

Network Working Group
Request for Comments: 3549
Category: Informational

J. Salim
Znyx Networks
H. Khosravi
Intel
A. Kleen
Suse
A. Kuznetsov
INR/Swsoft
July 2003

Netlink как протокол для служб IP

Linux Netlink as an IP Services Protocol

Статус документа

Этот документ содержит информацию, предназначенную для сообщества Internet, и не задает каких-либо стандартов Internet. Документ может распространяться без ограничений.

Авторские права

Copyright (C) The Internet Society (2003). All Rights Reserved.

Тезисы

Данный документ описывает интерфейс Netlink ОС Linux, который используется операционной системой для обмена сообщениями как между процессами ядра, так и между ядром и пользовательскими процессами. Основное внимание в документе уделяется описанию функциональности Netlink как протокола, связывающего компоненты FEC¹ и CPC², которые определяют работу сервиса IP. Прочие варианты использования Netlink, включая обмен сообщениями внутри ядра и между процессами IPC³, а также настройка конфигурации служб, не относящихся к IP (несетевые службы или сетевые службы других протоколов), в данном документе не рассматриваются.

Документ предназначен для создания информационного контекста на начальном этапе работы группы ForCES⁴ IETF.

Оглавление

| | |
|--|----|
| 1. Введение..... | 2 |
| 1.1. Определения..... | 2 |
| 1.1.1. Компоненты CPC..... | 2 |
| 1.1.2. Компоненты FEC..... | 2 |
| 1.1.2.1. Модель машины пересылки IP в Linux..... | 2 |
| 1.1.3. Службы IP..... | 3 |
| 2. Архитектура Netlink..... | 4 |
| 2.1. Логическая модель Netlink..... | 4 |
| 2.2. Формат сообщений..... | 5 |
| 2.3. Модель протокола..... | 5 |
| 2.3.1. Адресация служб..... | 6 |
| 2.3.2. Заголовок сообщений Netlink..... | 6 |
| 2.3.2.1. Механизмы создания протоколов..... | 7 |
| 2.3.2.2. Сообщение ACK в Netlink..... | 7 |
| 2.3.3. Шаблоны FE системных служб..... | 7 |
| 2.3.3.1. Сервисный модуль сетевого интерфейса..... | 7 |
| 2.3.3.2. Модуль службы адресов IP..... | 8 |
| 3. Определенные в данный момент IP-службы Netlink..... | 9 |
| 3.1. Служба NETLINK_ROUTE..... | 9 |
| 3.1.1. Модуль службы маршрутизации..... | 9 |
| 3.1.2. Модуль учета соседей..... | 11 |
| 3.1.3. Служба контроля трафика..... | 12 |

¹Forwarding Engine Component - компонента машины пересылки.

²Control Plane Component - компонента уровня управления.

³Inter-process communication - обмен информацией между процессами.

⁴Forwarding & Control Element Separation - разделение элементов пересылки и управления. Страница рабочей группы доступна по адресу <http://www.ietf.org/html.charters/forces-charter.html>. Работа группы завершена в марте 2015 г. Прим. перев.

| | |
|--|----|
| 3.2. Служба NETLINK_FIREWALL..... | 13 |
| 3.3. Служба NETLINK_ARPD..... | 15 |
| 4. Литература..... | 15 |
| 4.1. Нормативные документы..... | 15 |
| 4.2. Дополнительная литература..... | 15 |
| 5. Вопросы безопасности..... | 15 |
| 6. Благодарности..... | 16 |
| Приложение 1: Пример иерархии служб..... | 16 |
| Приложение 2: Пример протокола для IP-службы Foo..... | 16 |
| Приложение 2a: Взаимодействие с другими службами IP..... | 16 |
| Приложение 3: Примеры..... | 17 |

1. Введение

Концепция разделения служб IP на управление и пересылку впервые была реализована в начале 1990-х годов в сокетях маршрутизации BSD 4.4 [9]. В то время наибольшую важность представляло простое решение вопроса пересылки пакетов IP (v4) и управление таблицами пересылки IPv4 в CPC (с помощью консольного интерфейса или демона динамической маршрутизации).

Мир IP-сетей с тех давних пор существенно изменился. Linux Netlink, с точки зрения обеспечения сервиса и управления, кроме поддержки сокетов маршрутизации обеспечивает ряд дополнительных функций. Начиная с ядра Linux 2.1, сокет Netlink обеспечивает абстракцию служб IP для нескольких типов сервиса кроме классической пересылки IPv4 в соответствии с RFC 1812.

Мотивом для создания этого документа послужило отнюдь не желание описать весь набор служб, для которых можно использовать Netlink. Фактически многие типы сервиса (групповая маршрутизация, туннелирование, маршрутизация на основе правил и т. д.) просто не рассматриваются в данном документе. Не предназначен документ и для использования в качестве учебника по Netlink. Идея документа заключается в общем описании Netlink и более подробном рассмотрении обязательных компонент в контексте работы группы ForCES - IPv4 и QoS. Документ также служит предварительным описанием множества механизмов, изучение которых представляет интерес в рамках ForCES. Рассматривается подмножество функций, доступных в ядре версии 2.4.6, которая была последней во время подготовки данного документа. Кроме того, документ рассматривает лишь функции, связанные с IPv4.

Документ начинается с концептуальных определений, после чего приводится рассмотрение Netlink в свете этих определений.

1.1. Определения

CP¹ представляет собой среду исполнения, которая может иметь несколько субкомпонент, обозначаемых как CPC². Все CPC, обеспечивающие контроль для разных служб IP, будут выполняться посредством машины пересылки FE³. Такие отношения между компонентами означают возможность наличия нескольких CPC для одной физической CP, если они контролируют несколько служб IP. По сути, связь между CP и FE является абстракцией сервиса.

1.1.1. Компоненты CPC

Компоненты управляющего плана CPC включают сигнальные протоколы от динамических протоколов маршрутизации (например, OSPF [5]) до протоколов распространения тегов (например, CR-LDP [7]). Классические протоколы и операции управления также входят в эту категорию. Среди них такие механизмы, как SNMP [6], COPS [4] и фирменные средства настройки конфигурации CLI/GUI. Задача управляющего плана состоит в обеспечении среды исполнения для перечисленных действий с целью настройки конфигурации и управления второй компонентой элемента сети (NE⁴) - машиной пересылки FE. Результат настройки конфигурации определяет способ трактовки пакетов, проходящих через FE.

1.1.2. Компоненты FEC⁵

Машина пересылки FE представляет собой объект NE, который первым получает сетевые пакеты (из сети в NE).

Связанная с сервисом компонента FE просматривает пакет с целью обеспечения для него обработки, определенной компонентами CPC для данного типа сервиса IP. Различные службы будут использовать разные компоненты FEC. Сервисные модули могут объединяться в цепочки для поддержки более сложных типов сервиса (в рамках описанной ниже модели Linux FE).

Будучи созданной для поддержки конкретной службы, сервисная компонента FE будет по-прежнему соответствовать принципам модели пересылки.

1.1.2.1. Модель машины пересылки IP в Linux

На рисунке показана модель Linux FE для отдельного устройства. Единственной обязательной частью этой модели является модуль пересылки (Пересылка), соответствующий RFC 1812. Различные модули сетевого экранирования (FW), управления входящим и исходящим трафиком (TC⁶) не являются обязательными и могут даже использоваться для обхода модуля RFC 1812. Эти модули показаны в виде простых блоков на пути передачи данных и фактически могут представлять собой каскады из множества субмодулей. Дополнительную информацию о таких модулях вы найдете в [10] и [11].

¹Control Plane - плоскость (уровень) управления

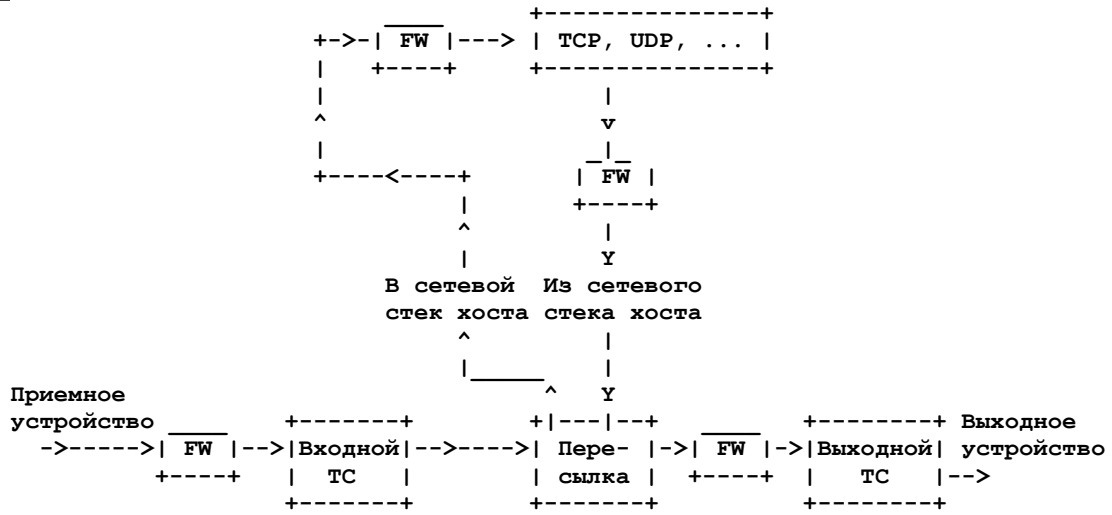
²Control Plane Components - компоненты плоскости управления.

³Forwarding Engine - машина пересылки.

⁴Network Element.

⁵Forwarding Engine Components - компоненты машины пересылки.

⁶Traffic Control - контроль трафика.



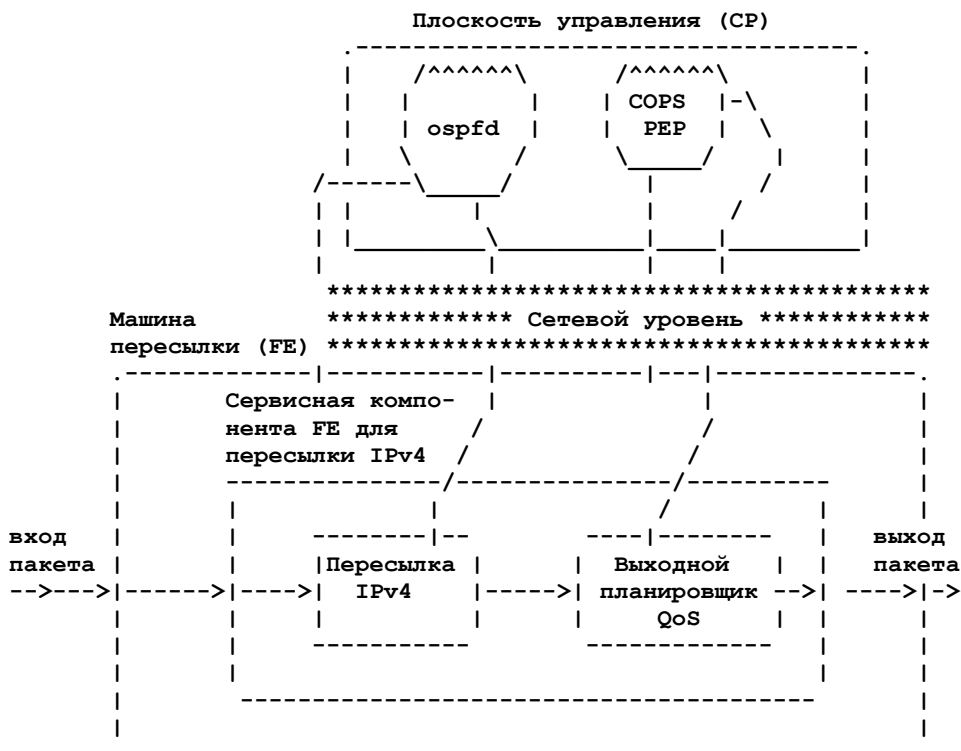
Пакеты, прибывающее на входное устройство, сначала проходят через модуль межсетевого экранирования (FW), который может отбрасывать (drop) и изменять (mangle) пакеты или выполнять с ними иные операции. После прохождения модуля FW входящие пакеты в зависимости от принятой политики, могут попадать во входной модуль контроля трафика TC, который выполняет операции по измерению и регулированию потоков входящего трафика. Пакеты могут отбрасываться входным модулем TC в зависимости от результатов измерения уровня трафика и принятой политики. После этого модуля пакет передается единственному обязательному модулю, который обеспечивает пересылку в соответствии с требованиями RFC 1812. Пакет может быть отброшен, если он не соответствует требованиям RFC 1812, RFC 1122, а также дополняющих их документов. Этот модуль является точкой выбора пути, из которой пакет, направленный принявшему его сетевому элементу NE, может быть передан сетевому стеку хоста.

Пакеты, не адресованные данному NE, могут проходить через submodule маршрутизации на базе правил (часть модуля пересылки), если такая маршрутизация поддерживается. Затем пакет передается следующему модулю сетевого экранирования, который может отбросить или изменить пакет в зависимости от настроек submodule и выбранной политики. После прохождения этого модуля пакет попадает в выходной фильтр контроля трафика (TC).

Выходной TC может отбрасывать пакеты с учетом политики, состояния очередей, уровня насыщения и правил управления скоростью исходящего потока. На этом этапе используются выходные очереди и задержки или отбрасывание пакета могут происходить как до его включения в очередь, так и после этого. Судьба пакета определяется выбранными для модуля алгоритмами и политикой.

1.1.3. Службы IP

Служба IP - это процессы обработки пакета IP внутри NE. Эти процессы определяются комбинацией CPC и FEC.



Занимаемое службой время начинается с момента прихода пакета в NE и заканчивается в момент, когда пакет покидает NE. Существенно, что поведение служб IP в этом контексте определяет конкретным хостом. Компоненты CP, запущенные на NE, определяют сквозной для всего пути контроль служб с помощью управляющих приложений и сигнальных протоколов. Такие распределенные компоненты CPC унифицируют сквозное представление служб IP. Как было отмечено выше, такие компоненты CP определяют поведение FE (и NE) по отношению к описываемому пакету.

Простым примером службы IP может служить классическая пересылка¹ IPv4. В этом случае управляющие компоненты (протоколы маршрутизации OSPF, RIP и т. п.) и фирменные средства настройки конфигурации CLI/GUI изменяют таблицы пересылки FE для того, чтобы обеспечить простой сервис по пересылке пакетов на следующий интервал (next hop). Обычно NE, обеспечивающие такой сервис, называют маршрутизаторами.

На приведенном рисунке показан простой пример реализации FE<->CPC для обеспечения классической пересылки IPv4 с некоторыми дополнительными функциями QoS для управления выходными очередями.

Демон ospfd управляет работой протокола OSPF, а COPS PEP² представляет собой дополнительную компоненту CPC. Компонента IPv4 FE включает модуль пересылки IPv4 и модуль выходного планировщика QoS. В качестве дополнительной службы может быть добавлен сервис пересылки на основе правил между модулем пересылки IPv4 и модулем планировщика QoS. Простейший классический вариант будет включать только модуль пересылки IPv4.

Опыт использования сетей говорит о важности добавления в маршрутизаторы новых типов сервиса, удовлетворяющих современным требованиям. Для решения этих задач были созданы и стандартизованы новые службы, которые могут выходить за пределы содержимого заголовков сетевого уровня. Однако, для обеспечивающих пересылку пакетов устройств NE по-прежнему используется термин «маршрутизатор». Новые службы (которые могут выходить за классические пределы заголовков L3) включают межсетевое экранирование, QoS с использованием Diffserv и RSVP, NAT, маршрутизацию на базе правил и т. п. Для таких служб создаются новые протоколы и средства управления.

Одним из крайних определений сервиса IP является «все, за что сервис-провайдеры могут взять деньги».

2. Архитектура Netlink

Управление компонентами IP-сервиса определяется с использованием шаблонов. Компоненты FEC и CPC участвуют в предоставлении услуг IP путем обмена данными с использованием таких шаблонов. FEC может непрерывно получать обновления от компоненты CPC, указывающие как предоставлять услуги (например, для пересылки пакетов IPv4, добавления, удаления или изменения маршрутов).

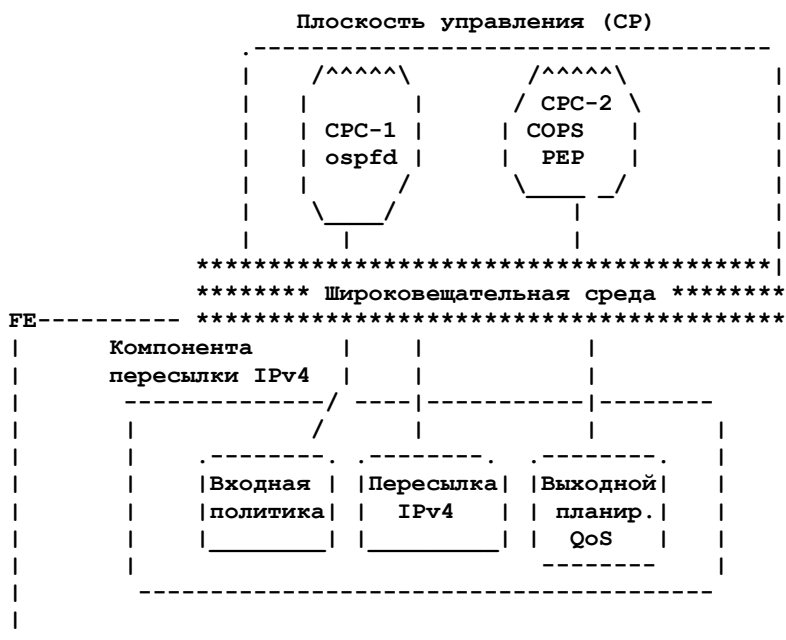
Взаимодействие между FEC и CPC в контексте Netlink определяется протоколом. Netlink предоставляет механизмы для CPC (в пользовательском пространстве) и FEC (в ядре), позволяющие им получить свои собственные определения для протокола. Это связано с тем, что пользовательское пространство и ядро находятся на разных уровнях безопасности. Следовательно, для обмена информацией между компонентами требуется протокол. Такой протокол обычно обеспечивается неким привилегированным сервисом, который имеет возможность копирования данных между различными уровнями безопасности. Будем называть такую службу сервисом Netlink. Этот сервис может также инкапсулироваться в протоколы транспортного уровня, если CPC и FEC выполняются на разных узлах. Компоненты FEC и CPC, используя механизмы Netlink, могут выбрать надежный протокол для обмена данными. По умолчанию Netlink не обеспечивает гарантированного обмена данными.

Отметим, что FEC и CPC могут располагаться на одном уровне защиты памяти и использовать системный вызов connect() для создания прямого пути и обмена информацией через этот путь. В данном документе этот механизм рассматриваться не будет - отметим лишь возможность его реализации. В данном документе предполагается, что FEC является частью ядра, а CPC размещается в пользовательском пространстве. Это не означает однако, что приведенная в документе информация относится лишь к случаю размещения этих компонент в разных областях защиты и не привязывает компоненты к одному узлу.

Отметим, что Netlink позволяет обеим компонентам участвовать в предоставлении сервиса IP.

2.1. Логическая модель Netlink

На рисунке показана простая диаграмма логических связей между компонентами FEC и CPC. В качестве примера использована FEC пересылки IPv4 (служба NETLINK_ROUTE, описанная ниже).



Netlink логически моделирует FEC и CPC в форме узлов, связанных между собой через широковещательную среду.

¹Forwarding

²Policy Enforcement Point - точка реализации политики.

Свойства среды обусловлены сервисом. В приведенном примере показана ширококестельная среда, принадлежащая к расширенному сервису пересылки IPv4.

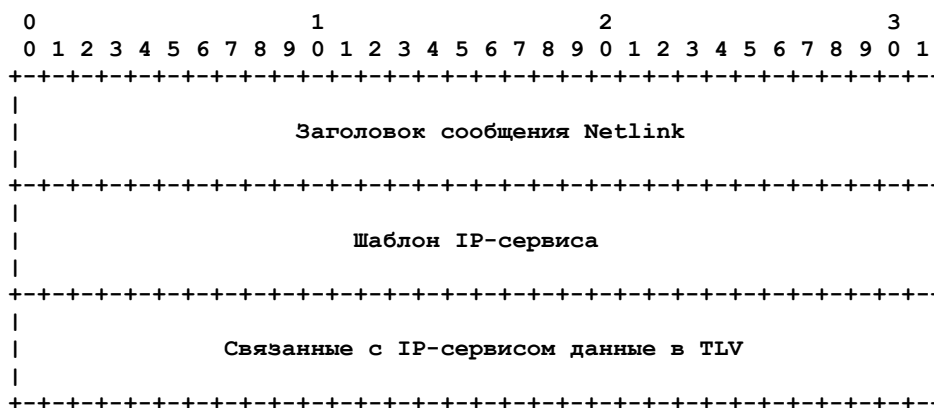
Узлы (CPC и FEC в рассматриваемом примере) подключены к среде передачи и регистрируются для получения сообщений определенных типов. CPC может подключаться к множеству сред, если это способствует более эффективному управлению сервисом. Все узлы (CPC и FEC) принимают пакеты из ширококестельной среды. Пакеты могут отбрасываться средой передачи, если они имеют непригодный формат или содержат ошибки. Отброшенные пакеты не поступают ни на один из узлов. Сервис Netlink может передавать отправителю сигналы об ошибках при обнаружении непригодных пакетов Netlink.

Передаваемые в среду пакеты могут быть ширококестельными, групповыми или индивидуальными. Узлы FEC и CPC регистрируют свою заинтересованность в сообщениях определенного типа для их обработки или мониторинга.

В приложениях 1 и 2 приведено более детальное рассмотрение этого взаимодействия.

2.2. Формат сообщений

В сообщениях Netlink существует три уровня - заголовок сообщения Netlink, шаблон IP-сервиса и связанные с IP-сервисом данные.



Сообщения Netlink используются для обмена данными между FEC и CPC, параметризации FEC, асинхронной передачи сведений о событиях FEC компонентам CPC и сбора/просмотра статистики (обычно с помощью CPC).

Заголовок сообщения Netlink используется для всех типов сервиса, тогда как шаблоны (IP Service Template) связаны с конкретными типами. Каждая служба IP передает данные параметризации (от CPC к FEC) или отклики (от FEC к CPC). Эти данные передаются в формате TLV¹ и являются уникальными для сервиса.

Отдельные компоненты сообщений Netlink подробно рассматриваются ниже.

2.3. Модель протокола

Здесь описано, как Netlink обеспечивает механизм ориентированного на службы взаимодействия между FEC и CPC.

2.3.1. Адресация служб

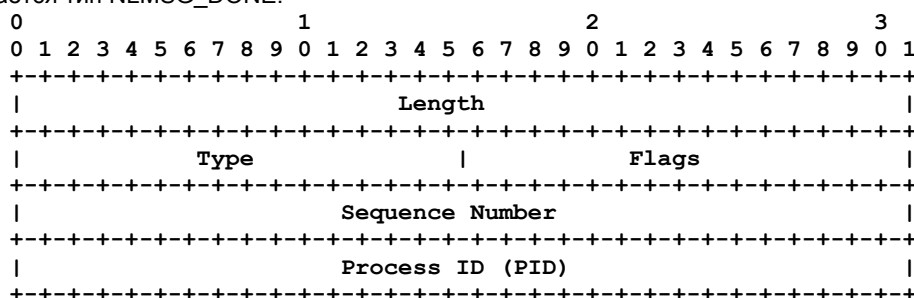
Для получения доступа сначала нужно соединиться с сервисом на FE. Соединение организуется путем системного вызова socket() для домена PF_NETLINK. Каждая компонента FEC указывается номером протокола. В результате вызова могут создаваться сокеты типа SOCK_RAW или SOCK_DGRAM, хотя Netlink не различает сокеты этих типов. Соединение с сокетом обеспечивает основу для адресации FE->CPC.

Затем организуется подключение к сервису (в любой момент в течение срока существования соединения) путем ввода обусловленной сервисом команды (от CPC к FEC, в основном для настройки конфигурации), команды сбора статистики или подписки (отказа) на уведомления о связанных с сервисом событиях. Закрытие сокета прерывает транзакцию.

Примеры рассматриваются в приложениях 1 и 2.

2.3.2. Заголовок сообщений Netlink

Сообщения Netlink представляют собой поток байтов с одним или несколькими заголовками Netlink и связанными с ними данными (payload). Если данных слишком много для одного сообщения, они могут быть разделены на несколько сообщений Netlink, которые обычно называют многокомпонентным сообщением². Для таких сообщений первый и последующие заголовки, за исключением последнего, содержат флаг NLM_F_MULTI. В заголовке последнего сообщения указывается тип NLMMSG_DONE.



¹Type-Length-Value - тип-размер-значение.

²Multipart message.

Формат заголовка сообщения Netlink показан на рисунке.

Заголовок включает описанные ниже поля.

Length - 32 бита

Размер сообщения в байтах с учетом заголовка.

Type - 16 битов

Это поле определяет тип содержимого в сообщении.

Фактически поле может включать один из стандартных идентификаторов типа, приведенных в таблице.

| | |
|--------------------|---|
| NLMSG_NOOP | Сообщение игнорируется. |
| NLMSG_ERROR | Сообщение сигнализирует об ошибке и поле данных содержит структуру nlmsgerr . Такие сообщения обычно передаются от FEC к CPC и могут рассматриваться как NACK ¹ . |
| NLMSG_DONE | Сообщение является последней частью многокомпонентного сообщения. |

Отдельные службы IP могут использовать добавочные типы сообщений, например сервис NETLINK_ROUTE задает несколько таких типов, включая RTM_NEWLINK, RTM_DELLINK, RTM_GETLINK, RTM_NEWADDR, RTM_DELADDR, RTM_NEWROUTE, RTM_DELROUTE и др.

Flags - 16 битов

Стандартные флаги, используемые в заголовках Netlink, приведены в таблице

| | |
|----------------------|--|
| NLM_F_REQUEST | Этот флаг должен устанавливаться для всех откликов (обычно они передаются из пользовательского пространства в ядро). |
| NLM_F_MULTI | Сообщение является частью (не последней) многокомпонентного сообщения. Для последней части указывается тип NLMSG_DONE. |
| NLM_F_ACK | Запрос на подтверждение при успехе. Обычно этот флаг устанавливается для сообщений из пользовательского пространства (CPC) в ядро (FEC). |
| NLM_F_ECHO | Возвратить «эхо» для данного запроса. Обычно этот флаг устанавливается для сообщений из пользовательского пространства (CPC) в ядро (FEC). |

В запросах GET для конфигурационной информации, передаваемых в FEC используются дополнительные флаги.

| | |
|---------------------|---|
| NLM_F_ROOT | Возвращать полную таблицу вместо одной записи. |
| NLM_F_MATCH | Возвращать все записи, соответствующие критерию, переданному в поле данных сообщения. |
| NLM_F_ATOMIC | Возвращать атомарную картину (atomic snapshot) таблицы, которая указана. Установка этого флага может требовать специальных привилегий, поскольку флаг способен прерывать сервис FE на достаточно продолжительное время. |

Подходящим макросом для поля флагов является

```
NLM_F_DUMP = NLM_F_ROOT OR NLM_F_MATCH
```

В запросах NEW также могут использоваться дополнительные флаги.

| | |
|----------------------|---|
| NLM_F_REPLACE | Заменить существующий объект конфигурации в соответствии с данным запросом. |
| NLM_F_EXCL | Не заменять существующий объект новым. |
| NLM_F_CREATE | Создать объект конфигурации, если его не существует. |
| NLM_F_APPEND | Добавить объект в конце списка имеющихся. |

Для тех, кто хорошо знаком с операциями на сокетах маршрутизации BSD, в таблице приведены эквиваленты таких операций.

| BSD | Netlink |
|---------------|---|
| ADD | NLM_F_CREATE или NLM_F_EXCL |
| CHANGE | NLM_F_REPLACE |
| Check | NLM_F_EXCL |
| APPEND | NLM_F_CREATE |

Sequence Number - 32 бита

Порядковый номер сообщения.

Process ID (PID) - 32 бита

Идентификатор процесса (PID), передающего сообщение. Значение PID используется ядром для мультиплексирования в нужный сокет. При передаче сообщений из ядра в пользовательское пространство устанавливается PID = 0.

2.3.2.1. Механизмы создания протоколов

Один из способов организации надежного протокола обмена между FEC и CPC является использование комбинации порядковых номеров, ACK² и таймеров повтора передачи. Порядковые номера и подтверждения ACK обеспечиваются Netlink, таймеры обеспечиваются ОС Linux.

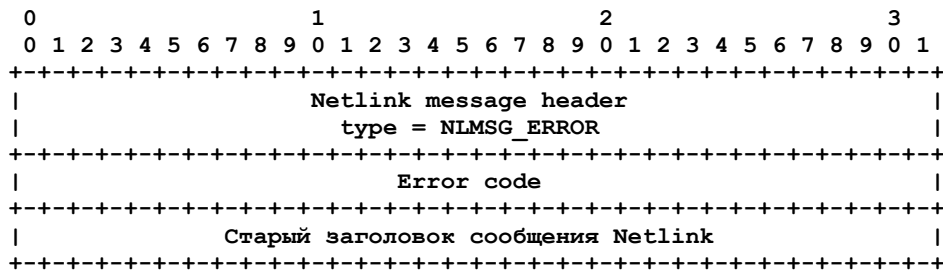
Можно также создать протокол heartbeat³ для обмена между FEC и CPC за счет использования флагов ECHO и сообщений типа NLMSG_NOOP.

¹Подтверждение отрицательного результата - Negative ACK. Прим. перев.

²Подтверждений.

³Пульс.

2.3.2.2. Сообщение ACK в Netlink



Эти сообщения используются как для передачи подтверждений (ACK), так и для передачи информации об отрицательном результате (NACK). Обычно такие сообщения передаются от FEC к CPC (в ответ на сообщение с запросом подтверждения). Однако CPC должны обеспечивать возможность передачи сообщений ACK в адрес FEC при наличии соответствующего запроса. Семантика этих сообщений зависит от сервиса IP.

Error code - integer (обычно 32 бита)

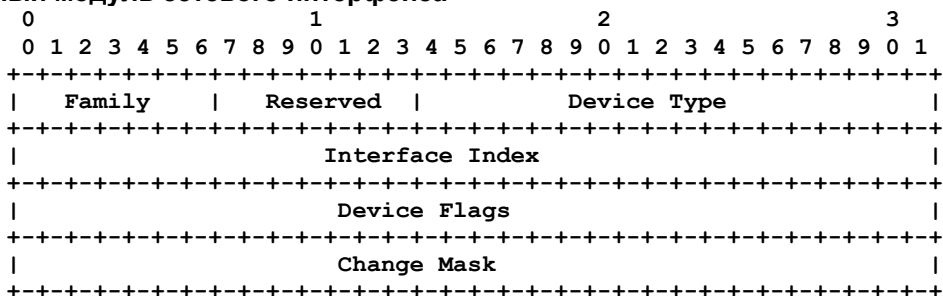
Нулевое значение кода ошибки говорит о том, что сообщение является подтверждением успеха (ACK). Такие сообщения содержат заголовок исходного сообщения Netlink, который может использоваться для сравнения (например, порядкового номера).

Отличный от нуля код говорит об отрицательном результате (NACK). В таких ситуациях данные Netlink, которые были переданы ядру, возвращаются вместе с исходным заголовком Netlink. Устанавливается также пригодное для вывода с помощью `reggor()` значение кода ошибки (не в заголовке сообщения, а в переменной окружения).

2.3.3. Шаблоны FE системных служб

Существуют системные службы, которые предлагают свой сервис для использования другими службами. Обычно они включают возможность настройки конфигурации, сбора статистики, прослушивание сведений об изменении общих ресурсов, управление адресами IP, каналные события и т. п. Данный раздел включает описание подобных служб для их логического разделения (несмотря на то, что все они доступны через FEC `NETLINK_ROUTE`). Причина этого заключается в том, что они существуют в `NETLINK_ROUTE` в силу исторически сложившихся причин (ошибки), связанных с тем, что сокеты BSD 4.4 Route реализованы как часть сокетов пересылки IPv4.

2.3.3.1. Сервисный модуль сетевого интерфейса



Эта служба обеспечивает возможность создания и удаления сетевых интерфейсов, а также получения информации о существующем интерфейсе. Интерфейс может быть физическим или виртуальным и не связан с сетевым протоколом (например, с помощью такого сообщения можно определить интерфейс `x.25`). Шаблон сообщения показан на рисунке.

Family - 8 битов

Это поле всегда имеет значение `AF_UNSPEC`.

Device Type - 16 битов

Определяет тип канала (Ethernet, туннель и т. п.). В данном документе рассматривается только IPv4, хотя тип канала не зависит от протокола L3.

Interface Index - 32 бита

Уникальный идентификатор интерфейса.

Device Flags - 32 бита

Флаги интерфейса, перечисленные в таблице.

| Флаг | Значение | Флаг | Значение |
|------------------------|--|-----------------------|---|
| IFF_UP | Интерфейс активизирован администратором | IFF_NOTRAILERS | Следует избегать использования трейлеров. |
| IFF_BROADCAST | Установлен корректный широковещательный адрес. | IFF_ALLMULTI | Принимать все пакеты с групповыми адресами. |
| IFF_DEBUG | Флаг режима отладки для интерфейса. | IFF_MASTER | Ведущий интерфейс для транка с распределением нагрузки. |
| IFF_LOOPBACK | Петлевой интерфейс (loopback). | IFF_SLAVE | Ведомый интерфейс для транка с распределением нагрузки. |
| IFF_POINTOPOINT | Интерфейс типа "точка-точка". | IFF_MULTICAST | Поддержка групповой адресации. |
| IFF_RUNNING | Интерфейс находится в работающем состоянии. | IFF_PORTSEL | Интерфейс может выбирать тип среды с помощью <code>ifmap</code> . |
| IFF_NOARP | Для интерфейса не требуется протокол ARP. | IFF_AUTOMEDIA | Активизирован автоматический выбор типа среды. |
| IFF_PROMISC | Интерфейс работает в режиме захвата (promiscuous). | IFF_DYNAMIC | Интерфейс создан в динамическом режиме. |

Change Mask - 32 бита

Зарезервированное поле, которое должно иметь значение 0xFFFFFFFF.

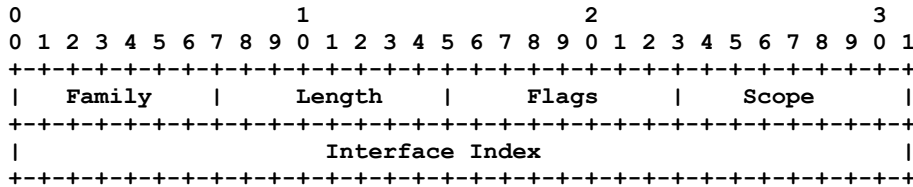
Применимые к данному сервису атрибуты перечислены в таблице.

| Атрибут | Описание | Атрибут | Описание |
|-----------------------|---|-------------------|--|
| IFLA_UNSPEC | Не определен. | IFLA_MTU | Значение MTU для устройства |
| IFLA_ADDRESS | Аппаратный адрес интерфейса на уровне L2. | IFLA_LINK | Значение ifindex для канала, к которому подключено устройство. |
| IFLA_BROADCAST | Аппаратный широковещательный адрес интерфейса на уровне L2. | IFLA_QDISC | Строка ASCII, указывающая имя дисциплины управления выходными очередями ¹ . |
| IFLA_IFNAME | Имя устройства (строка ASCII). | IFLA_STATS | Статистика для интерфейса. |

К данному типу сервиса относятся сообщения Netlink **RTM_NEWLINK**, **RTM_DELLINK** и **RTM_GETLINK**.

2.3.3.2. Модуль службы адресов IP

Эта служба обеспечивает возможность добавления и удаления адресов, а также получения сведений об IP-адресах, связанных с данным интерфейсом. Шаблон сообщения службы предоставления адресов² показан на рисунке.

**Family - 8 битов**

Идентификатор семейства адресов: **AF_INET** для IPv4 и **AF_INET6** для IPv6.

Length - 8 битов

Размер маски адреса.

Flags - 8 битов

| Флаг | Описание |
|-------------------------|---|
| IFA_F_SECONDARY | Вторичный адрес (псевдоним интерфейса) |
| IFA_F_PERMANENT | Постоянный адрес, установленный пользователем. Отсутствие этого флага говорит о динамическом выделении адреса (например, с помощью системы автоматической настройки конфигурации) |
| IFA_F_DEPRECATED | Недействующий (deprecated) адрес IP. |
| IFA_F_TENTATIVE | Предполагаемый (tentative) адрес IP. Процедура обнаружения дубликатов адресов находится в стадии разработки. |

Scope - 8 битов

Область действия адреса.

| | |
|-----------------------|---|
| SCOPE_UNIVERSE | Адрес глобального действия. |
| SCOPE_SITE | Адрес действует в пределах данного сайта (только для IPv6). |
| SCOPE_LINK | Адрес имеет смысл только для данного устройства. |
| SCOPE_HOST | Адрес имеет смысл только для данного хоста. |

Атрибуты сервиса перечислены в таблице.

| Атрибут | Описание | Атрибут | Описание |
|--------------------|-------------------------------------|----------------------|--|
| IFA_UNSPEC | Не определен. | IFA_BROADCAST | Широковещательный адрес для протокола RAW. |
| IFA_ADDRESS | Адрес интерфейса для протокола RAW. | IFA_ANYCAST | Anycast-адрес для протокола RAW. |
| IFA_LOCAL | Локальный адрес для протокола RAW. | IFA_CACHEINFO | Кэшированная информация об адресе. |
| IFA_LABEL | Имя интерфейса (строка ASCII). | | |

К данному типу сервиса относятся сообщения Netlink **RTM_NEWADDR**, **RTM_DELADDR** и **RTM_GETADDR**.

3. Определенные в данный момент IP-службы Netlink

Хотя, как было отмечено выше, существует множество других служб IP, использующих Netlink, в данном документе рассматривается лишь небольшая часть этих служб, интегрированных в ядро версии 2.4.6. К таким службам относятся **NETLINK_ROUTE**, **NETLINK_FIREWALL** и **NETLINK_ARPD**³.

3.1. Служба NETLINK_ROUTE

Эта служба позволяет СРС изменять таблицу маршрутизации IPv4 в машине пересылки FE. Кроме того, данный сервис может применяться СРС для получения данных об обновлении маршрутов и сбора статистики.

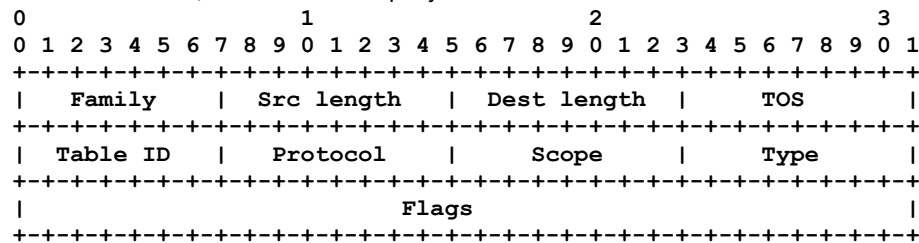
¹Egress root queuing discipline.

²Address provisioning service.

³На момент перевода документа (ядро 2.6.10) был определен целый ряд дополнительных служб, информацию о которых вы найдете в файле <linux/netlink.h>. *Прим. перев.*

3.1.1. Модуль службы маршрутизации

Эта служба обеспечивает возможность создания и удаления маршрутов, а также получения информации о сетевых маршрутах. Формат шаблона сообщения показан на рисунке.



Family - 8 битов

Идентификатор семейства адресов: **AF_INET** для IPv4 и **AF_INET6** для IPv6.

Src length - 8 битов

Размер префикса IP-адреса отправителя.

Dest length - 8 битов

Размер префикса IP-адреса получателя.

TOS - 8 битов

Восьмибитовое поле TOS (следует отказаться от него для освобождения места под DSCP).

Table ID - 8 битов

Идентификатор таблицы. Поддерживается до 255 таблиц маршрутизации.

| | | | |
|-------------------------|------------------------------------|-----------------------|--------------------|
| RT_TABLE_UNSPEC | Неуказанная таблица. | RT_TABLE_MAIN | Основная таблица. |
| RT_TABLE_DEFAULT | Используемая по умолчанию таблица. | RT_TABLE_LOCAL | Локальная таблица. |

Пользователь может выделять дополнительные значения в диапазоне¹ от **RT_TABLE_UNSPEC** (0) до **RT_TABLE_DEFAULT** (253).

Protocol - 8 битов

Указывает кто добавил маршрут в таблицу.

| Протокол | Источник маршрута | Протокол | Источник маршрута |
|------------------------|-----------------------------|----------------------|-----------------------|
| RTPROT_UNSPEC | Неизвестен. | RTPROT_BOOT | При загрузке системы. |
| RTPROT_REDIRECT | Из сообщения ICMP redirect. | RTPROT_STATIC | Администратор. |
| RTPROT_KERNEL | Ядро. | | |

Значения, превышающие **RTPROT_STATIC** (4)², не интерпретируются ядром и включены только с информационными целями. Эти значения могут использоваться, чтобы пометить источник маршрутной информации или различать разные демоны маршрутизации. Идентификаторы уже присвоенные демонам маршрутизации вы можете найти в файле <linux/rtnetlink.h>.

Scope - 8 битов

Область видимости маршрута (корректная дистанция до получателя).

| | |
|--------------------------|--|
| RT_SCOPE_UNIVERSE | Глобальный маршрут. |
| RT_SCOPE_SITE | Внутренний маршрут локальной автономной системы. |
| RT_SCOPE_LINK | Маршрут на данном канале (соединении). |
| RT_SCOPE_HOST | Маршрут на локальном хосте. |
| RT_SCOPE_NOWHERE | Получателя не существует. |

Значения в диапазоне от **RT_SCOPE_UNIVERSE** (0) до **RT_SCOPE_SITE** (200), не включая граничные, могут использоваться для пользовательских идентификаторов.

Type - 8 битов

Тип маршрута.

| Тип | Получатель |
|------------------------|---|
| RTN_UNSPEC | Неизвестный маршрут. |
| RTN_UNICAST | Шлюз или прямой маршрут. |
| RTN_LOCAL | Маршрут к локальному интерфейсу. |
| RTN_BROADCAST | Локальный широковещательный маршрут (передается как broadcast). |
| RTN_ANYCAST | Локальный anycast-маршрут (передается как unicast) |
| RTN_MULTICAST | Локальный групповой (multicast) маршрут. |
| RTN_BLACKHOLE | Маршрут для отбрасывания пакетов без уведомления (черная дыра). |
| RTN_UNREACHABLE | Недостижимый получатель. Пакеты отбрасываются с передачей отправителю сообщения ICMP о недоступности адресата. |
| RTN_PROHIBIT | Запрещенный маршрут. Пакеты отбрасываются с передачей отправителю сообщения ICMP о запрете доступа к адресату. |
| RTN_THROW | При использовании маршрутизации на базе правил указывает на продолжение просмотра маршрутов в другой таблице. При обычной маршрутизации пакеты отбрасываются с передачей отправителю сообщения ICMP о недоступности адресата. |

¹ Не включая граничные значения 0 и 253. Прим. перев.

² В файле <linux/rtnetlink.h> указано, что значение **RTPROT_STATIC** (4) также не интерпретируется ядром. Прим. перев.

| <i>Тип</i> | <i>Получатель</i> |
|---------------------|--|
| RTN_NAT | Правило трансляции сетевых адресов. |
| RTN_XRESOLVE | Указывает на внешний преобразователь (resolver). В настоящее время еще не реализовано. |

Flags - 32 бита

Дополнительная информация о маршруте.

| | |
|-----------------------|---|
| RTM_F_NOTIFY | При изменении маршрута пользователю передается уведомление. |
| RTM_F_CLONED | Маршрут клонирован из другого маршрута. |
| RTM_F_EQUALIZE | Маршрут допускает случайный выбор следующего интервала (next hop) в случае наличия нескольких путей (в настоящее время не реализовано). |

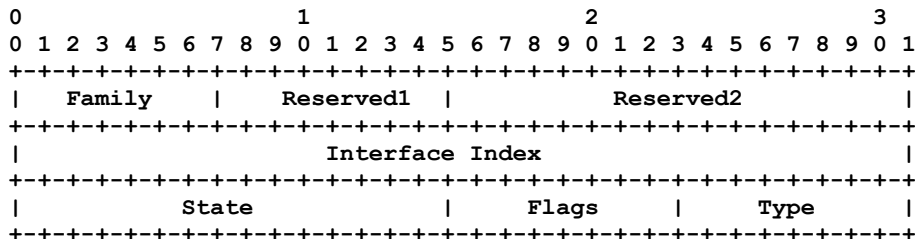
Имеющие отношение к данному сервису атрибуты перечислены в таблице.

| <i>Атрибут</i> | <i>Описание</i> |
|----------------------|--|
| RTA_UNSPEC | Игнорируется. |
| RTA_DST | Протокольный адрес источника маршрута. |
| RTA_SRC | Протокольный адрес конечной точки маршрута. |
| RTA_IIF | Индекс входного интерфейса. |
| RTA_OIF | Индекс выходного интерфейса. |
| RTA_GATEWAY | Протокольный адрес шлюза для маршрута. |
| RTA_PRIORITY | Приоритет маршрута. |
| RTA_PREFSRC | Предпочтительный адрес отправителя при наличии нескольких адресов. |
| RTA_METRICS | Присвоенная маршруту метрика (например, RTT, начальный размер окна TCP и т. п.). |
| RTA_MULTIPATH | Атрибуты следующего интервала для маршрута с множеством путей (Multipath route). |
| RTA_PROTOINFO | Атрибут маршрутизации, основанный на политике межсетевого экрана. |
| RTA_FLOW | Область маршрута (Route realm). |
| RTA_CACHEINFO | Кэшированная информация о маршруте. |

Для этого типа сервиса поддерживаются дополнительные сообщения Netlink **RTM_NEWROUTE**, **RTM_DELROUTE** и **RTM_GETROUTE**.

3.1.2. Модуль учета соседей

Этот сервис обеспечивает возможность добавления и удаления записей о соседях (например, ARP, IPv4 neighbor solicitation и т. п.), а также получения информации о существующих записях таблицы соседей. Шаблон сообщений этой службы показан на рисунке.

**Family - 8 битовое**

Идентификатор семейства адресов: **AF_INET** для IPv4 и **AF_INET6** для IPv6.

Interface Index - 32 бита

Уникальный индекс интерфейса.

State - 16 битовое

Битовая маска, которая может включать перечисленные в таблице биты состояния.

| | |
|-----------------------|--|
| NUD_INCOMPLETE | Продолжаются попытки преобразования адреса. |
| NUD_REACHABLE | Подтверждено наличием рабочей записи в кэше. |
| NUD_STALE | Просроченная запись из кэша. |
| NUD_DELAY | Сосед больше не достижим. Трафик передан, ожидается подтверждение. |
| NUD_PROBE | В настоящее время осуществляется запрос на обновление записи в кэше. |
| NUD_FAILED | Некорректная запись в кэше. |
| NUD_NOARP | Устройство, которое не выполняет обнаружения соседей (ARP). |
| NUD_PERMANENT | Статическая запись. |

Flags - 8 битовое

| | | | |
|------------------|------------------|-------------------|--------------------|
| NTF_PROXY | Запись проху ARP | NTF_ROUTER | Маршрутизатор IPv6 |
|------------------|------------------|-------------------|--------------------|

Применимые к этому сервису атрибуты перечислены в таблице.

Parent Qdisc - 32 бита

Используется для иерархической структуризации дисциплин очередей. Если это значение совпадает с идентификатором и TC_H_ROOT, данный экземпляр qdisc является называется корневым³ (старшим).

TCM Info - 32 бита

Для этого поля FE обычно устанавливает значение 1 за исключением тех случаев, когда экземпляр Qdisc уже используется (в этом случае в поле помещается значение счетчика использования данного экземпляра). При передаче со стороны CPC в направлении FEC это поле обычно имеет значение 0 за исключением тех случаев, когда оно используется в контексте фильтрации. В таких случаях 32-битовое поле делится на 16-битовые поля приоритета и протокола. Протоколы определены в исходных кодах ядра (файл <include/linux/if_ether.h>). Наиболее широко используемым протоколом является ETH_P_IP (протокол IP).

Значение приоритета используется для разрешения конфликтов при пересечении фильтрующих выражений. Базовые атрибуты этого типа сервиса перечислены в таблице.

| Атрибут | Описание |
|-------------|--|
| TCA_KIND | Каноническое имя компоненты FE. |
| TCA_STATS | Базовая статистика использования FEC. |
| TCA_RATE | Оценка скорости для FEC (расчет на основе текущего состояния). |
| TCA_XSTATS | Специфическая статистика FEC. |
| TCA_OPTIONS | Вложенные атрибуты, связанные с FEC. |

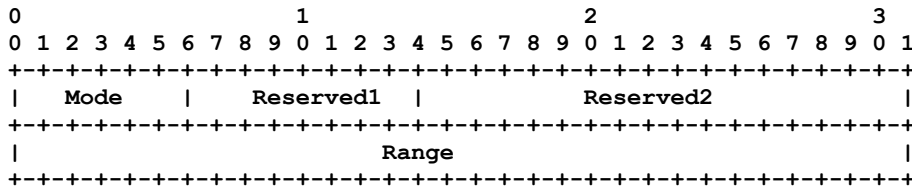
В приложении 3 дается пример конфигурации компоненты FE для дисциплины FIFO.

Для этого типа сервиса поддерживаются дополнительные сообщения Netlink **RTM_NEWQDISC**, **RTM_DELQDISC**, **RTM_GETQDISC**, **RTM_NEWTCCLASS**, **RTM_DELTCLASS**, **RTM_GETTCLASS**, **RTM_NEWTFILTER**, **RTM_DELTFILTER** и **RTM_GETTFILTER**.

3.2. Служба NETLINK_FIREWALL

Эта служба позволяет CPC принимать пакеты через сервисные модули межсетевого экрана IPv4 в FE, манипулировать этими пакетами и повторно передавать их. Правило межсетевого экрана является первым для активизации перенаправления пакетов. CPC информирует FEC о своем желании получить метаданные для пакета или реальные данные из него, а также сообщает максимальный размер данных, которые будут перенаправляться. Перенаправленные пакеты по-прежнему сохраняются в FEC, ожидая решения о своей судьбе от CPC. Решение может быть простой командой на восприятие или отбрасывание пакета (в этом случае решение применяется к пакету, все еще находящемуся в FEC) или включить измененный пакет, который должен быть передан взамен исходного.

Существует два типа сообщений, передаваемых от CPC к FEC - Mode (режим) и Verdict (решение). Сообщения типа Mode незамедлительно передаются FEC и сообщают о том, что CPC желает принимать от FEC. Сообщения типа Verdict передаются FEC после принятия решения о дальнейшей судьбе полученного пакета. Формат сообщений рассматривается ниже.



Опишем сначала сообщение, указывающее режим.

Mode - 8 битов

Определяет тип информации в пакетах, отправляемых CPC:

IPQ_COPY_META - копировать в CPC только метаданные для пакета.

IPQ_COPY_PACKET - копировать в CPC метаданные и содержимое поля данных пакета.

Range - 32 бита

В режиме **IPQ_COPY_PACKET** это значение определяет максимальный размер копируемых данных.

Пакет и связанные с ним метаданные, полученные из пользовательского пространства, показаны на рисунке.

³Root qdisc.

3.3. Служба NETLINK_ARPD

Этот сервис используется CPC для поддержки таблицы соседей в FE. Формат сообщений, передаваемых между FEC и CPC, описан параграфе, посвященном службе учета соседей (стр. 10).

Предполагается, что сервис CPC принимает участие в работе протоколов организации соседских отношений (neighbor solicitation protocol).

Сообщение типа RTM_NEWNEIGH передается CPC от FE для информирования CPC об изменениях, которые могут произойти с записью для этого соседа.

Сообщения RTM_GETNEIGH используются для получения информации о конкретном соседе.

4. Литература

4.1. Нормативные документы

- [1] Braden, R., Clark, D. and S. Shenker, "Integrated Services in the Internet Architecture: an Overview", RFC 1633, June 1994.
- [2] Baker, F., "Requirements for IP Version 4 Routers", [RFC 1812](#), June 1995.
- [3] Blake, S., Black, D., Carlson, M., Davies, E, Wang, Z. and W. Weiss, "An Architecture for Differentiated Services", [RFC 2475](#), December 1998.
- [4] Durham, D., Boyle, J., Cohen, R., Herzog, S., Rajan, R. and A. Sastry, "The COPS (Common Open Policy Service) Protocol", RFC 2748, January 2000.
- [5] Moy, J., "OSPF Version 2", STD 54, [RFC 2328](#), April 1998.
- [6] Case, J., Fedor, M., Schoffstall, M. and C. Davin, "Simple Network Management Protocol (SNMP)", STD 15, [RFC 1157](#), May 1990.
- [7] Andersson, L., Doolan, P., Feldman, N., Fredette, A. and B. Thomas, "LDP Specification", [RFC 3036](#), January 2001.
- [8] Bernet, Y., Blake, S., Grossman, D. and A. Smith, "An Informal Management Model for DiffServ Routers", RFC 3290, May 2002.

4.2. Дополнительная литература

- [9] G. R. Wright, W. Richard Stevens. "TCP/IP Illustrated Volume 2, Chapter 20", June 1995.
- [10] <http://www.netfilter.org>
- [11] <http://diffserv.sourceforge.net>

5. Вопросы безопасности

Netlink работает в доверенной среде на одном хосте с разделением ядра и пользовательского пространства. Средствами Linux обеспечивается возможность открывать сокет только для процессов с флагом возможностей CAP_NET_ADMIN (обычно процессы, запущенные пользователем root).

6. Благодарности

- 1) **Andi Kleen** за страницы руководства (man pages) для netlink и rtnetlink.
- 2) **Alexey Kuznetsov** за добавление модели службы доставки IP в Netlink. Исходный вариант символьного устройства Netlink создал Alan Cox.
- 3) **Jeremy Ethridge** за исполнение роли «не понимающего Netlink» и обзор документа с точки зрения его восприятия.

Приложение 1. Пример иерархии служб

На рисунке показан пример единичного IP-сервиса **foo** и взаимодействие компонент CP и FE для этой службы (метки 1-3).

Эта схема используется так же, как пример адресации CP<->FE. В этом приложении иллюстрируется только семантика адресации. В Приложении 2 эта схема рассматривается с точки зрения протокольного взаимодействия между компонентами CPC и FEC сервиса (метки 4-10).

Приложение 3. Примеры

```

0                               1                               2                               3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               Length (52)                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Type (RTM_NEWQDISC)           | Flags (NLM_F_EXCL |                     |
|                               |NLM_F_CREATE | NLM_F_REQUEST) |         |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               Sequence Number (произвольное значение)   |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               Process ID (0)                             |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Family (AF_INET) | Reserved1 | Reserved1 |                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               Interface Index (4)                       |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               Qdisc handle (0x1000001)                   |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               Parent Qdisc (0x1000000)                   |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               TCM Info (0)                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               Type (TCA_KIND) | Length(4)                |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               Value ("pfifo")                             |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               Type (TCA_OPTIONS) | Length(4)            |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               Value (limit=100)                         |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

В этом примере рассматривается простое конфигурационное сообщение Netlink, передаваемое от TC CPS выходной очереди TC FIFO. Этот алгоритм управления очередью основан на учете пакетов и отбрасывании их при достижении порогового значения 100. Предполагается, что очередь находится в иерархии с родителем Parent = 100:0, Classid = 100:1 и размещается на устройстве с ifindex = 4.

Адреса авторов**Jamal Hadi Salim**

Znyx Networks

Ottawa, Ontario

Canada

EMail: hadi@znyx.com**Hormuzd M Khosravi**

Intel

2111 N.E. 25th Avenue JF3-206

Hillsboro OR 97124-5961

USA

Phone: +1 503 264 0334

EMail: hormuzd.m.khosravi@intel.com**Andi Kleen**

SuSE

Stahlgruberring 28

81829 Muenchen

Germany

EMail: ak@suse.de**Alexey Kuznetsov**

INR/Swsoft

Moscow

Russia

E-Mail: kuznet@ms2.inr.ac.ru

Перевод на русский язык

Николай Малых

nmalykh@gmail.com

Полное заявление авторских прав

Copyright (C) The Internet Society (2003). All Rights Reserved.

This document and translations of it may be copied and furnished to others, and derivative works that comment on or otherwise explain it or assist in its implementation may be prepared, copied, published and distributed, in whole or in part, without restriction of any kind, provided that the above copyright notice and this paragraph are included on all such copies and derivative works. However, this document itself may not be modified in any way, such as by removing the copyright notice or references to the Internet Society or other Internet organizations, except as needed for the purpose of developing Internet standards in which case the procedures for copyrights defined in the Internet Standards process must be followed, or as required to translate it into languages other than English.

The limited permissions granted above are perpetual and will not be revoked by the Internet Society or its successors or assignees.

This document and the information contained herein is provided on an "AS IS" basis and THE INTERNET SOCIETY AND THE INTERNET ENGINEERING TASK FORCE DISCLAIMS ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

Подтверждение

Финансирование функций RFC Editor в настоящее время обеспечивается Internet Society.