité de service ».

En fait, vous pouvez appliquer vos propres politiques dans vos domaines locaux, mais des *autorisations au niveau service* devront être considérées lors de la connexion à d'autres domaines DS.

En ce moment, vous vous posez peut-être beaucoup de questions. DiffServ est plus vaste que ce que j'ai expliqué. En fait, vous pouvez comprendre que je ne peux pas résumer plus de trois RFC en 50 lignes :-).

14.3.4. Travailler avec DSMARK

Comme le spécifie la bibliographie concernant DiffServ, nous différencions les nœuds frontaliers et les nœuds intérieurs. Ce sont deux éléments importants dans le chemin qu'emprunte le trafic. Les deux réalisent une classification quand un paquet arrive. Le résultat peut être utilisé à différents endroits lors du processus DS avant que le paquet ne soit libéré vers le réseau. Cela est possible car le nouveau code DiffServ fournit une structure appelée *sk_buff*, incluant un nouveau champ appelé *skb->tcindex*. Ce champ mémorisera le résultat de la classification initiale et pourra être utilisé à plusieurs reprises dans le traitement DS.

La valeur *skb->tc_index* sera initialement configurée par le gestionnaire de mise en file d'attente DSMARK. Cette valeur sera extraite du champ DS de l'en-tête IP de tous les paquets reçus. En outre, le classificateur *cls_tcindex* lira tout ou une partie de la valeur *skb->tcindex* et l'utilisera pour sélectionner les classes.

Mais, avant tout, regardons la commande **qdisc DSMARK** et ses paramètres :

```
... dsmark indices INDICES [ default_index DEFAULT_INDEX ] [ set_tc_index ]
```

Que signifient ces paramètres ?

- *indices* : taille de la table des couples (masque,valeur). La valeur maximum est 2^n , où $n \geq 0$.
- *default_index* : index d'entrée par défaut de la table si aucune correspondance n'est trouvée.
- *set_tc_index* : indique au gestionnaire DSMARK de récupérer le champs DS et de l'enregistrer dans *skb->tc_index*.

Regardons DSMARK procéder.

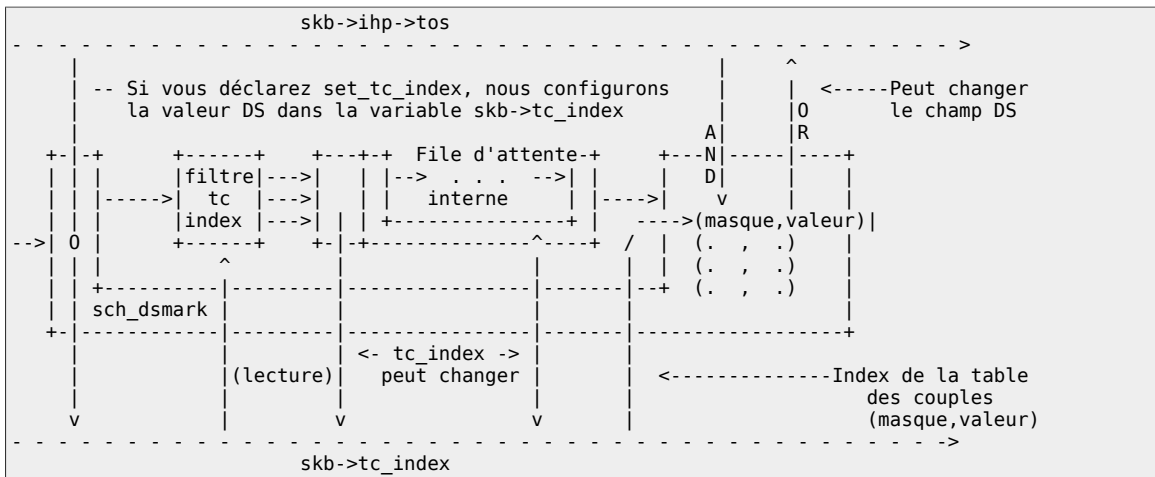
14.3.5. Comment SCH_DSMARK travaille ?

Ce gestionnaire de mise en file d'attente réalisera les étapes suivantes :

- Si vous avez déclaré l'option *set_tc_index* dans la commande **qdisc**, le champ DS est récupéré et mémorisé dans la variable *skb->tc_index*.
- Le classificateur est invoqué. Celui-ci sera exécuté et retournera un identificateur de classe (*class ID*) qui sera enregistré dans la variable *skb->tc_index*. Si aucun filtre correspondant n'est trouvé, nous considérons l'option *default_index* comme étant l'identificateur de classe à enregistrer. Si, ni *set_tc_index*, ni *default_index* n'ont été déclarés, alors les résultats peuvent être non prédictifs.
- Après avoir été envoyé dans le gestionnaire de file d'attente interne, où le résultat du filtre peut être réutilisé, l'identificateur de classe retourné par le gestionnaire est stocké dans la variable *skb->tc_index*. Cette valeur sera utilisée plus tard pour indexer la table masque-valeur. Le résultat de l'opération suivante sera assigné au paquet :

```
Nouveau_champ_DS = ( Ancien_champ_DS & masque ) | valeur
```

- La nouvelle valeur résultera donc d'un ET logique entre les valeurs du champ_DS et du masque, suivi d'un OU logique avec le paramètre valeur. Regardez la figure suivante pour comprendre tout ce processus :



Comment faire le marquage ? Il suffit de modifier le masque et la valeur associés à la classe que vous voulez marquer. Regardez la ligne de code suivante :

```
tc class change dev eth0 classid 1:1 dsmark mask 0x3 value 0xb8
```

Cela modifie le couple (masque,valeur) dans la table de hachage, et re-marque les paquets appartenant à la classe 1:1. Vous devez "changer" ces valeurs en raison des valeurs par défaut que le couple (masque, valeur) obtient initialement (voir le tableau ci-dessous).

Nous allons maintenant expliquer comment le filtre TC_INDEX travaille, et comment il s'intègre dans tout cela. En outre, le filtre TC_INDEX peut être utilisé dans des configurations autres que celles incluant les services DS.

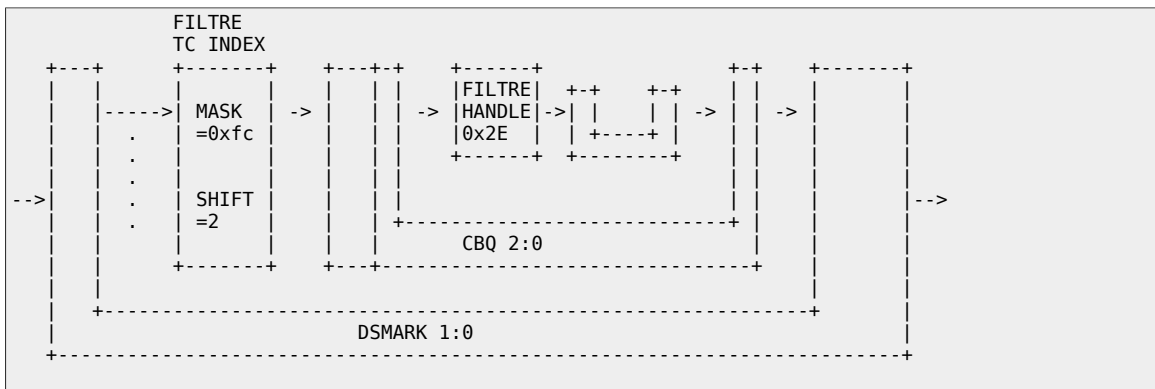
14.3.6. Le filtre TC_INDEX

Voici la commande de base pour déclarer un filtre TC_INDEX :

```
... tcindex [ hash SIZE ] [ mask MASK ] [ shift SHIFT ]
           [ pass_on | fall_through ]
           [ classid CLASSID ] [ police POLICE_SPEC ]
```

Ensuite, nous montrons l'exemple utilisé pour expliquer le mode opératoire de TC_INDEX. Soyez attentif aux mots en gras : `tc qdisc add dev eth0 handle 1:0 root dsmark indices 64 set_tc_index tc filter add dev eth0 parent 1:0 protocol ip prio 1 tcindex mask 0xfc shift 2` `tc qdisc add dev eth0 parent 1:0 handle 2:0 cbq bandwidth 10Mbit cell 8 avpkt 1000 mpu 64 # Classe du trafic EF` `tc class add dev eth0 parent 2:0 classid 2:1 cbq bandwidth 10Mbit rate 1500Kbit avpkt 1000 prio 1 bounded isolated allot 1514 weight 1 maxburst 10 # Gestionnaire de file d'attente fifo pour le trafic EF` `tc qdisc add dev eth0 parent 2:1 pfifo limit 5` `tc filter add dev eth0 parent 2:0 protocol ip prio 1 handle 0x2e tcindex classid 2:1 pass_on` (Ce code n'est pas complet. Ce n'est qu'un extrait de l'exemple EFCBQ inclus dans la distribution iproute2).

Avant tout, supposons que nous recevons un paquet marqué comme EF. Si vous lisez le RFC2598, vous verrez que DSCP recommande une valeur de 101110 pour le trafic EF. Cela signifie que le champ DS sera égal à 10111000 (rappelez-vous que les bits les moins significatifs de l'octet TOS ne sont pas utilisés dans DS) ou 0xb8 en notation hexadécimale.



Le paquet arrive alors avec la valeur du champ DS configurée à 0xb8. Comme je l'ai expliqué auparavant, le gestionnaire de mise en file d'attente dsmark, identifié par 1:0 dans cet exemple, récupère le champ DS et l'enregistre dans la variable `skb->tc_index`. L'étape suivante consistera à associer un filtre à ce gestionnaire de mise en file d'attente (la seconde ligne dans cet exemple). Les opérations suivantes seront réalisées :

```
Valeur1 = skb->tc_index & MASK
Clé = Valeur1 >> SHIFT
```

Dans cet exemple, MASK=0xFC et SHIFT=2.

```
Valeur1 = 10111000 & 11111100 = 10111000
Clé = 10111000 >> 2 = 00101110 -> 0x2E en hexadécimal
```

La valeur retournée correspondra à un identificateur de filtre du gestionnaire de file d'attente interne (dans l'exemple, identifier par 2:0). Si un filtre avec cet identificateur (id) existe, les conditions de contrôle et de performance seront vérifiées (au cas où le filtre inclurait ceci) et l'identificateur de classe sera retourné (dans notre exemple, classid 2:1) et stocké dans la variable `skb->tc_index`.

Si aucun filtre avec cet identificateur n'est trouvé, le résultat dépendra de la déclaration de l'option `fall_through`. Si tel est le cas, la valeur Clé est retournée comme identificateur de classe. Si cela n'est pas le cas, une erreur est retournée et le traitement continue avec les filtres restant. Faites attention si vous utilisez l'option `fall_through` ; ceci ne peut être fait que si une relation existe entre les valeurs de la variable `skb->tc_index` et les identificateurs de classe.

Les derniers paramètres à commenter sont `hash` et `pass_on`. Le premier est relié à la taille de la table de hachage. `Pass_on` sera utilisé pour indiquer d'essayer le filtre suivant dans le cas où aucun identificateur de classe égal au résultat du filtre ne serait trouvé. L'action par défaut est `fall_through` (regarder la table suivante).

Finalement, regardons quelles sont les valeurs possibles pour la configuration de tous ces paramètres TCINDEX :

Nom TC	Valeur	Défaut
Hash	1...0x10000	Dépendant de l'implémentation
Mask	0...0xffff	0xffff
Shift	0...15	0
Fall through / Pass_on	Flag	Fall_through
Classid	Major:minor	Rien
Police	Rien

Ce type de filtre est très puissant. Il est nécessaire d'explorer toutes les possibilités. En outre, ce filtre n'est pas seulement utilisé dans les configurations DiffServ. Vous pouvez l'utiliser comme n'importe quel autre filtre.

Je vous recommande de regarder les exemples DiffServ inclus dans la distribution `iproute2`. Je vous promets que j'essaierai de compléter ce texte dès que possible. Tout ce que j'ai expliqué est le résultat de nombreux tests. Merci de me dire si je me suis trompé quelque part.

14.4. Gestionnaire de mise en file d'attente d'entrée (*Ingress qdisc*)

Tous les gestionnaires de mise en file d'attente dont nous avons discuté jusqu'ici sont des gestionnaires de sortie. Chaque interface peut également avoir un gestionnaire de mise en file d'attente d'entrée qui n'est pas utilisé pour envoyer les paquets à l'extérieur du périphérique réseau. Au lieu de cela, il vous autorise à appliquer des filtres tc aux paquets entrants par l'interface, indépendamment de s'ils ont une destination locale ou s'ils sont destinés à être transmis.

Etant donné que les filtres tc contiennent une implémentation complète du Filtre à Seau de Jetons, et qu'ils sont également capables de s'appuyer sur l'estimation du flux fourni par le noyau, il y a beaucoup de fonctionnalités disponibles. Ceci vous permet de régler le trafic entrant de façon efficace, avant même qu'il n'entre dans la pile IP.

14.4.1. Paramètres & usage

Le gestionnaire de mise en file d'attente d'entrée ne nécessite pas de paramètres. Il diffère des autres gestionnaires dans le fait qu'il n'occupe pas la racine du périphérique. Attachez-le comme ceci :

```
# tc qdisc add dev eth0 ingress
```

Ceci vous autorise à avoir d'autres gestionnaires de sortie sur votre périphérique en plus du gestionnaire d'entrée.

Pour un exemple inventé sur la façon dont le gestionnaire d'entrée peut être utilisé, voir le chapitre Recettes de cuisine.

14.5. Random Early Detection (RED)

Ce chapitre est conçu comme une introduction au routage de dorsales (backbones). Ces liaisons impliquent souvent des bandes passantes supérieures à 100 mégabits/s, ce qui nécessite une approche différente de celle de votre modem ADSL à la maison.

Le comportement normal des files d'attente de routeurs est appelé "tail-drop" (NdT : élimine le reste). Le "tail-drop" consiste à mettre en file d'attente un certain volume de trafic et à éliminer tout ce qui déborde. Ce n'est pas du tout équitable et cela conduit à des retransmissions de synchronisation. Quand une retransmission de synchronisation a lieu, la brusque rafale de rejet d'un routeur qui a atteint sa limite entraînera une rafale de retransmissions retardée qui inondera à nouveau le routeur congestionné.

Dans le but d'en finir avec les congestions occasionnelles des liens, les routeurs de dorsales intègrent souvent des files d'attente de grande taille. Malheureusement, bien que ces files d'attente offrent un bon débit, elles peuvent augmenter sensiblement les temps de latence et entraîner un comportement très saccadé des connexions TCP pendant la congestion.

Ces problèmes avec le "tail-drop" deviennent de plus en plus préoccupants avec l'augmentation de l'utilisation d'applications hostiles au réseau. Le noyau Linux nous offre la technique RED, abréviation de Random Early Detect ou détection précoce directe.

RED n'est pas la solution miracle à tous ces problèmes. Les applications qui n'intègrent pas correctement la technique de "l'exponential backoff" obtiennent toujours une part trop grande de bande passante. Cependant, avec la technique RED elles ne provoquent pas trop de dégâts sur le débit et les temps de latence des autres connexions.

RED élimine statistiquement des paquets du flux avant qu'il n'atteigne sa limite "dure" (hard). Sur une dorsale congestionnée, cela entraîne un ralentissement en douceur de la liaison et évite les retransmissions de synchronisation. La technique RED aide aussi TCP à trouver une vitesse "équitable" plus rapidement : en permettant d'éliminer des paquets plus tôt, il conserve une file d'attente plus courte et des temps de latence mieux contrôlés. La probabilité qu'un paquet soit éliminé d'une connexion particulière est proportionnelle à la bande passante utilisée par cette connexion plutôt qu'au nombre de paquets qu'elle envoie.

La technique RED est une bonne gestion de file d'attente pour les dorsales, où vous ne pouvez pas vous permettre le coût d'une mémorisation d'état par session qui est nécessaire pour une mise en file d'attente vraiment équitable.

Pour utiliser RED, vous devez régler trois paramètres : Min, Max et burst. Min est la taille minimum de la file d'attente en octets avant que les rejets n'aient lieu, Max est le maximum "doux" (soft) en dessous duquel l'algorithme s'efforcera de rester, et burst est le nombre maximum de paquets envoyés "en rafale".

Vous devriez configurer Min en calculant le plus grand temps de latence acceptable pour la mise en file d'attente, multiplié par votre bande passante. Par exemple, sur mon lien ISDN à 64 Kbits/s, je voudrais avoir un temps de latence de base de mise en file d'attente de 200 ms. Je configure donc Min à 1600 octets (= $0,2 \times 64000 / 8$). Imposer une valeur Min trop petite va dégrader le débit et une valeur Min trop grande va dégrader le temps de latence. Sur une liaison lente, choisir un coefficient Min petit ne peut pas remplacer une réduction du MTU pour améliorer les temps de réponse.

Vous devriez configurer Max à au moins deux fois Min pour éviter les synchronisations. Sur des liens lents avec de petites valeurs de Min, il peut être prudent d'avoir Max quatre fois plus grand que Min ou plus.

Burst contrôle la réponse de l'algorithme RED aux rafales. Burst doit être choisi plus grand que min/avpkt (paquet moyen). Expérimentalement, j'ai trouvé que $(\text{min} + \text{min} + \text{max}) / (3 * \text{avpkt})$ marche bien.

De plus, vous devez configurer limit et avpkt. Limit est une valeur de sécurité : s'il y a plus de Limit octets dans la file, RED reprend la technique "tail-drop". Je choisis une valeur typique égale à 8 fois Max. Avpkt devrait être fixé à la taille moyenne d'un paquet. 1000 fonctionne correctement sur des liaisons Internet haut débit ayant un MTU de 1500 octets.

Lire [l'article sur la file d'attente RED](#)⁶ par Sally Floyd et Van Jacobson pour les informations techniques.

14.6. Generic Random Early Detection

Il n'y a pas grand chose de connu sur GRED. Il ressemble à RED avec plusieurs files d'attente internes, celles-ci étant choisies en se basant sur le champ tcindex de Diffserv. Selon une diapositive trouvée [ici](#)⁷, il possède les capacités *Distributed Weighted RED* de Cisco, ainsi que les capacités RIO de Clark.

Chaque file d'attente virtuelle peut avoir ses propres "Drop Parameters".

FIXME: Demandez à Jamal or Werner de nous en dire plus

14.7. Emulation VC/ATM

Ceci est l'effort principal de Werner Almesberger pour vous autoriser à construire des circuits virtuels au-dessus des sockets TCP/IP. Le circuit virtuel est un concept venant de la théorie des réseaux ATM.

Pour plus d'informations, voir la [page ATM sur Linux](#)⁸.

14.8. Weighted Round Robin (WRR)

Ce gestionnaire de mise en file d'attente n'est pas inclus dans les noyaux standards, mais peut être téléchargée à partir de [ce lien](#)⁹. Ce gestionnaire de mise en file d'attente n'a été testé qu'avec les noyaux 2.2, mais marchera probablement également avec les noyaux 2.4/2.5.

La file d'attente WRR partage la bande passante entre ses classes en utilisant la technique du tourniquet pondéré. Ceci est similaire à la file d'attente CBQ qui contient des classes sur lesquelles l'on peut associer arbitrairement des files d'attente. Toutes les classes qui ont suffisamment de demandes obtiendront la bande passante proportionnellement au poids associé des classes. Les poids peuvent être configurés manuellement en utilisant le programme tc. Ils peuvent également être configurés pour décroître automatiquement pour les classes transférant moins de données.

La file d'attente a un classificateur intégré qui assigne les paquets venant ou allant vers différentes machines à différentes classes. On peut utiliser soit l'adresse MAC soit l'adresse IP de l'adresse source ou de destination. L'adresse MAC ne peut cependant être utilisée que quand la boîte Linux est un pont ethernet. Les classes sont automatiquement assignées aux machines en fonction des paquets vus.

⁶ <http://www.aciri.org/floyd/papers/red/red.html>

⁷ <http://www.davin.ottawa.on.ca/ols/img22.htm>

⁸ <http://linux-atm.sourceforge.net/>

⁹ <http://wipl-wrr.dkik.dk/wrr/>

Ce gestionnaire de mise en file d'attente peut être très utile au site comme les résidences étudiantes où des individus sans liens particuliers partagent une connexion Internet. Un ensemble de scripts pour configurer un tel cas de figure pour ce genre de site est proposé dans la distribution WRR.

Cette section contient des « recettes de cuisine » qui peuvent vous aider à résoudre vos problèmes. Un livre de cuisine ne remplace cependant pas une réelle compréhension, essayez donc d'assimiler ce qui suit.

15.1. Faire tourner plusieurs sites avec différentes SLA (autorisations)

Vous pouvez faire cela de plusieurs manières. Apache possède un module qui permet de le supporter, mais nous montrerons comment Linux peut le faire pour d'autres services. Les commandes ont été reprises d'une présentation de Jamal Hadi, dont la référence est fournie ci-dessous.

Disons que nous avons deux clients avec : http, ftp et du streaming audio. Nous voulons leur vendre une largeur de bande passante limitée. Nous le ferons sur le serveur lui-même.

Le client A doit disposer d'au moins 2 mégabits, et le client B a payé pour 5 mégabits. Nous séparons nos clients en créant deux adresses IP virtuelles sur notre serveur.

```
# ip address add 188.177.166.1 dev eth0
# ip address add 188.177.166.2 dev eth0
```

C'est à vous d'associer les différents serveurs à la bonne adresse IP. Tous les démons courants supportent cela.

Nous pouvons tout d'abord attacher une mise en file d'attente CBQ à eth0 :

```
# tc qdisc add dev eth0 root handle 1: cbq bandwidth 10Mbit cell 8 avpkt 1000 \
mpu 64
```

Nous créons ensuite les classes pour nos clients :

```
# tc class add dev eth0 parent 1:0 classid 1:1 cbq bandwidth 10Mbit rate \
  2Mbit avpkt 1000 prio 5 bounded isolated allot 1514 weight 1 maxburst 21
# tc class add dev eth0 parent 1:0 classid 1:2 cbq bandwidth 10Mbit rate \
  5Mbit avpkt 1000 prio 5 bounded isolated allot 1514 weight 1 maxburst 21
```

Nous ajoutons les filtres pour nos deux classes :

```
##FIXME: Pourquoi cette ligne, que fait-elle ? Qu'est-ce qu'un
diviseur ?
##FIXME: Un diviseur est lié à une table de hachage et au nombre de
seaux -ahu
# tc filter add dev eth0 parent 1:0 protocol ip prio 5 handle 1: u32 divisor 1
# tc filter add dev eth0 parent 1:0 prio 5 u32 match ip src 188.177.166.1
  flowid 1:1
# tc filter add dev eth0 parent 1:0 prio 5 u32 match ip src 188.177.166.2
  flowid 1:2
```

Et voilà qui est fait.

FIXME: Pourquoi pas un filtre token bucket ? Y a-t-il un retour par défaut à pfifo_fast quelque part ?

15.2. Protéger votre machine des inondations SYN

D'après la documentation iproute d'Alexey adaptée à netfilter. Si vous utilisez ceci, prenez garde d'ajuster les nombres avec des valeurs raisonnables pour votre système.

Si vous voulez protéger tout un réseau, oubliez ce script, qui est plus adapté à un hôte seul.

Il apparaît que la toute dernière version de l'outil iproute2 est nécessaire pour que ceci fonctionne avec le noyau 2.4.0.

```
#!/bin/sh -x
#
# script simple utilisant les capacités de Ingress.
# Ce script montre comment on peut limiter le flux entrant des SYN.
# Utile pour la protection des TCP-SYN. Vous pouvez utiliser IPchains
# pour bénéficier de puissantes fonctionnalités sur les SYN.
#
# chemins vers les divers utilitaires
# À changer en fonction des vôtres
#
TC=/sbin/tc
IP=/sbin/ip
IPTABLES=/sbin/iptables
INDEV=eth2
#
# marque tous les paquets SYN entrant à travers $INDEV avec la valeur 1
#####
$IPTABLES -A PREROUTING -i $INDEV -t mangle -p tcp --syn \
  -j MARK --set-mark 1
#####
#
# installe la file d'attente ingress sur l'interface associée
#####
$TC qdisc add dev $INDEV handle ffff: ingress
#####
#
# Les paquets SYN ont une taille de 40 octets (320 bits), donc trois SYN
# ont une taille de 960 bits (approximativement 1Kbit) ; nous limitons donc
# les SYN entrants à 3 par seconde (pas vraiment utile, mais sert à
# montrer ce point -JHS
```

```
#####
$TC filter add dev $INDEV parent ffff: protocol ip prio 50 handle 1 fw \
police rate 1kbit burst 40 mtu 9k drop flowid :1
#####

#
echo "---- qdisc parameters Ingress -----"
$TC qdisc ls dev $INDEV
echo "---- Class parameters Ingress -----"
$TC class ls dev $INDEV
echo "---- filter parameters Ingress -----"
$TC filter ls dev $INDEV parent ffff:

#supprime la file d'attente ingress
#$TC qdisc del $INDEV ingress
```

15.3. Limiter le débit ICMP pour empêcher les dénis de service

Récemment, les attaques distribuées de déni de service sont devenues une nuisance importante sur Internet. En filtrant proprement et en limitant le débit de votre réseau, vous pouvez à la fois éviter de devenir victime ou source de ces attaques.

Vous devriez filtrer vos réseaux de telle sorte que vous n'autorisiez pas les paquets avec une adresse IP source non locale à quitter votre réseau. Cela empêche les utilisateurs d'envoyer de manière anonyme des cochonneries sur Internet.

La limitation de débit peut faire encore mieux, comme vu plus haut. Pour vous rafraîchir la mémoire, revoici notre diagramme ASCII :

```
[Internet] ---<E3, T3, n'importe quoi>--- [routeur Linux] --- [Bureau+FAI]
                                     eth1           eth0
```

Nous allons d'abord configurer les parties pré-requises :

```
# tc qdisc add dev eth0 root handle 10: cbq bandwidth 10Mbit avpkt 1000
# tc class add dev eth0 parent 10:0 classid 10:1 cbq bandwidth 10Mbit rate \
  10Mbit allot 1514 prio 5 maxburst 20 avpkt 1000
```

Si vous avez des interfaces de 100 Mbits ou plus, ajustez ces nombres. Maintenant, vous devez déterminer combien de trafic ICMP vous voulez autoriser. Vous pouvez réaliser des mesures avec tcpdump, en écrivant les résultats dans un fichier pendant un moment, et regarder combien de paquets ICMP passent par votre réseau. Ne pas oublier d'augmenter la longueur du "snapshot". Si la mesure n'est pas possible, vous pouvez consacrer par exemple 5% de votre bande passante disponible. Configurons notre classe :

```
# tc class add dev eth0 parent 10:1 classid 10:100 cbq bandwidth 10Mbit rate \
  100Kbit allot 1514 weight 800Kbit prio 5 maxburst 20 avpkt 250 \
  bounded
```

Cela limite le débit à 100 Kbits sur la classe. Maintenant, nous avons besoin d'un filtre pour assigner le trafic ICMP à cette classe :

```
# tc filter add dev eth0 parent 10:0 protocol ip prio 100 u32 match ip
  protocol 1 0xFF flowid 10:100
```

15.4. Donner la priorité au trafic interactif

Si beaucoup de données arrivent à votre lien ou en partent, et que vous essayez de faire de la maintenance via telnet ou ssh, cela peut poser problème : d'autres paquets bloquent vos frappes clavier. Cela ne serait-il pas mieux si vos paquets interactifs pouvaient se faufiler dans le trafic de masse ? Linux peut faire cela pour vous.

Comme précédemment, nous avons besoin de manipuler le trafic dans les deux sens. Evidemment, cela marche mieux s'il y a des machines Linux aux deux extrémités du lien, bien que d'autres UNIX soient capables de faire la même chose. Consultez votre gourou local Solaris/BSD pour cela.

Le gestionnaire standard pfifo_fast a trois "bandes" différentes. Le trafic de la bande 0 est transmis en premier, le trafic des bandes 1 et 2 étant traité après. Il est vital que votre trafic interactif soit dans la bande 0 ! Ce qui suit est adapté du (bientôt obsolète) Ipchains-HOWTO :

Il y a quatre bits rarement utilisés dans l'en-tête IP, appelés bits de Type de Service (TOS). Ils affectent la manière dont les paquets sont traités. Les quatre bits sont "Délai Minimum", "Débit Maximum", "Fiabilité Maximum" et "Coût Minimum". Seul un de ces bits peut être positionné. Rob van Nieuwkerk, l'auteur du code TOS-mangling dans ipchains, le configure comme suit :

```
Le "Délai Minimum" est particulièrement important pour moi. Je le
positionne à 1 pour les paquets interactifs sur mon routeur (Linux)
qui envoie le trafic vers l'extérieur. Je suis derrière un modem à
33,6 Kbps. Linux répartit les paquets dans trois files
d'attente. De cette manière, j'obtiens des performances acceptables
pour le trafic interactif tout en téléchargeant en même temps.
```

L'utilisation la plus commune est de configurer les connexions telnet et ftp à "Délai Minimum" et les données FTP à "Débit Maximum". Cela serait fait comme suit, sur mon routeur :

```
# iptables -A PREROUTING -t mangle -p tcp --sport telnet \
```

```
-j TOS --set-tos Minimize-Delay
# iptables -A PREROUTING -t mangle -p tcp --sport ftp \
-j TOS --set-tos Minimize-Delay
# iptables -A PREROUTING -t mangle -p tcp --sport ftp-data \
-j TOS --set-tos Maximize-Throughput
```

En fait, cela ne marche que pour les données venant d'un telnet extérieur vers votre ordinateur local. Dans l'autre sens, ça se fait tout seul : telnet, ssh, et consorts configurent le champ TOS automatiquement pour les paquets sortants.

Si vous avez un client incapable de le faire, vous pouvez toujours le faire avec netfilter. Sur votre machine locale :

```
# iptables -A OUTPUT -t mangle -p tcp --dport telnet \
-j TOS --set-tos Minimize-Delay
# iptables -A OUTPUT -t mangle -p tcp --dport ftp \
-j TOS --set-tos Minimize-Delay
# iptables -A OUTPUT -t mangle -p tcp --dport ftp-data \
-j TOS --set-tos Maximize-Throughput
```

15.5. Cache web transparent utilisant netfilter, iproute2, ipchains et squid

Cette section a été envoyée par le lecteur Ram Narula de "Internet for Education" (Internet pour l'éducation) (Thaïlande).

La technique habituelle pour réaliser ceci dans Linux est probablement l'utilisation d'ipchains APRES s'être assuré que le trafic sortant du port 80 (web) est routé à travers le serveur faisant fonctionner squid.

Il y a 3 méthodes communes pour être sûr que le trafic sortant du port 80 est routé vers le serveur faisant fonctionner squid et une quatrième est introduite ici.

La passerelle le fait.

Si vous pouvez dire à votre passerelle que les paquets sortants à destination du port 80 doivent être envoyés vers l'adresse IP du serveur squid.

MAIS

Ceci amènerait une charge supplémentaire sur le routeur et des routeurs commerciaux peuvent même ne pas supporter ceci.

Utiliser un commutateur Couche 4.

Les commutateurs Couche 4 peuvent manipuler ceci sans aucun problème.

MAIS

Le coût pour un tel équipement est en général très élevé. Typiquement, un commutateur couche 4 coûte normalement plus cher qu'un serveur classique + un bon serveur linux.

Utiliser le serveur cache comme passerelle réseau

Vous pouvez forcer TOUT le trafic à travers le serveur cache

MAIS

Ceci est plutôt risqué dans la mesure où Squid utilise beaucoup de ressources CPU, ce qui peut conduire à une baisse des performances de tout le réseau. Le serveur peut également ne plus fonctionner et personne sur le réseau ne pourra accéder à Internet si cela a lieu.

Routeur Linux+NetFilter.

En utilisant Netfilter, une autre technique peut être implémentée. Celle-ci consiste à utiliser Netfilter pour "marquer" les paquets à destination du port 80 et à utiliser iproute2 pour router les paquets "marqués" vers le serveur Squid.

```
-----|
| Implémentation |
|-----|

Adresses utilisées
10.0.0.1 naret (serveur NetFilter)
10.0.0.2 silom (serveur Squid)
10.0.0.3 donmuang (routeur connecté à Internet)
10.0.0.4 kaosarn (un autre serveur sur le réseau)
10.0.0.5 RAS
10.0.0.0/24 réseau principal
10.0.0.0/19 réseau total

-----|
| Schéma du réseau |
|-----|

Internet
|
donmuang
|
-----hub/commutateur-----
|                               |
```

```
naret silom kaosarn RAS etc.
```

Premièrement, faire en sorte que tout le trafic passe par naret en étant sûr que c'est la passerelle par défaut, à l'exception de silom. La passerelle par défaut de silom doit être donmuang (10.0.0.3) ou ceci créerait une boucle du trafic web.

(Tous les serveurs sur mon réseau avaient 10.0.0.1 comme passerelle par défaut qui était l'ancienne adresse du routeur donmuang. Cela m'a conduit à attribuer 10.0.0.3 comme adresse IP à donmuang et à donner 10.0.0.1 comme adresse IP à naret.)

```
Silom
-----
-configurer squid et ipchains
```

Pour la configuration du serveur Squid sur silom, soyez sûr que celui-ci supporte le cache/proxy transparent (transparent caching/proxying). Le port par défaut pour squid est en général 3128. Tout le trafic pour le port 80 doit donc être redirigé localement vers le port 3128. Ceci peut être fait en utilisant ipchains comme suit :

```
silom# ipchains -N allow1
silom# ipchains -A allow1 -p TCP -s 10.0.0.0/19 -d 0/0 80 -j REDIRECT 3128
silom# ipchains -I input -j allow1
```

Ou, avec netfilter:

```
silom# iptables -t nat -A PREROUTING -i eth0 -p tcp --dport 80 -j REDIRECT --to-port 3128
```

(note: vous pouvez avoir également d'autres entrées)

Pour plus d'informations sur la configuration du serveur Squid, se référer à la page FAQ Squid sur <http://squid.nlanr.net>).

Soyez sûr que "ip forwarding" est actif sur votre serveur et que la passerelle par défaut pour ce serveur est donmuang (PAS naret).

```
Naret
-----
- configurer squid et ipchains
- désactiver les messages icmp REDIRECT (si besoin)
```

1. "Marquer" les paquets à destination du port 80 avec la valeur 2

```
naret# iptables -A PREROUTING -i eth0 -t mangle -p tcp --dport 80 \
-j MARK --set-mark 2
```

2. Configurer iproute2 de sorte qu'il route les paquets avec la marque 2 vers silom

```
naret# echo 202 www.out >> /etc/iproute2/rt_tables
naret# ip rule add fwmark 2 table www.out
naret# ip route add default via 10.0.0.2 dev eth0 table www.out
naret# ip route flush cache
```

Si donmuang et naret sont sur le même réseau, naret ne doit pas envoyer de messages icmp REDIRECT. Ceux-ci doivent être désactivés par :

```
naret# echo 0 > /proc/sys/net/ipv4/conf/all/send_redirects
naret# echo 0 > /proc/sys/net/ipv4/conf/default/send_redirects
naret# echo 0 > /proc/sys/net/ipv4/conf/eth0/send_redirects
```

La configuration est terminée, vérifions-la.

```
Sur naret:

naret# iptables -t mangle -L
Chain PREROUTING (policy ACCEPT)
target      prot opt source                destination
MARK        tcp  --  anywhere              anywhere            tcp dpt:www MARK set 0x2

Chain OUTPUT (policy ACCEPT)
target      prot opt source                destination

naret# ip rule ls
0:          from all lookup local
32765:     from all fwmark      2 lookup www.out
32766:     from all lookup main
32767:     from all lookup default

naret# ip route list table www.out
default via 203.114.224.8 dev eth0

naret# ip route
10.0.0.1 dev eth0 scope link
10.0.0.0/24 dev eth0 proto kernel scope link src 10.0.0.1
127.0.0.0/8 dev lo scope link
default via 10.0.0.3 dev eth0

(soyez sûr que silom appartient à l'une des lignes ci-dessus. Dans ce cas,
c'est la ligne avec 10.0.0.0/24)
```


Récemment, j'ai séparé le serveur du routeur de sorte que la plupart des applications fonctionnent sur une machine différente de celle qui réalise le routage.

J'ai alors eu des problèmes en me connectant sur l'irc. Grosse panique ! Je vous assure que certains essais trouvaient que j'étais connecté à l'irc, me montrant même comme connecté sur l'irc mais je ne recevais pas le "motd" (message of the day, message du jour) de l'irc. J'ai vérifié ce qui pouvait être erroné et ai noté que j'avais déjà eu des soucis liés au MTU en contactant certains sites web. Je n'avais aucun souci pour les atteindre quand le MTU était à 1500, le problème n'apparaissant que lorsque le MTU était configuré à 296. Puisque les serveurs irc bloquent tout le trafic dont il n'ont pas besoin pour leurs opérations immédiates, ils bloquent aussi l'icmp.

J'ai manoeuvré pour convaincre les responsables d'un serveur web que ceci était la cause d'un problème, mais les responsables du serveur irc n'avaient pas l'intention de réparer ceci.

Donc, je devais être sûr que le trafic masqué sortant partait avec le mtu faible du lien externe. Mais, je voulais que le trafic ethernet local ait le MTU normal (pour des choses comme le trafic nfs).

Solution :

```
ip route add default via 10.0.0.1 mtu 296
```

(10.0.0.1 étant ma passerelle par défaut, l'adresse interne de mon routeur masquant)

En général, il est possible d'outrepasser la découverte du MTU de chemin en configurant des routes spécifiques. Par exemple, si seuls certains réseaux posent problèmes, ceci devrait aider :

```
ip route add 195.96.96.0/24 via 10.0.0.1 mtu 1000
```

15.7. Circonvenir aux problèmes de la découverte du MTU de chemin en imposant le MSS (pour les utilisateurs de l'ADSL, du câble, de PPPoE & PPTP)

Comme expliqué au-dessus, la découverte du MTU de chemin ne marche pas aussi bien que cela devrait être. Si vous savez qu'un saut de votre réseau a un MTU limité (<1500), vous ne pouvez pas compter sur la découverte du MTU de chemin pour le découvrir.

Outre le MTU, il y a encore un autre moyen de configurer la taille maximum du paquet, par ce qui est appelé le MSS (Maximum Segment Size, Taille Maximum du Segment). C'est un champ dans les options TCP du paquet SYN.

Les noyaux Linux récents, et quelques pilotes de périphérique PPPoE (notamment, l'excellent Roaring Penguin) implémentent la possibilité de 'fixer le MSS'.

Le bon côté de tout ceci est que, en positionnant la valeur MSS, vous dites à l'hôte distant de manière équivoque "n'essaie pas de m'envoyer des paquets plus grands que cette valeur". Aucun trafic ICMP n'est nécessaire pour faire fonctionner cela.

Malheureusement, c'est de la bidouille évidente -- ça détruit la propriété «bout-en-bout» de la connexion en modifiant les paquets. Ayant dit cela, nous utilisons cette astuce dans beaucoup d'endroit et cela fonctionne comme un charme.

Pour que tout ceci fonctionne, vous aurez besoin au moins de iptables-1.2.1a et de Linux 2.4.3 ou plus. La ligne de commande basique est :

```
# iptables -A FORWARD -p tcp --tcp-flags SYN,RST SYN -j TCPMSS --clamp-mss-to-pmtu
```

Ceci calcule le MSS approprié pour votre lien. Si vous vous sentez courageux ou que vous pensez être le mieux placé pour juger, vous pouvez aussi faire quelque chose comme ceci :

```
# iptables -A FORWARD -p tcp --tcp-flags SYN,RST SYN -j TCPMSS --set-mss 128
```

Ceci configure le MSS du paquet SYN à 128. Utilisez ceci si vous avez de la voix sur IP (VoIP) avec de tous petits paquets, et de grands paquets http qui provoquent des coupures dans vos communications vocales.

15.8. Le Conditionneur de Trafic Ultime : Faible temps de latence, Téléchargement vers l'amont et l'aval rapide

Note : ce script a récemment été mis à niveau et il ne marchait avant qu'avec les clients Linux de votre réseau ! Vous devriez donc le mettre à jour si vous avez des machines Windows ou des Macs dans votre réseau qui n'étaient pas capables de télécharger plus rapidement pendant que d'autres étaient en train de télécharger vers l'amont.

J'ai essayé de créer le Saint Graal :

Maintenir à tout moment un faible temps de latence pour le trafic interactif

Ceci signifie que la récupération ou la transmission de fichiers ne doivent pas perturber SSH ou même telnet. Ceci est la chose la plus importante, car même un temps de latence de 200ms est important pour pouvoir travailler confortablement.

Autoriser le 'surf' à des vitesses raisonnables pendant que l'on télécharge vers l'amont ou vers l'aval
Même si http est un trafic de masse, les autres trafics ne doivent pas trop le noyer.

Etre sûr que le téléchargement vers l'amont ne va pas faire du tort aux téléchargements vers l'aval et aux autres éléments autour

Le principal phénomène observé est la forte réduction de la vitesse de téléchargement vers l'aval quand il y a du trafic montant.

Il s'avère que tout ceci est possible, au prix d'une minuscule réduction de la bande passante. La présence de grandes files d'attente sur les dispositifs d'accès domestiques, comme le câble ou les modems DSL, explique pourquoi les téléchargements vers l'amont, vers l'aval et ssh se pénalisent mutuellement.

La prochaine partie explique en profondeur ce qui provoque les retards, et comment nous pouvons les corriger. Vous pouvez sans danger la passer et aller directement au script si vous vous fichez de la façon dont la magie opère.

15.8.1. Pourquoi cela ne marche t-il pas bien par défaut ?

Les FAI savent que leurs performances ne sont seulement jugées que sur la vitesse à laquelle les personnes peuvent télécharger vers l'aval. En plus de la bande passante disponible, la vitesse de téléchargement est lourdement influencée par la perte des paquets, qui gêne sérieusement les performances de TCP/IP. Les grandes files d'attente peuvent aider à prévenir la perte des paquets, et augmenter les téléchargements. Les FAI configurent donc de grandes files d'attente.

Ces grandes files d'attente endommagent cependant l'interactivité. Une frappe doit d'abord parcourir la file d'attente du flux montant, ce qui peut prendre plusieurs secondes, et aller jusqu'à l'hôte distant. Elle est alors traitée, conduisant à un paquet de retour qui doit traverser la file d'attente du flux descendant, localisée chez votre FAI, avant qu'elle n'apparaisse sur l'écran.

Cet HOWTO nous enseigne plusieurs manières de modifier et traiter la file d'attente mais, malheureusement, toutes les files d'attente ne nous sont pas accessibles. Les files d'attente du FAI sont sans limites et la file d'attente du flux montant réside probablement dans votre modem câble ou votre périphérique DSL. Il se peut que vous soyez capable ou non de le configurer. La plupart du temps, ce ne sera pas le cas.

Et après ? Etant donné que nous ne pouvons pas contrôler ces files d'attente, elles doivent disparaître et être transférées sur notre routeur Linux. Heureusement, ceci est possible.

Limiter la vitesse de téléchargement vers l'amont (upload)

En limitant notre vitesse de téléchargement vers l'amont à une vitesse légèrement plus faible que la vitesse réelle disponible, il n'y a pas de files d'attente mises en place dans notre modem. La file d'attente est maintenant transférée vers Linux.

Limiter la vitesse de téléchargement vers l'aval (download)

Ceci est légèrement plus rusé dans la mesure où nous ne pouvons pas vraiment influencer la vitesse à laquelle Internet nous envoie les données. Nous pouvons cependant rejeter les paquets qui arrivent trop vite, ce qui provoque le ralentissement de TCP/IP jusqu'au débit désiré. Comme nous ne voulons pas supprimer inutilement du trafic, nous configurons une vitesse de rafale ('burst') plus grande.

Maintenant que nous avons fait ceci, nous avons éliminé totalement la file d'attente du flux descendant (sauf pour de courtes rafales de données), et obtenu la possibilité de gérer la file d'attente du flux montant avec toute la puissance Linux.

Il nous reste à nous assurer que le trafic interactif se retrouve au début de la file d'attente du flux montant. Pour être sûr que le flux montant ne va pas pénaliser le flux descendant, nous déplaçons également les paquets ACK au début de la file d'attente. C'est ce qui provoque normalement un énorme ralentissement quand du trafic de masse est généré dans les deux sens. Les paquets ACK du trafic descendant rentrent en concurrence avec le trafic montant et sont donc ralentis.

Si nous avons fait tout ceci, nous obtenons les mesures suivantes en utilisant l'excellente connexion ADSL de xs4all, en Hollande :

```
Temps de latence de base :
round-trip min/avg/max = 14.4/17.1/21.7 ms

Sans le conditionneur de trafic, lors d'un téléchargement vers l'aval :
round-trip min/avg/max = 560.9/573.6/586.4 ms

Sans le conditionneur de trafic, lors d'un téléchargement vers l'amont :
round-trip min/avg/max = 2041.4/2332.1/2427.6 ms

Avec le conditionneur, lors d'un téléchargement vers l'amont à 220kbit/s :
round-trip min/avg/max = 15.7/51.8/79.9 ms

Avec le conditionneur, lors d'un téléchargement vers l'aval à 850kbit/s :
round-trip min/avg/max = 20.4/46.9/74.0 ms

Lors d'un téléchargement vers l'amont, les téléchargements vers l'aval s'effectuent à environ
80 % de la vitesse maximale disponible et 90% pour les téléchargements vers
l'amont. Le temps de latence augmente alors jusqu'à 850 ms ; je n'ai pas encore
déterminé la raison de ce phénomène.
```

Ce que vous pouvez attendre de ce script dépend largement de votre vitesse de lien réelle. Quand vous téléchargez vers l'amont à pleine vitesse, il y aura toujours un paquet devant votre frappe de clavier. Ceci est

la limite basse de votre temps de latence. Pour la calculer, divisez votre MTU par la vitesse du flux montant. Les valeurs classiques seront un peu plus élevées que ça. Diminuez votre MTU pour un meilleur effet !

Voici deux versions de ce script, une avec l'excellent HTB de Devik, et l'autre avec CBQ qui est présent dans chaque noyau Linux, contrairement à HTB. Les deux ont été testés et marchent correctement.

15.8.2. Le script (CBQ)

Marche avec tous les noyaux. A l'intérieur du gestionnaire de mise en file d'attente CBQ, nous plaçons deux SFQ pour être sûr que de multiples flux de masse ne vont pas mutuellement se pénaliser.

Le trafic descendant est réglementé en utilisant un filtre tc contenant un Token Bucket Filter.

Vous pourriez améliorer ce script en ajoutant 'bounded' aux lignes qui démarrent avec 'tc class add .. classid 1:20'. Si vous avez diminué votre MTU, diminuez aussi les nombres allot & avpkt !

```
#!/bin/bash

# La configuration ultime pour votre connexion Internet domestique
#
# Configurez les valeurs suivantes avec des valeurs légèrement inférieures que
# vos vitesses de flux montant et descendant. Exprimé en kilobits.
DOWNLINK=800
UPLINK=220
DEV=ppp0

# Nettoie les gestionnaires de sortie et d'entrés, cache les erreurs
tc qdisc del dev $DEV root 2> /dev/null > /dev/null
tc qdisc del dev $DEV ingress 2> /dev/null > /dev/null

##### Flux montant (uplink)

# installe CBQ à la racine

tc qdisc add dev $DEV root handle 1: cbq avpkt 1000 bandwidth 10mbit

# Le trafic est mis en forme à une vitesse de $UPLINK. Ceci évite
# d'énormes files d'attente dans votre modem DSL qui pénalisent le temps de
# latence.
# Classe principale

tc class add dev $DEV parent 1: classid 1:1 cbq rate ${UPLINK}kbit \
allot 1500 prio 5 bounded isolated

# classe de priorité supérieure 1:10:

tc class add dev $DEV parent 1:1 classid 1:10 cbq rate ${UPLINK}kbit \
allot 1600 prio 1 avpkt 1000

# la classe par défaut et pour le trafic de masse 1:20. Reçoit légèrement
# moins que le trafic et a une priorité plus faible :
# bulk and default class 1:20 - gets slightly less traffic,
# and a lower priority:

tc class add dev $DEV parent 1:1 classid 1:20 cbq rate [9*$UPLINK/10]kbit \
allot 1600 prio 2 avpkt 1000

# Les deux sont gérées par SFQ :
tc qdisc add dev $DEV parent 1:10 handle 10: sfq perturb 10
tc qdisc add dev $DEV parent 1:20 handle 20: sfq perturb 10

# Démarrage des filtres
# le bit Délai Minimum du champ TOS (ssh, PAS scp) est dirigé vers
# 1:10 :
tc filter add dev $DEV parent 1:0 protocol ip prio 10 u32 \
match ip tos 0x10 0xff flowid 1:10
# ICMP (ip protocol 1) est dirigé vers la classe interactive 1:10 de telle
# sorte que nous pouvons réaliser des mesures et impressionner nos
# amis :
tc filter add dev $DEV parent 1:0 protocol ip prio 11 u32 \
match ip protocol 1 0xff flowid 1:10

# Pour accélérer les téléchargements vers l'aval lors de la présence d'un
# flux montant, les paquets ACK sont placés dans la classe
# interactive :

tc filter add dev $DEV parent 1: protocol ip prio 12 u32 \
match ip protocol 6 0xff \
match u8 0x05 0x0f at 0 \
match u16 0x0000 0xffc0 at 2 \
match u8 0x10 0xff at 33 \
flowid 1:10

# Le reste est considéré 'non-interactif' cad 'de masse' et fini dans 1:20

tc filter add dev $DEV parent 1: protocol ip prio 13 u32 \
match ip dst 0.0.0.0/0 flowid 1:20

##### Flux descendant (downlink) #####
```

```
# Ralentir le flux descendant à une valeur légèrement plus faible que votre
# vitesse réelle de manière à éviter la mise en file d'attente chez notre
# FAI. Faites des tests pour voir la vitesse maximum à laquelle vous pouvez
# le configurer. Les FAI ont tendance à avoir *d'énormes* files d'attente
# pour s'assurer de la rapidité des gros téléchargements.
#
# attache la réglementation d'entrée (ingress policer) :

tc qdisc add dev $DEV handle ffff: ingress

# Filtre *tout* (0.0.0.0/0), rejette tout ce qui arrive trop
# rapidement :

tc filter add dev $DEV parent ffff: protocol ip prio 50 u32 match ip src \
  0.0.0.0/0 police rate ${DOWNLINK}kbit burst 10k drop flowid :1
```

Si vous voulez que ce script soit exécuté par ppp à la connexion, copiez-le dans `/etc/ppp/ip-up.d`.

Si les deux dernières lignes conduisent à une erreur, mettez à jour l'outil tc avec la dernière version !

15.8.3. Le script (HTB)

Le script suivant nous permet d'atteindre tous nos buts en utilisant la merveilleuse file d'attente HTB. Voir le chapitre correspondant. Cela vaut la peine de mettre à jour votre noyau !

```
#!/bin/bash

# La configuration ultime pour votre connexion Internet domestique
#
# Configurez les valeurs suivantes avec des valeurs légèrement inférieures que
# vos vitesses de flux montant et descendant. Exprimé en kilobits.
DOWNLINK=800
UPLINK=220
DEV=ppp0

#Nettoie les gestionnaires de sortie et d'entrés, cache les erreurs
tc qdisc del dev $DEV root 2> /dev/null > /dev/null
tc qdisc del dev $DEV ingress 2> /dev/null > /dev/null

##### Flux montant (uplink)

# installe HTB à la racine, le trafic ira par défaut vers 1:20 :

tc qdisc add dev $DEV root handle 1: htb default 20

# Le trafic est mis en forme à une vitesse de $UPLINK. Ceci évite
# d'énormes files d'attente dans votre modem DSL qui pénalisent le temps de
# latence.

tc class add dev $DEV parent 1: classid 1:1 htb rate ${UPLINK}kbit burst 6k

# la classe de haute priorité 1:10 :

tc class add dev $DEV parent 1:1 classid 1:10 htb rate ${UPLINK}kbit \
  burst 6k prio 1

# bulk & default class 1:20 - gets slightly less traffic,
# and a lower priority:

tc class add dev $DEV parent 1:1 classid 1:20 htb rate [9*$UPLINK/10]kbit \
  burst 6k prio 2

# Les deux sont gérées par SFQ :
tc qdisc add dev $DEV parent 1:10 handle 10: sfq perturb 10
tc qdisc add dev $DEV parent 1:20 handle 20: sfq perturb 10

# le bit Délai Minimum du champ TOS (ssh, PAS scp) est dirigé vers
# 1:10 :
tc filter add dev $DEV parent 1:0 protocol ip prio 10 u32 \
  match ip tos 0x10 0xff flowid 1:10

# ICMP (ip protocol 1) est dirigé vers la classe interactive 1:10 de telle
# sorte que nous pouvons réaliser des mesures et impressionner nos
# amis :
tc filter add dev $DEV parent 1:0 protocol ip prio 10 u32 \
  match ip protocol 1 0xff flowid 1:10

# Pour accélérer les téléchargements vers l'aval lors de la présence d'un
# flux montant, les paquets ACK sont placés dans la classe
# interactive :

tc filter add dev $DEV parent 1: protocol ip prio 10 u32 \
  match ip protocol 6 0xff \
  match u8 0x05 0x0f at 0 \
  match u16 0x0000 0xffc0 at 2 \
  match u8 0x10 0xff at 33 \
  flowid 1:10

# Le reste est considéré 'non-interactif' cad 'de masse' et fini dans 1:20
```

```
##### Flux descendant (downlink) #####
# Ralentir le flux descendant à une valeur légèrement plus faible que votre
# vitesse réelle de manière à éviter la mise en file d'attente chez notre
# FAI. Faites des tests pour voir la vitesse maximum à laquelle vous pouvez
# le configurer. Les FAI ont tendance à avoir *d'énormes* files d'attente
# pour s'assurer de la rapidité des gros téléchargements.
#
# attache la réglementation d'entrée (ingress policer) :

tc qdisc add dev $DEV handle ffff: ingress

# Filtre *tout* (0.0.0.0/0), rejette tout ce qui arrive trop
# rapidement :

tc filter add dev $DEV parent ffff: protocol ip prio 50 u32 match ip src \
  0.0.0.0/0 police rate ${DOWNLINK}kbit burst 10k drop flowid :1
```

Si vous voulez que ce script soit exécuté par ppp à la connexion, copiez-le dans /etc/ppp/ip-up.d.

Si les deux dernières lignes conduisent à une erreur, mettez à jour l'outil tc avec la dernière version !

15.9. Limitation du débit pour un hôte ou un masque de sous-réseau

Bien que ceci soit décrit en détail ailleurs et dans nos pages de manuel, cette question est souvent posée. Heureusement, il y a une réponse simple qui ne nécessite pas la compréhension complète du contrôle de trafic.

Ce script de trois lignes met en place la limitation du débit :

```
tc qdisc add dev $DEV root handle 1: cbq avpkt 1000 bandwidth 10mbit

tc class add dev $DEV parent 1: classid 1:1 cbq rate 512kbit \
  allot 1500 prio 5 bounded isolated

tc filter add dev $DEV parent 1: protocol ip prio 16 u32 \
  match ip dst 195.96.96.97 flowid 1:1
```

La première ligne installe un gestionnaire de mise en file d'attente basé sur des classes sur votre interface, et indique au noyau que, pour ses calculs, il peut la considérer comme une interface à 10 Mbits/s. Cependant, il n'y aura pas de grands dommages si vous indiquez une valeur erronée. Donner la vraie valeur permettra d'avoir des choses plus précises.

La seconde ligne crée une classe de 512kbit/s avec des valeurs par défaut raisonnables. Pour les détails, voir les pages de manuel cbq et [Chapitre 9, Gestionnaires de mise en file d'attente pour l'administration de la bande passante](#).

La dernière ligne indique quel trafic devra passer par la classe réalisant la mise en forme. Le trafic qui n'est sélectionné par cette règle n'est PAS mis en forme. Pour avoir des sélections plus compliquées (sous-réseaux, ports sources ou de destinations), voir [Section 9.6.2, « Toutes les commandes de filtres dont vous aurez normalement besoin »](#).

Si vous avez changé quelque chose et que vous vouliez recharger le script, exécutez la commande **tc qdisc del dev \$DEV root** pour supprimer votre configuration actuelle.

Le script peut être amélioré en ajoutant une dernière ligne optionnelle **tc qdisc add dev \$DEV parent 1:1 sfq perturb 10**. Voir [Section 9.2.3, « Mise en file d'attente stochastiquement équitable \(Stochastic Fairness Queueing\) »](#) pour savoir ce que cela fait.

15.10. Exemple d'une solution de traduction d'adresse avec de la QoS

Je m'appelle Pedro Larroy

```
<piotr%member.fsf.org>
```

. Je décris ici une configuration dans le cas où de nombreux utilisateurs seraient connectés à Internet à travers un routeur Linux qui possède une adresse publique et qui réalise de la traduction d'adresse (NAT). J'utilise cette configuration QoS pour fournir l'accès à 198 utilisateurs dans une cité universitaire dans laquelle je vis et où j'administre le réseau. Les utilisateurs sont de gros consommateurs de programmes "peer to peer" et un contrôle de trafic correct est nécessaire. J'espère que ceci servira d'exemples pratiques à tous les lecteurs intéressés par le lartc.

Dans un premier temps, la configuration sera réalisée pas à pas et, à la fin, j'expliquerai comment rendre ce processus automatique au démarrage. Le réseau utilisé pour cet exemple est un réseau local connecté à Internet à travers un routeur Linux ayant une adresse publique. L'ajout d'un ensemble de règles iptables permettrait facilement l'extension à plusieurs adresses IP publiques. Les éléments suivants sont nécessaires :

Linux 2.4.18 ou une version du noyau supérieure installée

Si vous utilisez le noyau 2.4.18, vous devrez appliquer la mise à jour HTB.

iproute

Soyez également sûr que le binaire "tc" est compatible avec HTB. Un binaire pré compilé est distribué avec HTB.

iptables

15.10.1. Commençons l'optimisation de cette rare bande passante

Tout d'abord, nous allons configurer des gestionnaires de mise en file d'attente dans lesquels nous allons classer le trafic. Nous créons un gestionnaire htb composé de 6 classes avec des priorités croissantes. Nous avons alors des classes qui obtiendront le débit alloué et qui pourront, de plus, utiliser la bande passante dont les autres classes n'ont pas besoin. Rappelons que les classes de plus hautes priorités (correspondant aux nombres prio les plus faibles) obtiendront en premier l'excès de bande passante. Notre liaison ADSL à un flux descendant de 2Mbits/s et un flux montant de 300 kbits/s. J'utilise un débit de seuil (ceil rate) de 240 kbits/s car, au-delà de cette limite, les problèmes de latence commencent à prendre de l'ampleur. Ceci est dû au remplissage d'un tampon placé quelque part entre nous et les hôtes distants.

Ajuster la variable CEIL à 75% de votre bande passante montante maximum et ajuster le nom de l'interface (eth0 dans la suite) à celle qui a l'adresse publique Internet. Exécutez ce qui suit dans un shell root :

```
CEIL=240
tc qdisc add dev eth0 root handle 1: htb default 15
tc class add dev eth0 parent 1: classid 1:1 htb rate ${CEIL}kbit ceil ${CEIL}kbit
tc class add dev eth0 parent 1:1 classid 1:10 htb rate 80kbit ceil 80kbit prio 0
tc class add dev eth0 parent 1:1 classid 1:11 htb rate 80kbit ceil ${CEIL}kbit prio 1
tc class add dev eth0 parent 1:1 classid 1:12 htb rate 20kbit ceil ${CEIL}kbit prio 2
tc class add dev eth0 parent 1:1 classid 1:13 htb rate 20kbit ceil ${CEIL}kbit prio 2
tc class add dev eth0 parent 1:1 classid 1:14 htb rate 10kbit ceil ${CEIL}kbit prio 3
tc class add dev eth0 parent 1:1 classid 1:15 htb rate 30kbit ceil ${CEIL}kbit prio 3
tc qdisc add dev eth0 parent 1:12 handle 120: sfq perturb 10
tc qdisc add dev eth0 parent 1:13 handle 130: sfq perturb 10
tc qdisc add dev eth0 parent 1:14 handle 140: sfq perturb 10
tc qdisc add dev eth0 parent 1:15 handle 150: sfq perturb 10
```

Nous avons juste créé une arborescence htb avec un seul niveau de profondeur. Quelque chose comme ceci :

```
+-----+
| racine 1: |
+-----+
|
+-----+
| classe 1:1 |
+-----+
| | | | |
+---+ +---+ +---+ +---+ +---+ +---+
|1:10| |1:11| |1:12| |1:13| |1:14| |1:15|
+---+ +---+ +---+ +---+ +---+ +---+
```

classid 1:10 htb rate 80kbit ceil 80kbit prio 0

Ceci est la classe de priorité la plus élevée. Les paquets de cette classe auront le plus faible délai et obtiendront en premier l'excès de bande passante. C'est donc une bonne idée de limiter le débit de seuil de cette classe. Nous enverrons dans cette classe les paquets qui ont un avantage à avoir un faible délai, tel que le trafic interactif : *ssh*, *telnet*, *dns*, *quake3*, *irc*, et les paquets avec le bit SYN activé.

classid 1:11 htb rate 80kbit ceil \${CEIL}kbit prio 1

Nous avons ici la première classe dans laquelle nous commençons à mettre du trafic de masse. Dans mon exemple, j'ai le trafic provenant de mon serveur web local et les requêtes pour les pages web : respectivement le port source 80 et le port destination 80.

classid 1:12 htb rate 20kbit ceil \${CEIL}kbit prio 2

Dans cette classe, je mettrai le trafic configuré avec le champ TOS "Débit Maximum" activé, ainsi que le reste du trafic provenant des *processus locaux* de mon routeur vers Internet. Les classes suivantes ne recevront donc que du trafic routé par cette machine.

classid 1:13 htb rate 20kbit ceil \${CEIL}kbit prio 2

Cette classe est pour le trafic des autres machines «NATées» (NdT : bénéficiant du service de traduction d'adresse) qui ont besoin d'une priorité plus grande dans leur trafic de masse.

classid 1:14 htb rate 10kbit ceil \${CEIL}kbit prio 3

Le trafic mail (SMTP, pop3,...) et les paquets configurés avec le champ TOS "Coût Minimum" seront envoyés dans cette classe.

classid 1:15 htb rate 30kbit ceil \${CEIL}kbit prio 3

Finalement, nous avons ici le trafic de masse des machines "NATées" se trouvant derrière le routeur. Les paquets liés à *kazaa*, *edonkey* et autres iront ici pour ne pas interférer avec les autres services.

15.10.2. Classification des paquets

Nous avons configuré le gestionnaire de mise en file d'attente, mais aucune classification de paquets n'a encore été faite. Pour l'instant, tous les paquets sortants passent par la classe 1:15 (car, nous avons utilisé : `tc qdisc add dev eth0 root handle 1: htb default 15`). Nous devons donc maintenant indiquer où doivent aller les paquets. Ceci est la partie la plus importante.

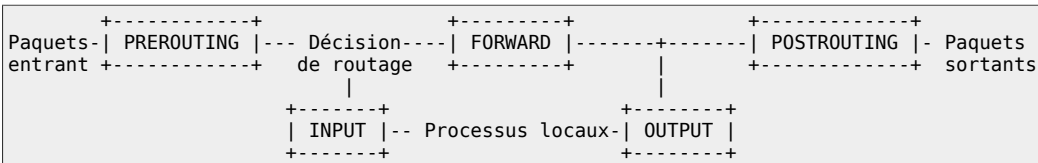
Nous allons maintenant configurer les filtres de tel sorte que nous puissions classer les paquets avec iptables. Je préfère vraiment le faire avec iptables car celui-ci est très souple et que nous avons un compteur

de paquets pour chaque règle. De plus, avec la cible RETURN, les paquets n'ont pas besoin de traverser toutes les règles. Nous exécutons les commandes suivantes :

```
tc filter add dev eth0 parent 1:0 protocol ip prio 1 handle 1 fw classid 1:10
tc filter add dev eth0 parent 1:0 protocol ip prio 2 handle 2 fw classid 1:11
tc filter add dev eth0 parent 1:0 protocol ip prio 3 handle 3 fw classid 1:12
tc filter add dev eth0 parent 1:0 protocol ip prio 4 handle 4 fw classid 1:13
tc filter add dev eth0 parent 1:0 protocol ip prio 5 handle 5 fw classid 1:14
tc filter add dev eth0 parent 1:0 protocol ip prio 6 handle 6 fw classid 1:15
```

Nous indiquons simplement au noyau que les paquets qui ont une valeur FWMARK spécifique (handle x fw) vont dans la classe spécifiée (classid x:x). Voyons maintenant comment marquer les paquets avec iptables.

Tout d'abord, nous devons comprendre comment les paquets traversent les filtres avec iptables :



Je suppose que toutes vos tables ont leur politique par défaut configurée à ACCEPT (-P ACCEPT), ce qui devrait être le cas si vous n'avez pas encore touché à iptables. Notre réseau privé est une classe B avec l'adresse 172.17.0.0/16 et notre adresse publique est 212.170.21.172.

Nous indiquons au noyau de faire de la traduction d'adresse NAT; les clients du réseau privé peuvent alors commencer à dialoguer avec l'extérieur.

```
echo 1 > /proc/sys/net/ipv4/ip_forward
iptables -t nat -A POSTROUTING -s 172.17.0.0/255.255.0.0 -o eth0 -j SNAT --to-source 212.170.21.172
```

Vérifions maintenant que les paquets transitent bien à travers 1:15 :

```
tc -s class show dev eth0
```

Vous pouvez commencer à marquer les paquets en ajoutant les règles dans la chaîne PREROUTING de la table mangle.

```
iptables -t mangle -A PREROUTING -p icmp -j MARK --set-mark 0x1
iptables -t mangle -A PREROUTING -p icmp -j RETURN
```

Vous devriez maintenant être capable de voir l'évolution du compteur de paquets quand vous pinguez des sites sur Internet depuis les machines du réseau privé. Vérifiez que le compteur de paquets augmente dans 1:10 :

```
tc -s class show dev eth0
```

Nous avons mis -j RETURN de manière à ce que les paquets ne traversent pas toutes les règles. Les paquets icmp ne scruteront pas les autres règles définies sous RETURN. Gardez ceci à l'esprit. Nous commençons maintenant à ajouter d'autres règles pour gérer les champs TOS :

```
iptables -t mangle -A PREROUTING -m tos --tos Minimize-Delay -j MARK --set-mark 0x1
iptables -t mangle -A PREROUTING -m tos --tos Minimize-Delay -j RETURN
iptables -t mangle -A PREROUTING -m tos --tos Minimize-Cost -j MARK --set-mark 0x5
iptables -t mangle -A PREROUTING -m tos --tos Minimize-Cost -j RETURN
iptables -t mangle -A PREROUTING -m tos --tos Maximize-Throughput -j MARK --set-mark 0x6
iptables -t mangle -A PREROUTING -m tos --tos Maximize-Throughput -j RETURN
```

Donnons la priorité aux paquets SSH :

```
iptables -t mangle -A PREROUTING -p tcp -m tcp --sport 22 -j MARK --set-mark 0x1
iptables -t mangle -A PREROUTING -p tcp -m tcp --sport 22 -j RETURN
```

Une bonne idée est de donner la priorité aux paquets initiant une connexion tcp, à savoir ceux qui ont le bit SYN activé :

```
iptables -t mangle -I PREROUTING -p tcp -m tcp --tcp-flags SYN,RST,ACK SYN -j MARK --set-mark 0x1
iptables -t mangle -I PREROUTING -p tcp -m tcp --tcp-flags SYN,RST,ACK SYN -j RETURN
```

Et ainsi de suite. Après la mise en place des règles dans la chaîne PREROUTING de la table "mangle", nous terminons par :

```
iptables -t mangle -A PREROUTING -j MARK --set-mark 0x6
```

Ainsi, le trafic non marqué est dirigé vers 1:15. En fait, cette dernière étape n'est pas nécessaire puisque la classe par défaut est 1:15. Un marquage est quand même réalisé de manière à avoir une cohérence pour l'ensemble de la configuration. De plus, il est utile d'avoir une comptabilité pour cette règle.

C'est une bonne idée de faire de même avec la chaîne OUTPUT. Répétez ces commandes avec -A OUTPUT à la place de PREROUTING (s/PREROUTING/OUTPUT/). Le trafic généré localement (sur le routeur Linux) sera alors également classifié. Je termine la chaîne OUTPUT par -j MARK --set-mark 0x3 de tel sorte que le trafic local ait une priorité plus grande.

15.10.3. Améliorer notre configuration

Toute notre configuration est maintenant opérationnelle. Prenez du temps pour regarder les graphes et observer où votre bande passante est la plus utilisée et cela correspond à vos souhaits. J'ai fait ceci pendant de nombreuses heures, ce qui m'a permis d'avoir une connexion Internet fonctionnant très bien. Autrement,

vous serez confronté en permanence à des "timeout" et des allocations de bande passante presque nulles pour les nouvelles connexions tcp.

Si vous repérez des classes qui sont pleines la plupart du temps, ce peut être une bonne idée de leur attacher un autre gestionnaire de mise en file d'attente de manière à ce que le partage de la bande passante soit plus équitable :

```
tc qdisc add dev eth0 parent 1:13 handle 130: sfq perturb 10
tc qdisc add dev eth0 parent 1:14 handle 140: sfq perturb 10
tc qdisc add dev eth0 parent 1:15 handle 150: sfq perturb 10
```

15.10.4. Rendre tout ceci actif au démarrage

Il est certain que ceci peut être fait de différentes façons. Dans mon cas, j'ai un shell script `/etc/init.d/packetfilter` qui accepte les arguments `[start | stop | stop-tables | start-tables | reload-tables]`. Celui-ci configure les gestionnaires de mise en file d'attente et charge les modules du noyau nécessaires et se comporte donc comme un démon. Le même script charge les règles iptables à partir de `/etc/network/iptables-rules`. Je vais l'embellir un peu et le rendrait disponible sur ma page web [ici](#)¹

¹ <http://omega.resa.es/piotr/files/packetfilter.tar.bz2>

Les ponts sont des périphériques qui peuvent être installés dans un réseau sans aucune reconfiguration. Un commutateur réseau est basiquement un pont multi-ports. Un pont est souvent un commutateur avec 2 ports. Cependant, Linux supporte très bien plusieurs interfaces dans un pont, le conduisant à fonctionner comme un vrai commutateur.

Les ponts sont souvent déployés quand on est confronté à un réseau défaillant qui a besoin d'être réparé sans aucune modification. Dans la mesure où un pont est un équipement de niveau 2, la couche sous la couche IP, les routeurs et serveurs ne sont pas conscients de son existence. Ceci signifie que vous pouvez bloquer ou modifier certains paquets de manière transparente ou mettre en forme le trafic.

Un autre élément positif est qu'un pont peut souvent être remplacé par un câble croisé ou un hub quand il tombe en panne.

L'aspect négatif est que la mise en place d'un pont peut engendrer beaucoup de confusion, à moins qu'il ne soit très bien configuré. Le pont n'apparaît pas dans les traceroute, mais pourtant des paquets disparaissent sans raison ou sont changés en allant d'un point A à un point B ('ce réseau est HANTE !). Vous devriez également vous demander si une organisation qui "ne veut rien changer" fait le bon choix.

Le pont Linux 2.4/2.5 est documenté sur [cette page](#)¹.

16.1. Etat des ponts et iptables

Au moment de Linux 2.4.20, le pont et iptables ne se "voient" pas l'un l'autre sans une aide. Si vous "pontez" les paquets de eth0 à eth1, ils ne "passent" pas par iptables. Ceci signifie que vous ne pouvez pas faire de filtrage, de traduction d'adresse (NAT), de désossage ou quoique ce soit d'autres. Ceci a été corrigé dans les versions 2.5.45 et supérieures.

Vous devriez également regarder 'ebtables', qui est encore un autre projet. Il vous permettra de faire des choses vraiment terribles comme MACNAT et 'brouting'. C'est vraiment effroyable.

16.2. Pont et mise en forme

Ca marche comme dans les réclames. Soyez sûr du côté attribué à chaque interface. Autrement, il se peut que vous mettiez en forme le trafic sortant au niveau de votre interface interne, ce qui ne marchera pas. Utilisez tcpdump si nécessaire.

16.3. Pseudo-pont avec du Proxy-ARP

Si vous voulez juste implémenter un pseudo pont, allez jusqu'à la section "Implémentez-le". Cependant, il est sage de lire un peu la façon dont il fonctionne en pratique.

Un pseudo pont travaille de manière un peu différente. Par défaut, un pont transmet les paquets sans les altérer d'une interface à une autre. Il ne regarde que l'adresse matérielle des paquets pour déterminer où ils doivent aller. Ceci signifie que vous pouvez "pontez" un trafic que Linux ne comprend pas, aussi longtemps qu'il y a une adresse matérielle.

Un "pseudo pont" travaille différemment et ressemble plus à un routeur caché qu'à un pont. Mais, comme un pont, il a un impact faible sur l'architecture du réseau.

Le fait qu'il ne soit pas un pont présente l'avantage que les paquets traversent réellement le noyau, et peuvent être filtrés, modifiés, redirigés ou reroutés.

Un pont réel peut également réaliser ces tours de force, mais il a besoin d'un code spécial, comme le Ethernet Frame Diverter ou la mise à jour mentionnée au-dessus.

Un autre avantage d'un pseudo pont est qu'il ne transmet pas les paquets qu'il ne comprend pas, nettoyant ainsi votre réseau de beaucoup de cochonneries. Dans le cas où vous auriez besoin de ces cochonneries (comme les paquets SAP ou Netbeui), utilisez un vrai pont.

16.3.1. ARP & Proxy-ARP

Quand un hôte veut dialoguer avec un autre hôte sur le même segment physique, il envoie un paquet du Protocole de Résolution d'Adresse (ARP) qui, en simplifiant quelque peu, est lu comme ceci : "Qui a 10.0.0.1, le dire à 10.0.0.7". En réponse à ceci, 10.0.0.1 renvoie un petit paquet "ici".

10.0.0.7 envoie alors des paquets à l'adresse matérielle mentionnée dans le paquet "ici". Il met dans un cache cette adresse matérielle pour un temps relativement long et, après l'expiration du cache, repose sa question.

Quand on construit un pseudo pont, on configure le pont pour qu'il réponde à ces paquets ARP, les hôtes du réseau envoyant alors leurs paquets au pont. Le pont traite alors ces paquets et les envoie à l'interface adaptée.

Donc, en résumé, quand un hôte d'un côté du pont demande l'adresse matérielle d'un hôte se situant de l'autre côté, le pont répond avec un paquet qui dit "transmets le moi".

De cette façon, tout le trafic de données est transmis à la bonne place et il traverse toujours le pont.

16.3.2. Implémentez-le

Les versions anciennes du noyau linux permettait de faire du proxy ARP uniquement à une granularité sous réseaux. Ainsi, pour configurer un pseudo pont, il fallait spécifier les bonnes routes vers les deux côtés du

¹ <http://bridge.sourceforge.net/>

pont, et également créer les règles proxy-ARP correspondantes. C'était pénible, déjà par la quantité de texte qu'il fallait taper, puis parce qu'il était facile de se tromper et créer des configurations erronées, où le pont répondait à des requêtes pour des réseaux qu'il ne savait pas router.

Avec Linux 2.4 (et peut-être bien le 2.2), cette possibilité a été retirée et a été remplacée par une option dans le répertoire /proc, appelée "proxy-arp". La procédure pour construire un pseudo pont est maintenant :

1. Assigner une adresse à chaque interface, la "gauche" et la "droite"
2. Créer des routes pour que votre machine connaisse quels hôtes résident à gauche et quels hôtes résident à droite
3. Activer le proxy-ARP sur chaque interface `echo 1 > /proc/sys/net/ipv4/conf/ethL/proxy_arp` `echo 1 > /proc/sys/net/ipv4/conf/ethR/proxy_arp` où L et R désignent les numéros de l'interface du côté gauche (Left) et de celle du côté droit (Right)

N'oubliez pas également d'activer l'option `ip_forwarding` ! Quand on convertit un vrai pont, il se peut que vous trouviez cette option désactivée dans la mesure où il n'y en a pas besoin pour un pont.

Une autre chose que vous devriez considérer lors de la conversion est que vous aurez besoin d'effacer le cache arp des ordinateurs du réseau. Le cache arp peut contenir d'anciennes adresses matérielles du pont qui ne sont plus correctes.

Sur un Cisco, ceci est réalisé en utilisant la commande `'clear arp-cache'` et, sous linux, en utilisant `'arp -d ip.adresse'`. Vous pouvez aussi attendre l'expiration manuelle du cache, ce qui peut être plutôt long.

Il se peut que vous découvriez également que votre réseau était mal configuré si vous avez/aviez l'habitude de spécifier les routes sans les masques de sous-réseau. Dans le passé, certaines versions de **route** pouvaient correctement deviner le masque ou, au contraire, se tromper sans pour autant vous le notifier. Quand vous faites du routage chirurgical comme décrit plus haut, il est **vital** que vous vérifiez vos masques de sous-réseau.

Si votre réseau commence à devenir vraiment gros ou si vous commencez à considérer Internet comme votre propre réseau, vous avez besoin d'outils qui routent dynamiquement vos données. Les sites sont souvent reliés entre eux par de multiples liens, et de nouveaux liens surgissent en permanence.

L'Internet utilise la plupart du temps les standards OSPF (RFC 2328) et BGP4 (RFC 1771). Linux supporte les deux, par le biais de `gated` et `zebra`.

Ce sujet est pour le moment hors du propos de ce document, mais laissez-nous vous diriger vers des travaux de référence :

Vue d'ensemble :

Cisco Systems [Cisco Systems Designing large-scale IP Internetworks](#)¹

Pour OSPF :

Moy, John T. "OSPF. The anatomy of an Internet routing protocol" Addison Wesley. Reading, MA. 1998.

Halabi a aussi écrit un très bon guide sur la conception du routage OSPF, mais il semble avoir été effacé du site Web de Cisco.

Pour BGP :

Halabi, Bassam "Internet routing architectures" Cisco Press (New Riders Publishing). Indianapolis, IN. 1997.

Il existe aussi

Cisco Systems

[Using the Border Gateway Protocol for Interdomain Routing](#)²

Bien que les exemples soient spécifiques à Cisco, ils sont remarquablement semblables au langage de configuration de Zebra :-)

17.1. Configurer OSPF avec Zebra

Pedro Larroy Tovar <piotr%member.fsf.org>

Contactez-moi³ si les informations qui suivent ne sont pas exactes ou si vous avez des suggestions. **Zebra**⁴ est un formidable logiciel de routage dynamique écrit par Kunihiro Ishiguro, Toshiaki Takada et Yasuhiro Ohara. Configurer OSPF avec zebra est simple et rapide mais, en pratique, il y a de nombreux paramètres dans le cas où vous auriez des besoins spécifiques. OSPF est l'abréviation de Open Shortest Path First et quelques une de ses fonctionnalités sont :

hiérarchique

Les réseaux sont regroupés par *zones (areas)*, qui sont interconnectées par une *zone épine dorsale* qui sera appelée *zone 0*. Tout le trafic passe par la zone 0 et tous les routeurs de cette zone ont les informations de routage de toutes les autres zones.

convergence rapide

Les routes sont propagées très rapidement, comparé à RIP par exemple.



Note de publication

Avec le protocole OSPF, ce sont les états de liens ou *Link States* qui sont propagés, et non pas les routes. La quantité de données échangée entre routeurs OSPF est très inférieure à celle échangée entre routeurs RIP pour lesquels la totalité de la table de routage est transmise périodiquement.

Tout changement d'un ou plusieurs états de liens provoque un recalcul de la topologie complète de l'aire OSPF. La rapidité de convergence d'OSPF est alors liée à l'algorithme de Dijkstra pour le calcul des routes. Si n est le nombre de sommets du graphe associé à la base de données topologique du domaine OSPF, la rapidité de convergence est d'un ordre qui varie de $O(n)$ à $O(n^2)$ en fonction du degré de maillage du réseau.

économie de bande passante

Utilise la multi-distribution à la place de la diffusion, ce qui évite de submerger les autres hôtes avec des informations de routage sans intérêt pour eux. La multi-distribution réduit ainsi le débit sur le réseau. De même, *les routeurs internes* (ceux dont toutes les interfaces sont situées dans la même zone) n'obtiennent pas d'informations sur les autres zones.



Note de publication

Les routeurs internes à une zone possèdent une entrée spécifique correspondant à chacun des réseaux d'une autre zone si l'on n'a pas synthétisé de super-réseau (fonction `summary`) au niveau des *Area Border routers*.

Les routeurs avec des interfaces dans plus d'une zone sont appelés *Area Border Routers*. Ils possèdent les informations de topologie sur les zones auxquelles ils sont connectés.

¹ <http://www.cisco.com/univercd/cc/td/doc/cisintwk/idg4/nd2003.htm>

² <http://www.cisco.com/univercd/cc/td/doc/cisintwk/ics/icsbgp4.htm>

³ <mailto:piotr%member.fsf.org>

⁴ <http://www.quagga.net>

Utilisation intensive de CPU

OSPF est basé sur l'algorithme de Dijkstra **Shortest Path First**⁵, qui est coûteux en temps de calcul comparé aux autres algorithmes de routage.



Note de publication

Les charges CPU et mémoire ne sont pas liées exclusivement à l'exécution de l'algorithme Dijkstra. Elles dépendent beaucoup de la gestion de la base de données des états de liens (*Link States*).

Ce n'est pas forcément mauvais, dans la mesure où le plus court chemin est calculé uniquement pour chaque zone. Donc, pour les réseaux de petite à moyenne taille, ce ne sera pas un problème ; vous ne vous en rendez pas compte.

Information d'état de lien

OSPF prend en compte les caractéristiques spécifiques des réseaux et interfaces, telles que la bande passante, les défauts de liens et le coût monétaire.

Protocole ouvert et logiciel sous license GPL

OSPF est un protocole ouvert et Zebra est un logiciel sous license GPL, ce qui représente un avantage évident par rapport aux protocoles et logiciels propriétaires.

17.1.1. Prérequis

Noyau Linux :

Compilé avec CONFIG_NETLINK_DEV and CONFIG_IP_MULTICAST (Je ne sais pas si d'autres éléments sont également nécessaires).

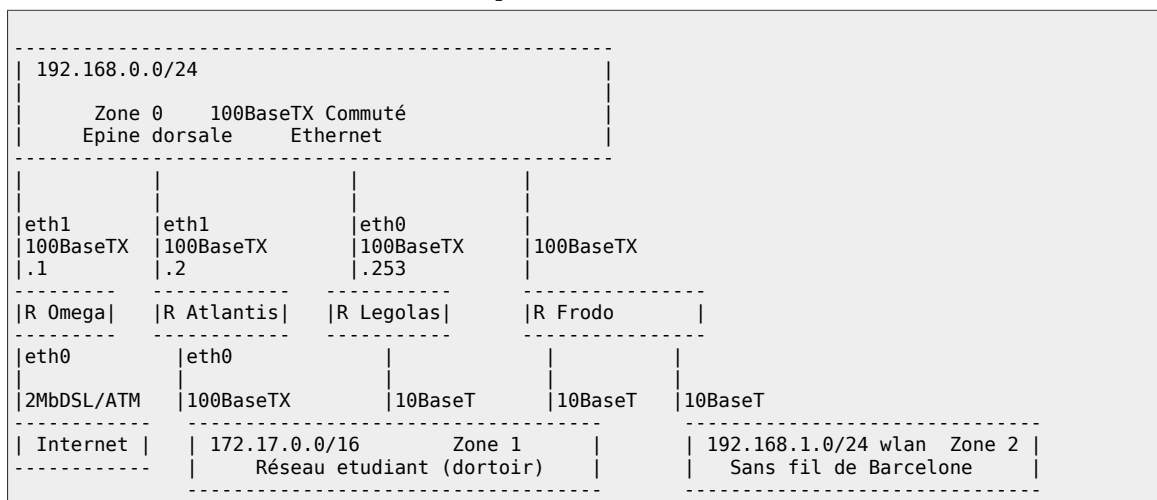
Iproute

Zebra

Récupérez-le avec votre gestionnaire de paquet favori ou à partir de <http://www.quagga.net>.

17.1.2. Configurer Zebra

Prenons le réseau suivant comme exemple :



Ne soyez pas effrayé par ce diagramme, Zebra réalise la plus grande partie du travail automatiquement ; ce qui ne demandera aucun travail de saisie des routes avec Zebra. Il serait pénible de maintenir toutes ces routes à la main au quotidien. La chose la plus importante à maîtriser clairement, c'est la topologie du réseau. Faites particulièrement attention à la zone 0, puisque c'est la plus importante. Dans un premier temps, configurez Zebra en éditant `zebra.conf` et en l'adaptant à vos besoins :

```

hostname omega
password xxx
enable password xxx
!
! Interface's description.
!
!interface lo
! description test of desc.
!
interface eth1
multicast
!
! Static default route
!
ip route 0.0.0.0/0 212.170.21.129
!
  
```

⁵ <http://www soi wide ad jp/class/99007/slides/13/07.html>

```
log file /var/log/zebra/zebra.log
```

Debian nécessite également l'édition de `/etc/zebra/daemons` pour qu'ils soient lancés au démarrage :

```
zebra=yes
ospfd=yes
```

Nous devons maintenant éditer `ospfd.conf` si vous utilisez encore IPv4 ou `ospfd6.conf` si vous travaillez avec IPv6. Mon fichier `ospfd.conf` ressemble à ceci :

```
hostname omega
password xxx
enable password xxx
!
router ospf
 network 192.168.0.0/24 area 0
 network 172.17.0.0/16 area 1
!
! log stdout
log file /var/log/zebra/ospfd.log
```

Ceci indique à ospf la topologie de notre réseau.

17.1.3. Exécuter Zebra

Nous devons maintenant démarrer Zebra soit à la main en tapant "zebra -d", soit avec un script comme "`/etc/init.d/zebra start`". En regardant attentivement les logs de `ospfd`, on peut voir les éléments suivants :

```
2002/12/13 22:46:24 OSPF: interface 192.168.0.1 join AllSPFRouters Multicast group.
2002/12/13 22:46:34 OSPF: SMUX_CLOSE with reason: 5
2002/12/13 22:46:44 OSPF: SMUX_CLOSE with reason: 5
2002/12/13 22:46:54 OSPF: SMUX_CLOSE with reason: 5
2002/12/13 22:47:04 OSPF: SMUX_CLOSE with reason: 5
2002/12/13 22:47:04 OSPF: DR-Election[1st]: Backup 192.168.0.1
2002/12/13 22:47:04 OSPF: DR-Election[1st]: DR 192.168.0.1
2002/12/13 22:47:04 OSPF: DR-Election[2nd]: Backup 0.0.0.0
2002/12/13 22:47:04 OSPF: DR-Election[2nd]: DR 192.168.0.1
2002/12/13 22:47:04 OSPF: interface 192.168.0.1 join AllDRouters Multicast group.
2002/12/13 22:47:06 OSPF: DR-Election[1st]: Backup 192.168.0.2
2002/12/13 22:47:06 OSPF: DR-Election[1st]: DR 192.168.0.1
2002/12/13 22:47:06 OSPF: Packet[DD]: Negotiation done (Slave).
2002/12/13 22:47:06 OSPF: nsm_change_status(): scheduling new router-LSA origination
2002/12/13 22:47:11 OSPF: ospf_intra_add_router: Start
```

Ignorez le message `SMUX_CLOSE` pour l'instant dans la mesure où il concerne SNMP. Nous pouvons voir que 192.168.0.1 est *routeur désigné (Designated Router)* et que 192.168.0.2 est le *routeur désigné de sauvegarde (Backup Designated Router)*.

Nous pouvons également interagir avec zebra et ospfd en exécutant :

```
$ telnet localhost zebra
$ telnet localhost ospfd
```

Voyons comment les routes se sont propagées en se connectant à zebra :

```
root@atlantis:~# telnet localhost zebra
Trying 127.0.0.1...
Connected to atlantis.
Escape character is '^]'.

Hello, this is zebra (version 0.92a).
Copyright 1996-2001 Kunihiro Ishiguro.

User Access Verification

Password:
atlantis> show ip route
Codes: K - kernel route, C - connected, S - static, R - RIP, O - OSPF,
B - BGP, > - selected route, * - FIB route

K>* 0.0.0.0/0 via 192.168.0.1, eth1
C>* 127.0.0.0/8 is directly connected, lo
O 172.17.0.0/16 [110/10] is directly connected, eth0, 06:21:53
C>* 172.17.0.0/16 is directly connected, eth0
O 192.168.0.0/24 [110/10] is directly connected, eth1, 06:21:53
C>* 192.168.0.0/24 is directly connected, eth1
atlantis> show ip ospf border-routers
===== OSPF router routing table =====
R 192.168.0.253 [10] area: (0.0.0.0), ABR
via 192.168.0.253, eth1
[10] area: (0.0.0.1), ABR
via 172.17.0.2, eth0
```

ou directement avec `iproute` :

```
root@omega:~# ip route
212.170.21.128/26 dev eth0 proto kernel scope link src 212.170.21.172
192.168.0.0/24 dev eth1 proto kernel scope link src 192.168.0.1
172.17.0.0/16 via 192.168.0.2 dev eth1 proto zebra metric 20
```

```
default via 212.170.21.129 dev eth0 proto zebra
root@omega:~#
```

Nous pouvons voir les routes Zebra, qui n'étaient pas présentes auparavant. Il est vraiment agréable de voir apparaître les routes quelques secondes après le lancement de zebra et ospfd. Vous pouvez vérifier la connectivité avec les autres hôtes en utilisant ping. Les routes zebra sont automatiques. Vous pouvez ajouter un autre routeur au réseau, configurez Zebra et voilà !

Astuce ; vous pouvez utiliser :

```
tcpdump -i eth1 ip[9] == 89
```

pour analyser les paquets OSPF. Le numéro du protocole OSPF est 89 et le champ du protocole est le 9ième octet de l'en-tête ip.

OSPF possède de nombreux paramètres, spécialement pour les grands réseaux. Dans de prochains développements du HOWTO, nous montrerons des méthodes de réglages fins d'OSPF.

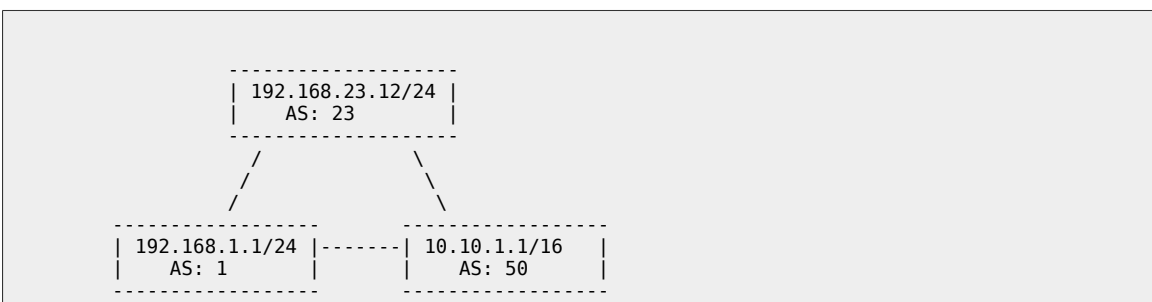
17.2. Configurer BGP4 avec Zebra

Le Border Gateway Protocol Version 4 (BGP4) est un protocole de routage dynamique décrit dans la RFC 1771. Il permet la distribution des informations de connectivité, c'est à dire les tables de routage, vers d'autres nœuds BGP4 actifs. Il peut être utilisé comme un EGP ou un IGP. Dans le mode EGP, chaque nœud doit avoir son propre numéro de système autonome (*Autonomous System (AS)*). BGP4 supporte ?????????? et l'aggrégation de routes (réunir plusieurs routes en une seule). > The Border Gateway Protocol Version 4 (BGP4) is a dynamic routing > protocol described in RFC 1771. It allows the distribution of > reachability information, i.e. routing tables, to other BGP4 > enabled nodes. It can either be used as EGP or IGP, in EGP mode > each node must have its own Autonomous System (AS) number. > BGP4 supports Classless Inter Domain Routing (CIDR) and route > aggregation (merge multiple routes into one).

17.2.1. schéma réseau (Exemple)

Le schéma réseau suivant est utilisé pour les exemples à suivre. AS 1 et 50 ont plusieurs voisins mais nous avons seulement besoin de configurer 1 et 50 comme nos voisins. Les nœuds communiquent entre eux par des tunnels dans cet exemple, mais ce n'est pas une obligation.

Note : les numéros AS utilisés dans cet exemple sont réservés. Veuillez obtenir vos propres numéros AS du RIPE si vous installez des liens officiels.



17.2.2. Configuration (Exemple)

La configuration suivante est écrite pour le nœud 192.168.23.12/24 et elle sera facile à adapter pour les autres nœuds.

Elle commence par des éléments généraux comme le nom de l'hôte, les mots de passe et les options de debug :

```
! hostname
hostname anakin

! login password
password xxx

! enable password (super user mode)
enable password xxx

! path to logfile
log file /var/log/zebra/bgpd.log

! debugging: be verbose (can be removed afterwards)
debug bgp events
debug bgp filters
debug bgp fsm
debug bgp keepalives
debug bgp updates
```

La liste de contrôle d'accès (*Access list*) est utilisée pour limiter la redistribution aux réseaux privés (RFC 1918).

```
! RFC 1918 networks
access-list local_nets permit 192.168.0.0/16
```

```
access-list local_nets permit 172.16.0.0/12
access-list local_nets permit 10.0.0.0/8
access-list local_nets deny any
```

L'etape suivante consiste à configurer chaque AS :

```
! Own AS number
router bgp 23

! IP address of the router
bgp router-id 192.168.23.12

! announce our own network to other neighbors
network 192.168.23.0/24

! advertise all connected routes (= directly attached interfaces)
redistribute connected

! advertise kernel routes (= manually inserted routes)
redistribute kernel
```

Chaque section 'router bgp' contient une liste de voisins auxquels le routeur est connecté :

```
neighbor 192.168.1.1 remote-as 1
neighbor 192.168.1.1 distribute-list local_nets in
neighbor 10.10.1.1 remote-as 50
neighbor 10.10.1.1 distribute-list local_nets in
```

17.2.3. Vérification de la configuration

Note : vtysh est un multiplexeur qui connecte toutes les interfaces utilisateur de zebra ensemble.

```
anakin# sh ip bgp summary
BGP router identifier 192.168.23.12, local AS number 23
2 BGP AS-PATH entries
0 BGP community entries

Neighbor      V    AS MsgRcvd MsgSent   TblVer  InQ  OutQ Up/Down   State/PfxRcd
10.10.0.1     4    50     35     40         0    0   0 00:28:40      1
192.168.1.1   4     1  27574  27644         0    0   0 03:26:04     14

Total number of neighbors 2
anakin#
anakin# sh ip bgp neighbors 10.10.0.1
BGP neighbor is 10.10.0.1, remote AS 50, local AS 23, external link
  BGP version 4, remote router ID 10.10.0.1
  BGP state = Established, up for 00:29:01
  ....
anakin#
```

Voyons quelles routes nous avons obtenues de nos voisins :

```
anakin# sh ip ro bgp
Codes: K - kernel route, C - connected, S - static, R - RIP, O - OSPF,
       B - BGP, > - selected route, * - FIB route

B>* 172.16.0.0/14 [20/0] via 192.168.1.1, tun0, 2d10h19m
B>* 172.30.0.0/16 [20/0] via 192.168.1.1, tun0, 10:09:24
B>* 192.168.5.10/32 [20/0] via 192.168.1.1, tun0, 2d10h27m
B>* 192.168.5.26/32 [20/0] via 192.168.1.1, tun0, 10:09:24
B>* 192.168.5.36/32 [20/0] via 192.168.1.1, tun0, 2d10h19m
B>* 192.168.17.0/24 [20/0] via 192.168.1.1, tun0, 3d05h07m
B>* 192.168.17.1/32 [20/0] via 192.168.1.1, tun0, 3d05h07m
B>* 192.168.32.0/24 [20/0] via 192.168.1.1, tun0, 2d10h27m
anakin#
```

Ce chapitre est une liste des projets ayant une relation avec le routage avancé et la mise en forme du trafic sous Linux. Certains de ces liens mériteraient des chapitres spécifiques, d'autres sont très bien documentés, et n'ont pas besoin de HOWTO en plus.

Implémentation VLAN 802.1Q pour Linux [\(site\)](#)¹

VLAN est une façon très sympa de diviser vos réseaux d'une manière plus virtuelle que physique. De bonnes informations sur les VLAN pourront être trouvées [ici](#)². Avec cette implémentation, votre boîte Linux pourra dialoguer VLAN avec des machines comme les Cisco Catalyst, 3Com: {Corebuilder, Netbuilder II, SuperStack II switch 630}, Extreme Ntwks Summit 48, Foundry: {ServerIronXL, FastIron}.

Implémentation alternative VLAN 802.1Q pour Linux [\(site\)](#)³

Une implémentation alternative de VLAN pour Linux. Ce projet a démarré suite au désaccord avec l'architecture et le style de codage du projet VLAN 'établi', avec comme résultat une structure de l'ensemble plus clair. Mise à jour : a été inclus dans le noyau 2.4.14 (peut-être dans le 2.4.13).

Un bon HOWTO à propos des VLAN peut être trouvé [ici](#)⁴.

Mise à jour : a été incluse dans le noyau à partir de la version 2.4.14 (peut-être 13).

Serveur Linux Virtuel (Linux Virtual Server) [\(site\)](#)⁵

Ces personnes sont très talentueuses. Le Serveur Virtuel Linux est un serveur à haute disponibilité, hautement évolutif, construit autour d'une grappe (cluster) de serveurs, avec un équilibreur de charge tournant sur le système d'exploitation Linux. L'architecture du cluster est transparente pour les utilisateurs finaux, qui ne voient qu'un simple serveur virtuel.

En résumé, que vous ayez besoin d'équilibrer votre charge ou de contrôler votre trafic, LVS aura une manière de le faire. Certaines de leurs techniques sont positivement diaboliques !. Par exemple, ils permettent à plusieurs machines d'avoir une même adresse IP, mais en désactivant l'ARP dessus. Seule la machine LVS qui a, elle, l'ARP actif, décide de l'hôte qui manipulera le paquet entrant. Celui-ci est envoyé avec la bonne adresse MAC au serveur choisi. Le trafic sortant passe directement par le routeur, et non par la machine LVS, qui, par conséquent n'a pas besoin de voir vos 5Gbit/s de données allant sur Internet. Cette machine LVS ne peut alors pas être un goulot d'étranglement.

L'implémentation de LVS nécessite une mise à jour pour les noyaux 2.0 et 2.2, alors qu'un module Netfilter est disponible dans le 2.4. Il n'y a donc pas besoin de mise à jour pour cette version du noyau. Le support 2.4 est encore en développement. Battez-vous donc avec et envoyez vos commentaires ou vos mises à jour.

CBQ.init [\(site\)](#)⁶

Configurer CBQ peut être un peu intimidant, spécialement si votre seul souhait est de mettre en forme le trafic d'ordinateurs placés derrière un routeur. CBQ.init peut vous aider à configurer Linux à l'aide d'une syntaxe simplifiée.

Par exemple, si vous voulez que tous les ordinateurs de votre réseau 192.168.1.0/24 (sur eth1 10 Mbits) aient leur vitesse de téléchargement limitée à 28 Kbits, remplissez le fichier de configuration de CBQ.init avec ce qui suit :

```
DEVICE=eth1,10Mbit,1Mbit
RATE=28Kbit
WEIGHT=2Kbit
PRIO=5
RULE=192.168.1.0/24
```

Utiliser simplement ce programme si le 'comment et pourquoi' ne vous intéresse pas. Nous utilisons CBQ.init en production et il marche très bien. On peut même faire des choses plus avancées, comme la mise en forme dépendant du temps. La documentation est directement intégrée dans le script, ce qui explique l'absence d'un fichier README.

Scripts faciles de mise en forme Chronox [\(site\)](#)⁷

Stephan Mueller (smueller@chronox.de) a écrit deux scripts utiles, "limit.conn" et "shaper". Le premier vous permet de maîtriser une session de téléchargement, comme ceci :

```
# limit.conn -s SERVERIP -p SERVERPORT -l LIMIT
```

Il fonctionne avec Linux 2.2 et 2.4.

Le second script est plus compliqué et peut être utilisé pour mettre en place des files d'attente différentes basées sur les règles iptables. Celles-ci sont utilisées pour marquer les paquets qui sont alors mis en forme.

¹ <http://scry.wanfear.com/~greear/vlan.html>

² ftp://ftp.netlab.ohio-state.edu/pub/jain/courses/cis788-97/virtual_lans/index.htm

³ <http://vlan.sourceforge.net>

⁴ http://scry.wanfear.com/~greear/vlan/cisco_howto.html

⁵ <http://www.LinuxVirtualServer.org/>

⁶ <ftp://ftp.equinox.gu.net/pub/linux/cbq/>

⁷ <http://www.chronox.de>

Implémentation du Protocole Redondant Routeur Virtuel (site)⁸

Ceci est purement pour la redondance. Deux machines avec leurs propres adresses IP et MAC créent une troisième adresse IP et MAC virtuelle. Bien que destiné à l'origine uniquement aux routeurs, qui ont besoin d'adresses MAC constantes, cela marche également pour les autres serveurs.

La beauté de cette approche est l'incroyable facilité de la configuration. Pas de compilation de noyau ou de nécessité de mise à jour, tout se passe dans l'espace utilisateur.

Lancer simplement ceci sur toutes les machines participant au service :

```
# vrrpd -i eth0 -v 50 10.0.0.22
```

Et vous voilà opérationnel ! 10.0.0.22 est maintenant géré par l'un de vos serveurs, probablement le premier à avoir lancé le démon vrrp. Déconnectez maintenant cet ordinateur du réseau et très rapidement, l'adresse 10.0.0.22 et l'adresse MAC seront gérées par l'un des autres ordinateurs.

J'ai essayé ceci et il a été actif et opérationnel en 1 minute. Pour une raison étrange, ma passerelle par défaut a été supprimée. Cependant, l'option -n permet de prévenir cela.

Voici une défaillance en "direct" :

```
64 bytes from 10.0.0.22: icmp_seq=3 ttl=255 time=0.2 ms
64 bytes from 10.0.0.22: icmp_seq=4 ttl=255 time=0.2 ms
64 bytes from 10.0.0.22: icmp_seq=5 ttl=255 time=16.8 ms
64 bytes from 10.0.0.22: icmp_seq=6 ttl=255 time=1.8 ms
64 bytes from 10.0.0.22: icmp_seq=7 ttl=255 time=1.7 ms
```

Pas *un* paquet ping n'a été perdu ! Après 4 paquets, j'ai déconnecté mon P200 du réseau, et mon 486 a pris le relais, ce qui est visible par l'augmentation du temps de latence.

⁸ <http://w3.arobas.net/~jetienne/vrrpd/index.html>

<http://snafu.freedom.org/linux2.2/iproute-notes.html>

Contient beaucoup d'informations techniques, et de commentaires sur le noyau.

<http://www.davin.ottawa.on.ca/ols/>

Transparents de Jamal Hadi Salim, un des auteurs du contrôleur de trafic de Linux.

<http://defiant.coinet.com/iproute2/ip-cref/>

Version HTML de la documentation LaTeX d'Alexeys ; explique une partie d'iproute2 en détails.

<http://www.aciri.org/floyd/cbq.html>

Sally Floyd a une bonne page sur CBQ, incluant ses publications originales. Aucune n'est spécifique à Linux, mais il y a un travail de discussion sur la théorie et l'utilisation de CBQ. Contenu très technique, mais une bonne lecture pour ceux qui sont intéressés.

Differentiated Services on Linux

This [document](#)¹ par Werner Almesberger, Jamal Hadi Salim et Alexey Kuznetsov. Décrit les fonctions DiffServ du noyau Linux, entre autres les gestionnaires de mise en file d'attente TBF, GRED, DSMARK et le classificateur tcindex.

http://ceti.pl/~ekravietz/cbq/NET4_tc.html

Un autre HOWTO, en polonais ! Vous pouvez cependant copier/coller les lignes de commandes, elles fonctionnent de la même façon dans toutes les langues. L'auteur travaille en collaboration avec nous et sera peut être bientôt un auteur de sections de cet HOWTO.

IOS Committed Access Rate²

Des gens de Cisco qui ont pris la louable habitude de mettre leur documentation en ligne. La syntaxe de Cisco est différente mais les concepts sont identiques, sauf qu'on fait mieux, et sans matériel coutant le prix d'une voiture :-)

TCP/IP Illustrated, volume 1, W. Richard Stevens, ISBN 0-201-63346-9

Sa lecture est indispensable si vous voulez réellement comprendre TCP/IP, et de plus elle est divertissante.

¹ <ftp://icaftp.epfl.ch/pub/linux/diffserv/misc/dsid-01.txt.gz>

² <http://www.cisco.com/univercd/cc/td/doc/product/software/ios111/cc111/car.htm>

Notre but est de faire la liste de toutes les personnes qui ont contribué à ce HOWTO, ou qui nous ont aidés à expliquer le fonctionnement des choses. Alors qu'il n'existe pas actuellement de tableau d'honneur Netfilter, nous souhaitons saluer les personnes qui apportent leur aide.

- Junk Alins
- Joe Van Andel
- Michael T. Babcock
- Christopher Barton
- Peter Bieringer
- Ard van Breemen
- Ron Brinker
- ?ukasz Bromirski
- Lennert Buytenhek
- Esteve Camps
- Ricardo Javier Cardenes
- Stef Coene
- Don Cohen
- Jonathan Corbet
- Gerry N5JXS Creager
- Marco Davids
- Jonathan Day
- Martin aka devik Devera
- Hannes Ebner
- Derek Fawcus
- David Fries
- Stephan "Kobold" Gehring
- Jacek Glinkowski
- Andrea Glorioso
- Thomas Graaf
- Sandy Harris
- Nadeem Hasan
- Erik Hensema
- Vik Heyndrickx
- Spauldo Da Hippie
- Koos van den Hout
- Stefan Huelbrock <shuelbrock%datasystems.de>
- Ayotunde Itayemi
- Alexander W. Janssen <yalla%ynfonatic.de>
- Andreas Jellinghaus <aj%dungeon.inka.de>
- Gareth John <gdjohn%zepler.org>
- Dave Johnson

- Martin Josefsson <gandalf%wlug.westbo.se>
- Andi Kleen <ak%suse.de>
- Andreas J. Koenig <andreas.koenig%anima.de>
- Pawel Krawczyk <kravietz%alfa.ceti.pl>
- Amit Kucheria <amitk@ittc.ku.edu>
- Edmund Lau <edlau%ucf.ics.uci.edu>
- Philippe Latu <philippe.latu%inetdoc.net>
- Arthur van Leeuwen <arthurv1%sci.kun.nl>
- Jose Luis Domingo Lopez
- Robert Lowe
- Jason Lunz <j@cc.gatech.edu>
- Stuart Lynne <sl@fireplug.net>
- Alexey Mahotkin <alexm@formulabez.ru>
- Predrag Malicevic <pmalic@ieee.org>
- Patrick McHardy <kaber@trash.net>
- Andreas Mohr <andi%lisas.de>
- James Morris <jmorris@intercode.com.au>
- Andrew Morton <akpm%zip.com.au>
- Wim van der Most
- Stephan Mueller <smueller@chronox.de>
- Togan Muftuoglu <toganm@yahoo.com>
- Chris Murray <cmurray@stargate.ca>
- Patrick Nagelschmidt <dto@gmx.net>
- Ram Narula <ram@princess1.net>
- Jorge Novo <jnovo@educanet.net>
- Patrik <ph@kurd.nu>
- P?l Osgy?ny <oplab%westel900.net>
- Lutz Preßler <Lutz.Pressler%SerNet.DE>
- Jason Pyeron <jason%pyeron.com>
- Rod Roark <rod%sunsetsystems.com>
- Pavel Roskin <proski@gnu.org>
- Rusty Russell <rusty%rustcorp.com.au>
- Mihai RUSU <dizzy%roedu.net>
- Rob Pitman <rob%pitman.co.za>
- Jamal Hadi Salim <hadi%cyberus.ca>
- Ren? Serral <rserral%ac.upc.es>
- David Sauer <davids%penguin.cz>
- Sheharyar Suleman Shaikh <sss23@drexel.edu>
- Stewart Shields <MourningBlade%bigfoot.com>

- Nick Silberstein <nhsilber@yahoo.com>
- Konrads Smelkov <konrads@interbaltika.com>
- William Stearns
- Andreas Steinmetz <ast%domdv.de>
- Matthew Strait <straitm%mathcs.carleton.edu>
- Jason Tackaberry <tack@linux.com>
- Charles Tassell <ctassell%isn.net>
- Glen Turner <glen.turner%aarnet.edu.au>
- Tea Sponsor: Eric Veldhuyzen <eric%terra.nu>
- Thomas Walpuski <thomas%bender.thinknerd.de>
- Song Wang <wsong@ece.uci.edu>
- Chris Wilson
- Lazar Yanackiev
- Pedro Larroy
 - Chapitre 15, section 10: Exemple d'une solution de traduction d'adresse avec de la QoS
 - Chapitre 17, section 1: Configurer OSPF avec Zebra